

A whole-genome scan for association with invasion success in the fruit fly *Drosophila suzukii* using contrasts of allele frequencies corrected for population structure

Laure Olazcuaga¹, Anne Loiseau¹, Hugues Parrinello², Mathilde Paris³, Antoine Fraimout¹, Christelle Guedot⁴, Lauren M. Diepenbrock⁵, Marc Kenis⁶, Jinping Zhang⁷, Xiao Chen⁸, Nicolas Borowieck⁹, Benoit Facon¹⁰, Heidrun Vogt¹¹, Donald K. Price¹², Heiko Vogel¹³, Benjamin Prud'homme³, Arnaud Estoup^{*,1,14} and Mathieu Gautier^{1,14,*}

¹INRA, UMR CBGP (INRA IRD Cirad Montpellier SupAgro), Montferrier-sur-Lez, France

²MGX, Biocampus Montpellier, CNRS, INSERM, Université de Montpellier, Montpellier, France

³Aix Marseille Université, CNRS, IBDM, Marseille, France

⁴Department of Entomology, University of Wisconsin, Madison, WI

⁵Department of Entomology and Plant Pathology, NC State University

⁶CABI, Delémont, Switzerland

⁷MoA-CABI Joint Laboratory for Bio-safety, Chinese Academy of Agricultural Sciences, BeiXiaGuan, Haidian Qu, China

⁸College of Plant Protection, Yunnan Agricultural University, Kunming 650201, Yunnan Province, China

⁹UMR INRA-CNRS-Université Côte d'Azur Sophia Agrobiotech Institute, Sophia Antipolis, France

¹⁰UMR Peuplements Végétaux et Bioagresseurs en Milieu Tropical, INRA, Saint-Pierre, La Réunion, France

¹¹Julius Kühn-Institut (JKI), Federal Research Centre for Cultivated Plants, Institute for Plant Protection in Fruit Crops and Viticulture, Dossenheim, Germany

¹²School of Life Sciences, University of Nevada, Las Vegas, Las Vegas, NV

¹³Department of Entomology, Max Planck Institute for Chemical Ecology, Jena, Germany

¹⁴These authors are joint senior authors on this work

***Corresponding author:** E-mail: mathieu.gautier@inra.fr and arnaud.estoup@inra.fr

Abstract

Evidence is accumulating that evolutionary changes are not only common during biological invasions but may also contribute directly to invasion success. The genomic basis of such changes is still largely unexplored. Yet, understanding the genomic response to invasion may help to predict the conditions under which invasiveness can be enhanced or suppressed. Here we characterized the genome response of the spotted wing drosophila *Drosophila suzukii* during the worldwide invasion of this pest insect species, by conducting a genome-wide association study to identify genes involved in adaptive processes during invasion. Genomic data from 22 population samples were analyzed to detect genetic variants associated with the status (invasive versus native) of the sampled populations based on a newly developed statistic, we called C_2 , that contrasts allele frequencies corrected for population structure. This new statistical framework has been implemented in an upgraded version of the program BAYPASS. We identified a relatively small set of single nucleotide polymorphisms (SNPs) that show a highly significant association with the invasive status of populations. In particular, two genes *RhoGEF64C* and *cpo*, the latter contributing to natural variation in several life-history traits (including diapause) in *Drosophila melanogaster*, contained SNPs significantly associated with the invasive status in the two separate main invasion routes of *D. suzukii*. Our methodological approaches can be applied to any other invasive species, and more generally to any evolutionary model for species characterized by non-equilibrium demographic conditions for which binary covariables of interest can be defined at the population level.

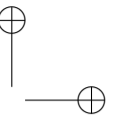
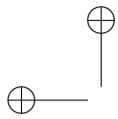
Key words: Biological invasion, *Drosophila suzukii*, GWAS, BAYPASS, Pool-Seq.

1 **Introduction**

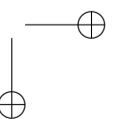
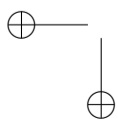
2 Managing and controlling introduced species
3 require an understanding of the ecological
4 and evolutionary processes that underlie
5 invasions. Biological invasions are also of
6 more general interest because they constitute
7 natural experiments that allow investigation
8 of evolutionary processes on contemporary
9 timescales. Colonizers are known to experience
10 differences in biotic interactions, climate,
11 availability of resources, and disturbance regimes
12 relative to populations in their native regions,
13 often with opportunities for colonizers to evolve
14 changes in resource allocation which favor their
15 success (Balanya *et al.*, 2006; Dlugosch *et al.*,
16 2015; Lee and Gelembiuk, 2008). Adaptive
17 evolutionary shifts in response to novel selection
18 regimes may therefore be central to initial
19 establishment and spread of invasive species
20 after introduction (Colautti and Barrett, 2013;
21 Colautti and Lau, 2015). In agreement with
22 this adaptive evolutionary shift hypothesis,
23 experimental evidence is accumulating that
24 evolutionary changes are not only common
25 during invasions but also may contribute directly
26 to invasion success (Bock *et al.*, 2015; Colautti
27 and Lau, 2015; Ellstrand and Schierenbeck,
28 2000; Facon *et al.*, 2011; Lee, 2002; Ochocki and
29 Miller, 2017; Williams *et al.*, 2016). However,
30 despite an increase in theoretical and empirical

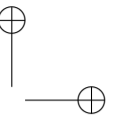
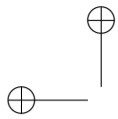
31 studies on the evolutionary biology of invasive
32 species in the past decade, the genetic basis of
33 evolutionary adaptations during invasions is still
34 largely unexplored (Barrett, 2015; Reznick *et al.*,
35 2019; Welles and Dlugosch, 2018).

36 The spotted wing drosophila, *Drosophila*
37 *suzukii*, represents an attractive biological model
38 to study invasive processes. This pest species,
39 native to South East Asia, initially invaded
40 North America and Europe, simultaneously in
41 2008, and subsequently La Réunion Island (Indian
42 Ocean) and South America, in 2013. Unlike most
43 Drosophilids, this species lays eggs in unripe
44 fruits by means of its sclerotized ovipositor. In
45 agricultural areas, it causes dramatic losses in
46 fruit production, with a yearly cost exceeding
47 one billion euros worldwide (e.g., Asplen *et al.*,
48 2015; Cini *et al.*, 2012). The rapid spreading
49 of *D. suzukii* in America and Europe suggests
50 its remarkable ability to adapt or to acclimate
51 to new environments and host plants. Using
52 evolutionarily neutral molecular markers, Adrien
53 *et al.* (2014) and Fraimout *et al.* (2017) finely
54 deciphered the routes taken by *D. suzukii* in
55 its invasion worldwide. Interestingly, both studies
56 showed that North American (plus Brazil) and
57 European (plus La Réunion Island) populations
58 globally represent separate invasion routes, with
59 different native source populations and multiple
60 introduction events in both invaded regions
61 (Fraimout *et al.*, 2017). These two major
62 and separate invasion pathways provide the



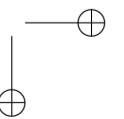
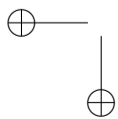
63 opportunity to evaluate replicate evolutionary 95 a chromosomal inversion polymorphism between
64 trajectories. Finally, *D. sukikii* is a good model 96 the native and introduced range.
65 species for finely interpreting genomic signals 97 Identifying loci underlying invasion success
66 of interest due to the availability of genome 98 can be considered in the context of whole-
67 assemblies for this species (Chiu *et al.*, 2013; 99 genome scan for association with population-
68 Ometto *et al.*, 2013; Paris *et al.*, 2019) along with 100 specific covariate. These approaches, also known
69 the large amount of genomic and gene annotation 101 as Environmental Association Analysis (EAA),
70 resources available in its closely related model 102 have received considerable attention in recent
71 species *D. melanogaster* (Hoskins *et al.*, 2015). 103 years (e.g., Coop *et al.*, 2010; de Villemereuil
72 In this context, advances in high-throughput 104 and Gaggiotti, 2015; Frichot *et al.*, 2013; Gautier,
73 sequencing technologies together with population 105 2015). Most of the methodological developments
74 genomics statistical methods offer novel 106 have focused on properly accounting for the
75 opportunities to disentangle responses to selection 107 covariance structure among population allele
76 from other forms of evolution. These advances 108 frequencies that is due to the shared demographic
77 are thus expected to provide insights into the 109 history of the populations. This neutral covariance
78 genomic changes that might have contributed to 110 structure may indeed confound the relationship
79 the success in a new environment (reviewed in 111 between the across population variation in allele
80 Bock *et al.*, 2015; Welles and Dlugosch, 2018). 112 frequencies and the covariates of interest (Coop
81 Hence, comparing the structuring of genetic 113 *et al.*, 2010; Frichot *et al.*, 2013, 2015; Gautier,
82 diversity on a whole genome scale among invasive 114 2015). Yet, defining relevant environmental
83 populations and their source populations might 115 characteristics or traits as proxy for invasion
84 allow the characterization of the types of genetic 116 success remains challenging and might even
85 variation involved in adaptation during invasion 117 be viewed as the key aim. Therefore, we
86 of new areas and their potential ecological 118 propose to simply summarize invasion success
87 functions. For example, Puzey and Vallejo-Marin 119 into a binary variable corresponding to the
88 (2014) used whole genome resequencing data to 120 population's historical status (i.e., invasive or
89 scan for shifts in site frequency spectra to detect 121 native) based on previous studies. By extension,
90 positive selection in introduced populations 122 functional annotation of the associated variants
91 of monkey-flower (*Mimulus guttatus*). Regions 123 identified may provide insights into candidate
92 putatively under selection were associated with 124 traits underlying invasion success (Estoup *et al.*,
93 flowering time and abiotic and biotic stress 125 2016; Li *et al.*, 2008; Wu *et al.*, 2019).
94 tolerance and included regions associated with

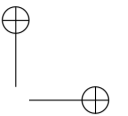
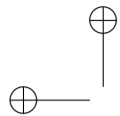




126 The Bayesian hierarchical model initially 158 sampled in both the invasive (n=16 populations)
127 proposed by Coop *et al.* (2010), later extended 159 and native (n=6 populations) ranges of the
128 in Gautier (2015) and implemented in the 160 species. We then estimated the C_2 statistics
129 software BAYPASS, represents one of the most 161 associated with the invasive vs. native status
130 flexible and powerful frameworks to carry 162 of the populations on a worldwide scale or
131 out EAA since it efficiently accounts for the 163 considering separately each of the two invasion
132 correlation structure among allele frequencies in 164 routes (European and American) as characterized
133 the sampled populations. Although association 165 by Fraimout *et al.* (2017). Our aim was to identify
134 analyses may be carried out with categorical or 166 genomic regions and genes involved in adaptive
135 binary covariables (see the example of *Littorina* 167 processes underlying the invasion success of *D.*
136 population ecotypes in Gautier, 2015), the 168 *suzukii*.

137 assumed linear relationship with allele frequencies 169 **New Approaches**
138 is not entirely satisfactory and may even be 170 To identify single nucleotide polymorphisms
139 problematic when dealing with small data sets or 171 (SNPs) associated with a population-specific
140 if one wishes to disregard some populations. 172 binary trait, such as the invasive versus native
141 In the present study, we developed a non- 173 status of *D. suzukii* populations, we developed
142 parametric counterpart for the association model 174 a new statistic, we called C_2 . The C_2 statistic
143 implemented in BAYPASS (Gautier, 2015). This 175 was designed to contrast SNP allele frequencies
144 new approach relies on a contrast statistic, 176 between the two groups of populations specified by
145 we named C_2 , that compares the standardized 177 the binary trait while accounting for the possibly
146 population allele frequencies (i.e., the allele 178 complex evolutionary history of the different
147 frequencies corrected for the population structure) 179 populations. Indeed, the shared population
148 between the two groups of populations specified 180 history is a major (neutral) contributor to allele
149 by the binary covariable of interest. We evaluated 181 frequency differentiation across populations (e.g.
150 the performance of this statistic on simulated data 182 Bonhomme *et al.*, 2010; Gunther and Coop, 2013)
151 and used it to characterize the genome response of 183 that may confound association signals (e.g. Coop
152 *D. suzukii* during its worldwide invasion. To that 184 *et al.*, 2010; Gautier, 2015).
153 end, we generated Pool-Seq data (e.g., Gautier 185 We here relied on the multivariate normal
154 *et al.*, 2013; Schlotterer *et al.*, 2014) consisting 186 approximation introduced by Coop *et al.* (2010)
155 of whole-genome sequences of pools of individual 187 and further extended by Gautier (2015) to model
156 DNA (from n=50 to n=100 individuals per 188 population allele frequencies and to define the
157 pool) representative of 22 worldwide populations 189 C_2 contrast statistic. More precisely, consider a





190 sample made of J populations (each with a label 219 scaled allele frequencies that are corrected for
 191 $j=1, \dots, J$) that have been characterized for I bi- 220 both the population structure (summarized by $\mathbf{\Omega}$)
 192 allelic SNPs (each with a label $i=1, \dots, I$), with 221 and the across-population (e.g., ancestral) allele
 193 the reference allele arbitrarily defined (e.g., by 222 frequency (π_i).
 194 randomly drawing the ancestral or the derived 223 The C_2 contrast statistic is then simply defined
 195 state). Let α_{ij} represent the (unobserved) allele 224 as the mean squared difference of the sum of
 196 frequency of the reference allele at SNP i in 225 standardized allele frequencies of the two groups of
 197 population j . As previously defined and discussed 226 populations defined according to the binary trait
 198 (Coop *et al.*, 2010; Gautier, 2015), we introduced 227 modalities:
 199 an instrumental allele frequency α_{ij}^* (for each SNP 228
 200 i and population j) taking values on the real line 229 where $\mathbf{c} = c_{j(1..J)}$ is a vector of the trait values
 201 such that $\alpha_{ij} = \min(1, \max(0, \alpha_{ij}^*))$. 230 observed for each population j such that $c_j = 1$
 202 (respectively $c_j = -1$) if population j displays the
 203 (Coop *et al.* (2010) and Gautier 231 first (respectively second) trait modality. One may
 204 (2015), a multivariate Gaussian (prior) 232 also define $c_j = 0$ to exclude a given population j
 205 distribution of the vector $\boldsymbol{\alpha}_i^* = \{\alpha_{ij}^*\}_{1..J}$ is 233 from the comparison.

206 Following Coop *et al.* (2010) and Gautier 234
 207 (2015), a multivariate Gaussian (prior) 235
 208 distribution of the vector $\boldsymbol{\alpha}_i^* = \{\alpha_{ij}^*\}_{1..J}$ is 236
 209 then assumed for each SNP i : 237

$$\boldsymbol{\alpha}_i^* | \mathbf{\Lambda}, \pi_i \sim N_J(\pi_i \mathbf{1}_J; \pi_i(1 - \pi_i)\mathbf{\Omega}) \quad (1)$$

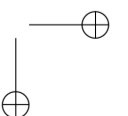
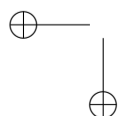
206 where $\mathbf{1}_J$ is the all-one vector of length J ; $\mathbf{\Omega}$ is the 235
 207 (scaled) covariance matrix of the population allele 236
 208 frequencies which captures information about 237
 209 their shared demographic history; and π_i is the 238
 210 weighted mean frequency of the SNP i reference 239
 211 allele. If $\mathbf{\Omega}$ is used to build a tree or an 240
 212 admixture graph (Pickrell and Pritchard, 2012), 241
 213 π_i corresponds to the root allele frequency. We 242
 214 further define for each SNP i the vector $\ddot{\boldsymbol{\alpha}}_i$ of 243
 215 standardized (instrumental) allele frequencies in 244
 216 the J populations as: 245

$$\ddot{\boldsymbol{\alpha}}_i = \Gamma_{\Omega}^{-1} \left\{ \frac{\alpha_{ij} - \pi_i}{\sqrt{\pi(1 - \pi_i)}} \right\}_{(1..J)} \quad (2)$$

217 where Γ_{Ω} results from the Cholesky decomposition 248
 218 of Ω (i.e., $\Omega = \Gamma_{\Omega}^t \Gamma_{\Omega}$). The vector $\ddot{\boldsymbol{\alpha}}_i$ thus contains 249

234 According to our model, the J elements of $\ddot{\boldsymbol{\alpha}}_i$
 235 are independent and identically distributed as
 236 a standard Gaussian distribution under the null
 237 hypothesis of only neutral marker differentiation.
 238 The C_2 statistic is thus expected to follow a χ^2
 239 distribution with one degree of freedom.

240 The estimation of the C_2 statistic was
 241 performed here under the hierarchical Bayesian
 242 model implemented using a Markov-Chain Monte
 243 Carlo (MCMC) algorithm in the BAYPASS
 244 software (Gautier, 2015). However, such a multi-
 245 level modeling approach shrinks the estimated
 246 posterior means of the C_2 toward their prior
 247 means, as already noticed in Gautier (2015)
 248 for the estimation of the SNP-specific XtX
 249 differentiation statistic defined as $XtX = \ddot{\boldsymbol{\alpha}}_i^t \ddot{\boldsymbol{\alpha}}_i$



(Gunther and Coop, 2013). To ensure proper calibration of both the C_2 and XtX estimates thus relied on the scaled posterior means of the $\hat{\alpha}_{ij}$'s, denoted $\hat{\underline{\alpha}}_i$ and computed as:

$$\hat{\underline{\alpha}}_i = \left\{ \frac{\hat{\alpha}_{ij} - \mu_{\hat{\alpha}}}{\sigma_{\hat{\alpha}}} \right\}_{(1 \dots J)} \quad (4)$$

where $\hat{\alpha}_{ij}$ is the posterior means of α_{ij} and $\mu_{\hat{\alpha}}$ (respectively $\sigma_{\hat{\alpha}}$) is the mean (respectively standard deviation) of the $I \times J$ $\hat{\alpha}_{ij}$'s ($\mu_{\hat{\alpha}} \simeq 0$ usually). The following estimators of XtX and C_2 , denoted for each SNP i as $\widehat{XtX^*}(i)$ and $\widehat{C}_2(i)$ respectively, were then obtained as:

$$\begin{aligned} \widehat{XtX^*}(i) &= \hat{\underline{\alpha}}_i^t \hat{\underline{\alpha}}_i \\ \widehat{C}_2(i) &= \frac{1}{\mathbf{c}^t \mathbf{c}} \left(\hat{\underline{\alpha}}_i^t \mathbf{c} \right)^2 \end{aligned} \quad (5)$$

Under the null hypothesis, $\widehat{XtX^*}(i) \sim \chi_J^2$ and $\widehat{C}_2(i) \sim \chi_1^2$ allowing one to rely on standard decision making procedures, e.g. based on p-values or more preferably on q-values to control for multiple-testing issues (Storey and Tibshirani, 2003).

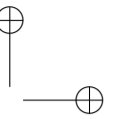
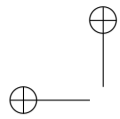
Results

Simulation-based evaluation of the performance of our novel statistical framework

To evaluate the performances of the C_2 contrast statistic for the identification of SNP associated with binary population-specific covariables, we simulated 100 data sets under the evolutionary scenario depicted in Figure 1A. Each simulated data set consisted of 5,000 SNPs genotyped for 320 individuals belonging to 16 differentiated populations subjected to two different contrasting

environmental constraints, denoted *ec1* and *ec2* in Figure 1A. The *ec1* constraint was aimed at mimicking adaptation of eight pairs of geographically differentiated populations to two different ecotypes (e.g., host plant) replicated in different geographic areas. Conversely, the *ec2* might be viewed as replicated local adaptive constraints with a first type *a* specifying a large native area with several geographically differentiated populations (here six), and a second type *b* specifying invasive areas with differentiated populations originating from various regions of the native area (i.e., not related to the same extent to their contemporary native populations). It should be noted that the two *ec1* types were evenly distributed in the population tree while for *ec2*, the type *b* was over-represented in 10 populations (Figure 1A). During the adaptive phase, the fitness of individuals in the environment of their population of origin was determined by their genotypes at 25 SNPs for *ec1* and 25 SNPs for *ec2* constraints (hereafter referred to as *ec1* and *ec2* selected SNPs, respectively). Overall, the realized F_{ST} (Weir and Cockerham, 1984) ranged from 0.110 to 0.122 (0.116 on average) across the data sets, a level of differentiation similar to that observed in our worldwide *D. sukikii* sample (see below).

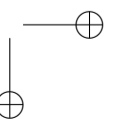
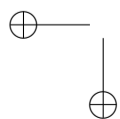
We further estimated with BAYPASS (Gautier, 2015) the C_2 statistics for each *ec1* or *ec2* contrasting environmental constraints together with the corresponding Bayes Factors (BF) as an

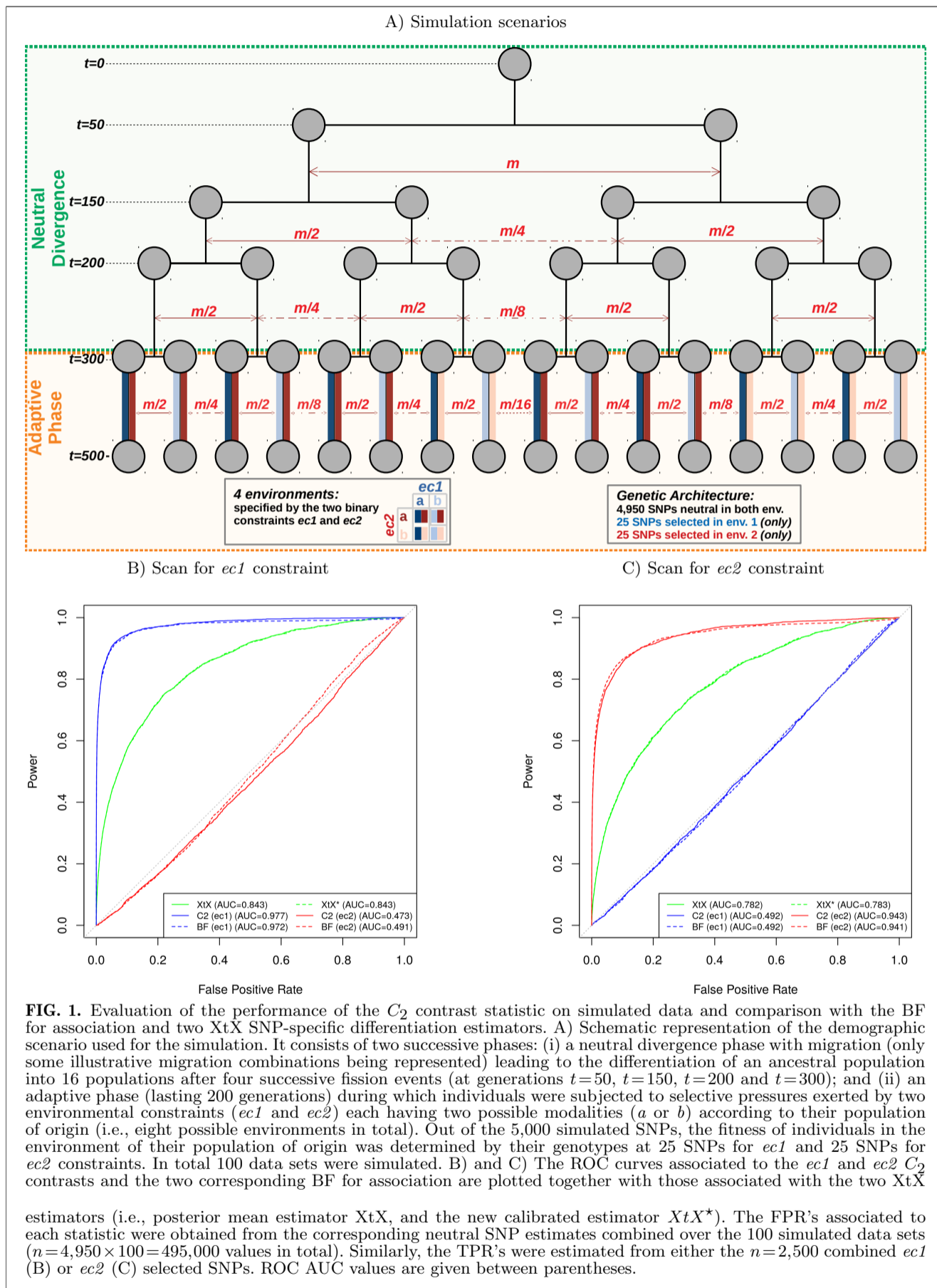


310 alternative measure of the support for association. 342 environmental constraint. In other words, no
311 For comparison purposes, we also estimated the 343 selection signal was identified by the C_2 statistic
312 SNP XtX differentiation statistic, using both the 344 computed for the *ec2* (respectively *ec1*) contrast
313 posterior mean estimator (Gautier, 2015) and the 345 on *ec1* (respectively *ec2*) selected SNPs, resulting
314 \widehat{XtX}^* estimator described above. Note however 346 in ROC AUC close to the value of 0.5 obtained
315 that, as an overall (covariate-free) differentiation 347 with a random classifier.
316 statistic, the XtX does not distinguish outlier 348
317 SNPs responding to the *ec1* constraint from those 349
318 responding to the *ec2* constraint. 350

319 Based on the status of each simulated SNPs 351 correlation between both statistics were fairly
320 (i.e., neutral, and *ec1* or *ec2* selected) and 352 high (Pearson's r equal to 0.983 and 0.923 for
321 combining results in the 100 simulated data sets, 353 *ec1* and *ec2*, respectively). Yet, one practical
322 standard receiver operating curves (ROCs) were 354 advantage of the C_2 statistic was its good
323 computed (Grau *et al.*, 2015) and plotted in 355 calibration with respect to the null hypothesis
324 Figure 1B (respectively 1C) for the six statistics. 356 of no association, the corresponding p-values
325 This allowed comparing for various thresholds 357 (assuming a χ^2 distribution with 1 degree of
326 covering their range of variation of the different 358 freedom) being close to uniform (Figure S1).

327 statistics, the power to detect *ec1* (respectively 359
328 *ec2*) selected SNPs (i.e., the proportion of true 360 highly correlated (Pearson's $r=0.998$) with
329 positives among the corresponding selected SNPs) 361 almost confounded ROC curves, but only the
330 as a function of the false positive rates (FPR, i.e., 362 \widehat{XtX}^* was properly calibrated (Figure S2). Their
331 the proportion of positives among neutral SNPs). 363 performances were however clearly worse than
332 The C_2 statistic was found efficient to detect 364 those obtained with the C_2 (and BF) statistics.
333 SNPs affected by *ec1* and *ec2* environmental 365 This was in part explained by their inability to
334 constraints, the area under the ROC curve 366 discriminate between the two types of selected
335 (AUC) being equal to 0.977 (Figure 1B) and 367 SNPs, selected SNPs overly differentiated in *ec2*
336 0.943 (Figure 1C), respectively. The unbalanced 368 generating false positives in the identification of
337 population representation of the two *ec2* types 369 *ec1* SNPs (Figure 1B) and vice versa. Accordingly,
338 had a limited impact on the performance of the C_2 370 ROC AUC in Figure 1B for the XtX were also
339 statistic to identify the underlying selected SNPs. 371 smaller than in Figure 1C, *ec1* selected SNPs
340 In addition, the C_2 statistics clearly discriminated 372 being more differentiated than those in *ec2* due
341 the selected SNPs according to their underlying 373 to the simulated design. Yet, the power of the

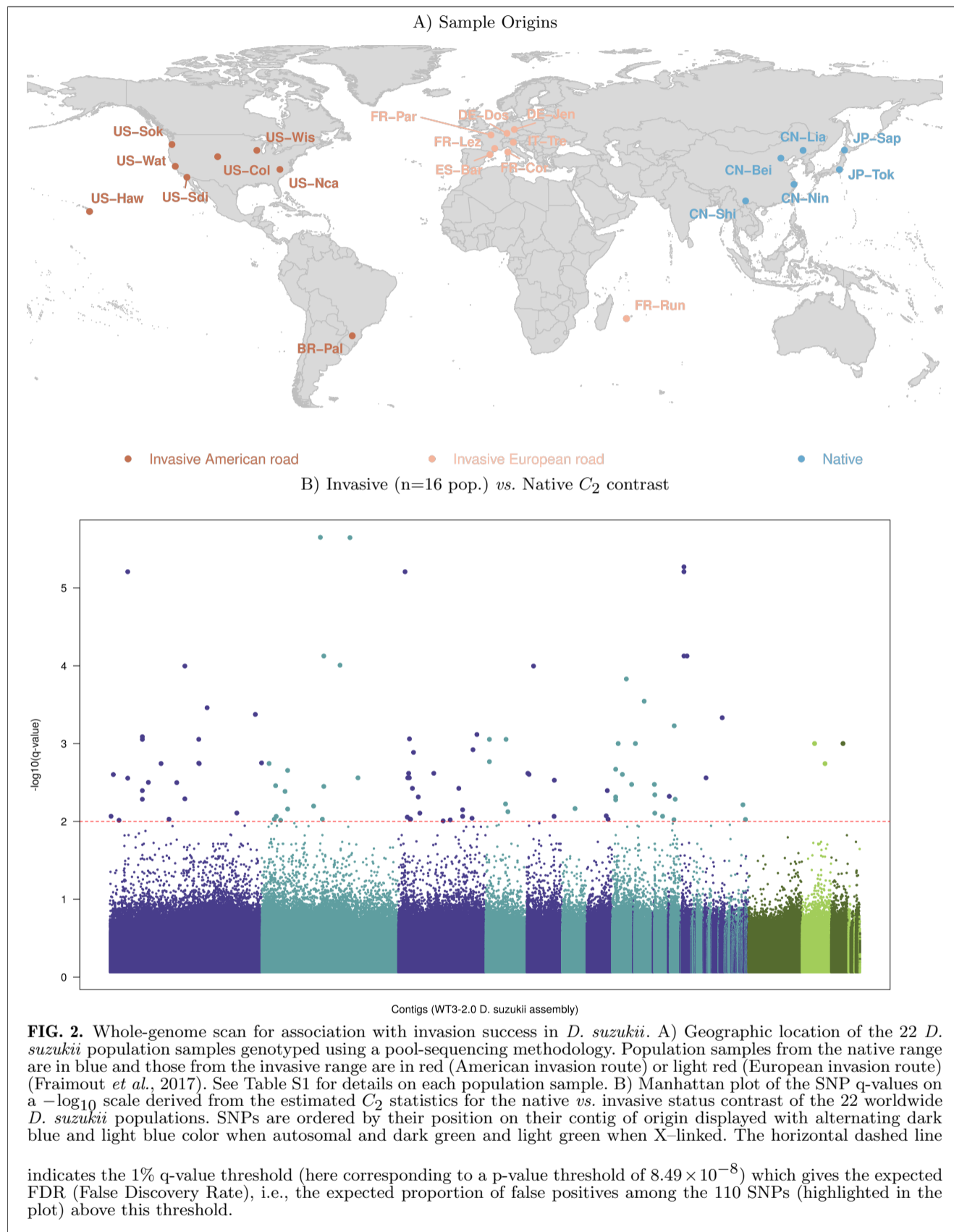


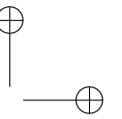
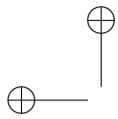


374 XtX statistic to detect *ec1* or *ec2* selected SNPs 406 sub-sampled into 154 autosomal and 26 X-linked
375 remained substantially smaller than that of the 407 data sets (of ca. 75,000 SNPs each) for further
376 corresponding C_2 contrast statistics. For instance, 408 analyses.
377 at the 1% p-value significance threshold, the power 409 The overall differentiation was estimated using
378 to detect *ec1* (respectively *ec2*) selected SNPs 410 the recently developed F_{ST} estimator for Pool-Seq
379 was equal to 72.6% (respectively 59.1%) with the 411 data (Hivert *et al.*, 2018). It ranged from 8.86% to
380 C_2 statistic and only 17.1% (respectively 10.4%) 412 9.02% (8.95% on average) for the autosomal data
381 with the $\widehat{XtX^*}$ estimator, even when considering 413 sets and from 17.6% to 17.8% (17.8% on average)
382 for the latter, a unilateral test to only target 414 for the X-chromosome data sets. Although a
383 overly differentiated SNPs. Note that, as expected 415 higher genetic differentiation is expected for the
384 from the good calibration of the $\widehat{XtX^*}$ statistic, 416 X-chromosome even under equal contribution of
385 similar results were obtained when considering 417 males and females to demography, the almost
386 empirical p-value thresholds computed from the 418 twice higher overall differentiation observed for
387 distribution of the XtX statistics estimated from 419 the X chromosome compared to autosomes might
388 neutral SNPs. 420 have been accentuated by unbalanced sex-ratio

389 Genome-wide scan for association with
390 invasion success in *D. suzukii*

391 To identify genomic regions associated with 423
392 the invasion success of *D. suzukii*, we carried 424 demography was beyond the scope of the present
393 out a genome scan, based on the C_2 statistic, 425 study, but for our purposes, this finding justified
394 to contrast the patterns of genetic diversity 426 to perform separate genome scans on autosomal
395 among 22 populations originating from either 427 and X-linked SNPs.
396 the native (n=6 populations) or invaded areas 428 We ran BAYPASS on the different data sets
397 (n=16 populations) (Figure 2A). To that end 429 to estimate, for every SNPs, the C_2 statistic
398 we sequenced pools of 50 to 100 individuals 430 that contrasts the allele frequencies of native
399 representative of each population (Table S1) 431 and invasive populations, while accounting for
400 and mapped the resulting sequencing reads 432 their shared population history as summarized
401 onto the newly released WT3-2.0 *D. suzukii* 433 in the scaled covariance matrix Ω . Interestingly,
402 genome assembly (Paris *et al.*, 2019). These 434 the estimated Ω matrices for autosomal and X-
403 Pool-Seq data allowed the characterization of 435 linked SNPs resulted in a similar structuring of
404 11,564,472 autosomal and 1,966,184 X-linked 436 the genetic diversity across the 22 populations
405 SNPs segregating in the 22 populations that were 437 (Figure S3), which may rule out selective

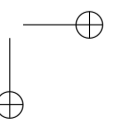
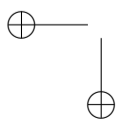




438 forces as the main driver of the differences of 470 XtX measure of overall differentiation. The (two-
439 global differentiation levels observed between the 471 sided) p-values derived from the latter were
440 two chromosome types. As expected from the 472 also well behaved (Figure S4B) and allowed the
441 simulation results, the distribution of the p-values 473 computation of q-values to control for multiple
442 derived from the C_2 statistics was well-behaved, 474 testing. As shown in Figure S5B, at the same 1%
443 being close to uniform for higher p-values (Figure 475 q-value threshold for XtX , 71 out of the 101 C_2
444 S4A). To account for multiple testing issues, we 476 significant SNPs were significantly differentiated
445 used the *qvalue* R package (Storey and Tibshirani, 477 but they represented only a small proportion of
446 2003) to compute the individual SNP q-values 478 the 35,546 significantly differentiated SNPs. This
447 plotted in Figure 2B. 479 is not surprising since invasion success is obviously

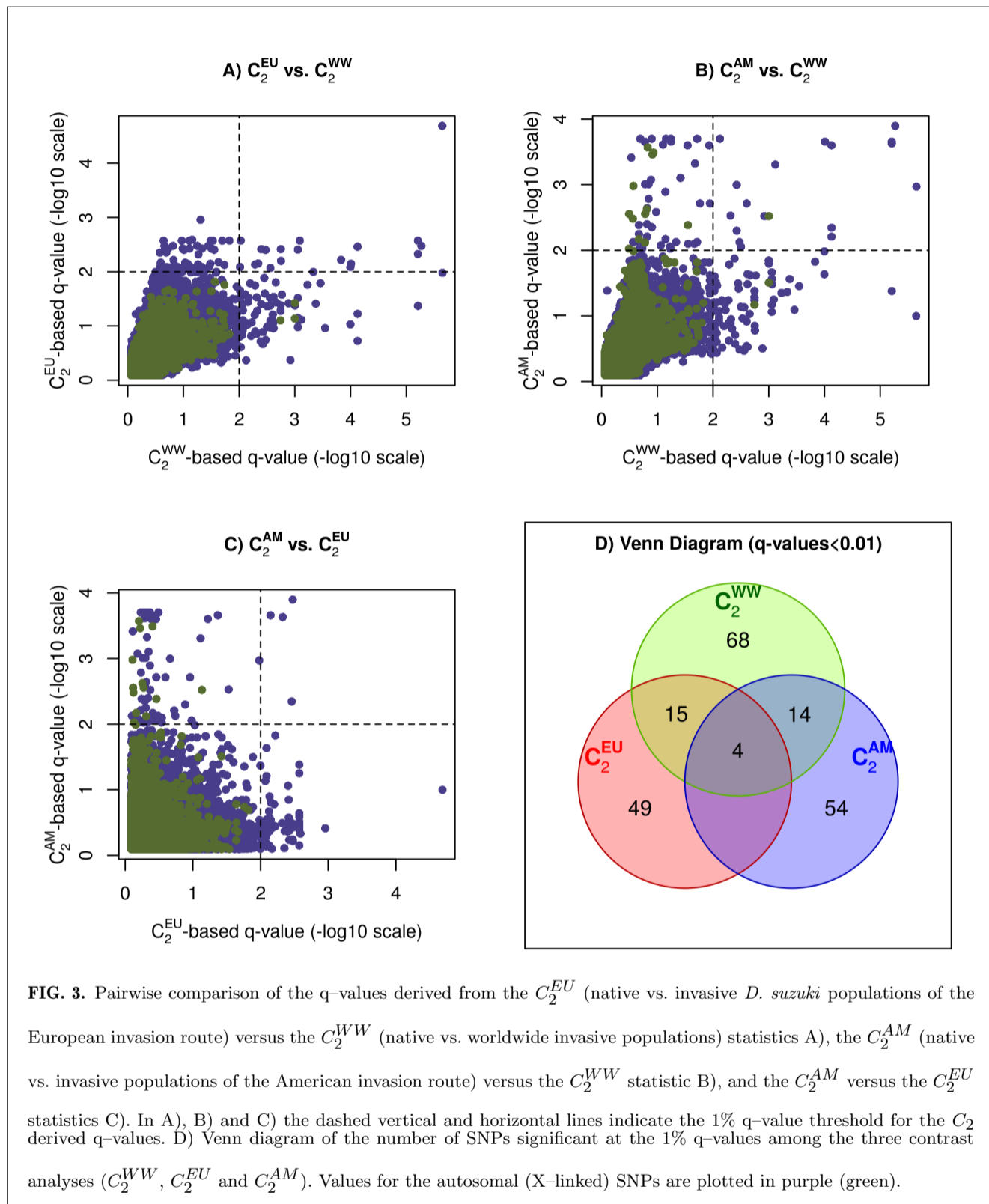
448 A striking feature of the resulting Manhattan 480 not the only selective constraint exerted on the 22
449 plot was the lack of clustering of SNPs with 481 worldwide populations considered here.

450 high q-values which might be related to a small 482 The North-American (plus Brazil) and
451 extent of linkage disequilibrium (LD) across the 483 European (plus La Réunion Island) populations
452 *D. sukukii* populations, as expected from their 484 globally represent separate invasion routes that
453 large effective population sizes (Fraitout *et al.*, 485 can be considered as two independent invasive
454 2017). We identified 101 SNPs (including three 486 replicates (Figure 2A). Interestingly enough,
455 X-linked) that were significant at the 1% q- 487 this feature of historical invasion fits well with
456 value threshold (i.e., 1% of these 101 SNPs are 488 the overall pattern of structuring of genetic
457 expected to be false positives). As a matter 489 diversity inferred from the Ω matrix estimated
458 of comparison, we also estimated the BF for 490 with our Pool-Seq data (see above and Figure
459 association of the (standardized) population allele 491 S3). To identify signals common or specific
460 frequencies with the native or invasive status of 492 to each invasion routes, we estimated the C_2
461 the population, i.e., under a parametric regression 493 statistic associated with the invasive vs. native
462 model (Gautier, 2015) (Figure S5A). Out of 494 status focusing either on the native and invasive
463 the 101 significant SNPs previously identified, 495 populations of the European invasion route
464 80 displayed a $BF > 20$ db, the threshold for 496 (C_2^{EU}), or native and invasive populations of
465 decisive evidence according to the Jeffreys' rule 497 the American invasion route (C_2^{AM}). Note that
466 (Jeffreys, 1961). However, in total, 6,406 SNPs 498 the two invasion routes were both represented
467 displayed a $BF > 20$ db probably as a consequence 499 by eight invasive populations, suggesting similar
468 of these BF's not accounting for multiple testing 500 power for the two C_2^{EU} and C_2^{AM} statistics. As
469 issue. We also compared the C_2 statistic to the 501 observed above, the distribution of p-values



502 derived from C_2^{EU} and C_2^{AM} were found well 534 signals among the two invasion routes (i.e., the
503 behaved (Figures S4C and S4D, respectively) 535 informative populations are distributed among the
504 and hence q-values to control for multiple 536 two routes). Most interestingly, four SNPs were
505 testing could be confidently computed. The 537 found significant at the 1% q-value threshold in
506 cross-comparisons of the C_2 statistics considering 538 the three contrast analyses (C_2^{EU} , C_2^{AM} and C_2^{WW})
507 the 22 worldwide populations (hereafter denoted 539 and might thus be viewed as strong candidates
508 C_2^{WW}), the C_2^{EU} and the C_2^{AM} are plotted in 540 for association with the global worldwide invasion
509 Figures 3A (C_2^{EU} versus C_2^{WW}), 3B (C_2^{AM} versus 541 success of *D. suzukii*.
510 C_2^{WW}) and 3C (C_2^{AM} versus C_2^{EU}). 542

511 In total, 204 SNPs (detailed in Table S2) 543 Annotation of candidate SNPs
512 were significant in at least one of the three 544 For annotation purposes, we relied on genomic
513 contrasts at the 1% q-value threshold. The overlap 545 resources available in *D. melanogaster*, a model
514 among the three different sets of significant 546 species closely related to *D. suzukii*. More
515 SNPs was summarized in the Venn diagram 547 specifically we extracted from the WT3-2.0 *D.*
516 displayed in Figure 3D. Among the 68 SNPs 548 *suzukii* genome assembly 5 kb long genomic
517 significant for the C_2^{EU} , 15 were also significant 549 sequences surrounding each of the 204 SNPs
518 for C_2^{WW} and 49 were not significant in the 550 identified above and aligned them onto the
519 other tests. Likewise, among the 72 SNPs found 551 *dmel6* reference genome (Hoskins *et al.*, 2015)
520 significant for the C_2^{AM} , 14 were also significant 552 using the BLAT algorithm implemented in the
521 for C_2^{WW} and 54 were not significant in the 553 program *pblat* (Wang and Kong, 2019). The
522 other tests. Hence, the majority of the significant 554 gene annotation available from the UCSC genome
523 SNPs identified with either the C_2^{EU} or the 555 browser allowed us to map 169 SNPs out of the 204
524 C_2^{AM} contrasts might be viewed as specific 556 SNPs onto 130 different *D. melanogaster* genes,
525 to one of the two invasion routes, the signal 557 145 SNPs lying within the gene sequences and 24
526 being lost in the global worldwide comparison 558 less than 2.5 kb apart (our predefined threshold;
527 for a substantial proportion of them. This is 559 Table S2). Only one of the four SNPs significant
528 presumably due to a reduced power resulting 560 for the three contrasts (C_2^{WW} , C_2^{EU} and C_2^{AM})
529 from the addition of non-informative populations 561 could not be assigned to a *D. melanogaster* gene,
530 when computing the C_2^{WW} statistic. Conversely, 562 because its derived 5 kb long sequences aligned
531 68 SNPs found significant with C_2^{WW} were neither 563 onto a *D. melanogaster* sequence located 10 kb
532 significant with C_2^{EU} nor C_2^{AM} contrasts. These 564 away from the closest annotated gene.
533 SNPs might correspond to partially convergent 565 Most of the 130 identified genes (80%) were
represented by a single SNP, a feature in

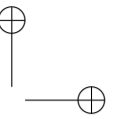
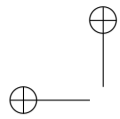


566 agreement with the visual lack of clustering of 569 that 14 of the 130 genes (ca. 11%) were long
 567 SNPs with strong signal already observed in the 570 non-coding RNA. We however decided to focus
 568 Manhattan plot (Figure 2B). It should be noticed 571 on the 26 genes that were represented by at least

<i>D. melanogaster</i> Gene (Full Name)	Position on <i>dmel6</i> (in kb)	Number of significant SNPs			
		All C_2 (dist. in bp)	C_2^{WW}	C_2^{EU}	C_2^{AM}
Der-1 (Derlin-1)	chr2L:1,974-1,975	2 (236)	1	-	1
Gdi (GDP dissociation inhibitor)	chr2L:9,492-9,495	4 (342)	4	4	-
lncRNA:CR45693 (long non-coding RNA)	chr2L:14,51-14,512	2 (14)	2	1	-
Tpr2 (tetratricopeptide repeat protein 2)	chr2L:16,492-16,507	2 (8)	-	2	-
Ret (Ret oncogene)	chr2L:21,182-21,199	2 (70)	2	-	-
tou (toutatis)	chr2R:11,579-11,616	2 (18)	1	-	2
jeb (jelly belly)	chr2R:12,091-12,119	2 (14)	2	-	-
CG5065	chr2R:16,608-16,625	2 (13)	-	2	-
bab2 (bric a brac 2)	chr3L:1,140-1,177	2 (11189)	1	-	1
axo (axotactin)	chr3L:4,630-4,687	2 (25886)	-	1	1
RhoGEF64C (ρ guanine nucl. exch. fact. at 64C)	chr3L:4,693-4,796	2 (8)	2	1	1
CG7509	chr3L:4,803-4,805	2 (5)	-	2	-
Con (connectin)	chr3L:4,938-4,976	2 (616)	1	1	-
Ets65A (Ets at 65A)	chr3L:6,098-6,124	2 (27998)	1	1	-
lncRNA:CR45759 (long non-coding RNA)	chr3L:6,787-6,787	4 (106)	-	-	4
ome (omega)	chr3L:14,673-14,748	2 (1)	2	-	-
sa (spermatocyte arrest)	chr3L:21,405-21,407	2 (61)	1	1	-
yellow-e (yellow-e)	chr3R:13,410-13,415	3 (33)	3	-	1
cv-c (crossveinless c)	chr3R:14,392-14,482	4 (2737)	1	-	3
osa (osa)	chr3R:17,688-17,718	2 (29)	-	-	2
cpo (couch potato)	chr3R:17,944-18,016	3 (193)	3	2	3
Rh3 (rhodopsin 3)	chr3R:20,081-20,082	2 (5709)	2	1	-
Ctl2 (choline transporter-like 2)	chr3R:29,123-29,128	2 (3)	-	-	2
Syt12 (synaptotagmin 12)	chrX:13,359-13,368	3 (65)	1	-	2
Ac13E (adenylyl cyclase 13E)	chrX:15,511-15,554	4 (19)	-	-	4
Axs (abnormal X segregation)	chrX:16,680-16,684	2 (11)	-	-	2

Table 1. Description of the 26 orthologous *D. melanogaster* genes represented by at least two of the 204 SNPs found significant for one of the three contrast analyses, C_2^{WW} (6 native vs. 16 invasive populations), C_2^{EU} (6 native vs. 8 invasive populations of the European invasion route) and C_2^{AM} (6 native vs. 8 invasive populations of the American invasion route). The third column gives the overall number of significant SNPs (at the 1% q-value threshold) and their maximal spacing in bp (on the *D. suzukii* assembly). Columns 4 to 6 gives the number of significant SNPs for each of the three contrast analyses.

572 two SNPs significant in one of the three contrast 585 *RhoGEF64C* with one SNP and *cpo* with two
573 analyses; see Table 1 for details. The significant 586 SNPs. Such convergent signals of association
574 SNPs underlying the different genes tended to be 587 with invasive status in the two independent
575 very close, spanning a few bp (span > 1kb for only 588 invasion routes were particularly convincing. The
576 five genes). In particular, we observed doublet 589 median allele frequencies (computed from raw
577 variants (i.e., adjacent SNPs in complete LD) 590 read counts) for the reference allele underlying the
578 within three genes (*cpo*, *ome* and *lnc:CR45759*). 591 corresponding *RhoGEF64C* significant SNP was
579 Among these 26 candidate genes, 10 and 12 592 0.09 (from 0.00 to 0.44) in the native populations
580 might be considered as specific to the European 593 compared to 0.93 (from 0.90 to 0.98) and 0.87
581 and American invasion routes, respectively, since 594 (from 0.59 to 1.00) in the invasive populations
582 they did not contain any SNP significant for the 595 of the European and American invasion routes,
583 alternative contrasts. Only two genes contained 596 respectively (Table S2). Similarly, the two SNPs
584 SNPs significant in all three contrast analyses: 597 significant for the three contrast analyses in the

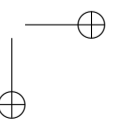
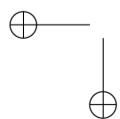


598 *cpo* gene actually formed a doublet with a median 630 may seem surprising since the binary trait under
599 reference allele frequency of 0.20 (from 0.02 to 631 study (invasive versus native) is complex in the
600 0.33) in the native populations compared to 632 sense that numerous biological differences may
601 0.99 (from 0.91 to 1.00, excluding the outlying 633 characterize invasive and native populations. The
602 Hawaiian population) in the invasive populations 634 invasion process itself, including the associated
603 of the European and American invasion routes, 635 selective pressures and the genetic composition
604 respectively (Table S2). Finally, for both the 636 of the source populations, may actually differ
605 genes *RhoGEF64C* and *cpo*, all *D. suzukii* 637 depending on the considered invaded areas. Hence
606 extended sequences underlying the corresponding 638 the small number of SNPs showing strong signals
607 SNPs aligned within potentially rapidly evolving 639 of association with the invasive status may stem
608 intronic sequences. These sequences nevertheless 640 from the integrative nature of our analysis over
609 displayed substantial similarities with other 641 a large number of invasive populations from
610 related drosophila species, as shown in Figure S6 642 different invasion routes. The genomic features
611 for the gene *cpo*. 643 that may be identified under this evolutionary

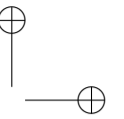
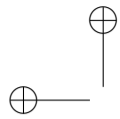
612 Discussion 644

613 We characterized the genome response of 645 genetic changes instrumental to invasions shared
614 *D. suzukii* during its worldwide invasion by 646 by a majority of populations. Accordingly, it
615 conducting a genome-wide scan for association 647 is worth noting that the independent contrast
616 with the invasive or native status of the sampled 648 analyses of the two main invasion routes
617 populations. To that end, we relied on the 649 (i.e. the American and the European routes)
618 newly developed C_2 statistic that was aimed at 650 point to substantially different subsets of SNPs
619 identifying significant allele frequencies differences 651 significantly associated with the invasive status
620 between two contrasting groups of populations 652 of the populations. This suggests that the source
621 while accounting for their overall correlation 653 populations and some aspects of the invasion
622 structure due to the shared population history. 654 process differ in the two invaded areas. This could
623 Our approach identified genomic regions and 655 however also reflect the presumably polygenic
624 candidate genes most likely involved in adaptive 656 nature of the traits underlying invasion success
625 processes underlying the invasion success of *D.* 657 since the evolutionary trajectories of complex
626 *suzukii*. 658 traits may rely on different combination of

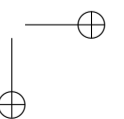
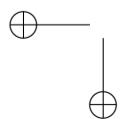
627 Overall, we found that a relatively small number 659 favorable genetic variants.
628 of SNPs were significantly associated with the 660 The availability of a high quality genome
629 invasive status of *D. suzukii* populations. This 661 assembly of *D. suzukii* (Paris *et al.*, 2019) and

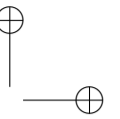
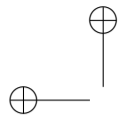


662 a large amount of genomic resources for its 694 species *Drosophila montana* (Kankare *et al.*, 2010;
663 sister model species *D. melanogaster* allowed 695 Schmidt *et al.*, 2008). Moreover, indirect action
664 identifying a set of genes associated with the 696 of selection on diapause, by means of genetic
665 invasive status of populations. A subset of 697 correlations involving *cpo* genetic variation, was
666 those genes was associated with physiological 698 found on numerous other life-history traits in *D.*
667 functions and traits previously documented in *D.* 699 *melanogaster* (Schmidt and Paaby, 2008; Schmidt
668 *melanogaster*, but for most of them, functional 700 *et al.*, 2005). Specifically, compared to diapausing
669 and phenotypic studies turned out to be limited. 701 populations, non-diapausing populations had a
670 Their putative role in explaining the invasion 702 shorter development time and higher early
671 success thus remained largely elusive. To avoid 703 fecundity, but also lower rates of larval and adult
672 too speculative interpretations (Pavlidis *et al.*, 704 survival and lower levels of cold resistance.
673 2012), we will not elaborate further on the 705 Both theoretical (Roughgarden, 1971) and
674 candidate genes. Yet, we did notice that long 706 experimental (Mueller and Ayala, 1981) evidence
675 non-coding RNAs represent more than 10% (14 707 show that traits typical for colonization (i.e., the
676 out of 130) of our candidate genes, a feature 708 so-called r-traits; Charlesworth, 1994), such as
677 which may underline a critical role of variants 709 a non-diapausing phenotype, are selected when
678 involved in gene regulation to promote short- 710 a population evolves in a new habitat with low
679 term response to adaptive constraints during 711 densities and low levels of competition. Common
680 invasion. Also, two genes *RhoGEF64C* and *cpo* 712 garden studies are needed to assess potential
681 contained SNPs that were found to be highly 713 differences in key life history traits (including
682 significantly associated with the invasive status 714 diapause induction and correlated traits) between
683 in both the European and American invasion 715 native and invasive populations of *D. suzukii* and
684 routes. While the function of the *RhoGEF64C* 716 to evaluate to which extent these are related to
685 gene has so far not been extensively studied, 717 the identified variants (including those within the
686 several functional and phenotypic studies in other 718 *cpo* gene) differentiating the native and invasive
687 *Drosophila* species identified genetic variation 719 populations of this species.
688 in the *cpo* gene associated with traits possibly 720 The C_2 statistic we developed in the present
689 important for invasion success. For instance, *cpo* 721 study appears particularly well suited to search
690 genetic variation was found to contribute to 722 for association with population-specific binary
691 natural variation in diapause in *D. melanogaster* 723 traits. Apart from the invasive vs. native
692 populations of a North American cline and 724 status we studied in *D. suzukii*, numerous
693 in populations from the more distantly related 725 examples can be found where adaptive constraints



726 may be formulated in terms of contrasting 758 context of large data sets since it allows to deal
727 binary population features, including individual 759 with multiple testing issues by controlling for FDR
728 resistance or sensibility to pathogens or host- 760 (Francois *et al.*, 2016), via, e.g., the estimation of
729 defense systems (e.g., Eoche-Bosy *et al.*, 2017), 761 q-values (Storey and Tibshirani, 2003).
730 high vs. low altitude adaptation (e.g., Foll *et al.*, 762 To estimate the C_2 statistic, we needed to
731 2014), ecotypes of origin (e.g., Roesti *et al.*, 763 correct allele frequencies for population structure.
732 2015; Westram *et al.*, 2014), or domesticated 764 To that end, we relied on the Bayesian hierarchical
733 vs. wild status (e.g., Alberto *et al.*, 2018). In 765 model implemented in the software BAYPASS
734 our simulation study, the performance of the 766 that has several valuable properties including (i)
735 C_2 statistic was similar to that of a standard 767 the accurate estimation of the scaled covariance
736 BF obtained after assuming a linear relationship 768 matrix of population allele frequencies (Ω), (ii)
737 between the (standardized) population allele 769 the integration over the uncertainty of the
738 frequencies and their corresponding binary status. 770 across population allele frequencies (π parameter),
739 It is worth stressing, however, that C_2 has several 771 and (iii) the inclusion of additional layers of
740 critical advantages over BF, as well as over any 772 complexities such as the sampling of reads
741 other decision criterion that may be derived from 773 from (unobserved) allele counts in Pool-Seq data
742 a parametric modeling. 774 (Gautier, 2015). A cost of Bayesian hierarchical
743 From a practical point of view, the C_2 775 modeling is however to shrink the posterior
744 estimation does not require inclusion of any other 776 means of the model parameters and related
745 model parameters making it more robust when 777 statistics such as the C_2 and XtX differentiation
746 dealing with data sets including a small number of 778 statistics (Gelman *et al.*, 2003). To ensure proper
747 populations (e.g., <8 populations), the later type 779 calibration of the two corresponding estimates,
748 of data sets often leading to unstable estimates 780 we then needed to rely on the rescaled posterior
749 of BF (unpublished results). In addition, it may 781 means of the standardized allele frequencies.
750 easily be derived from only a subset of the 782 This empirical procedure proved efficient in
751 populations under study (as we did here when 783 providing well behaved p-values while avoiding
752 computing the C_2^{EU} and C_2^{AM} contrasts specific 784 computationally intensive calibration procedure
753 to each of the two invasion routes), while using 785 based on the analysis of pseudo-observed data sets
754 the complete design to capture more accurate 786 simulated under the generative model (Gautier,
755 information about the shared population history. 787 2015). Still, this did not allow accounting for the
756 Last, the χ^2 calibration of the C_2 under the null 788 uncertainty of the allele frequencies estimation
757 hypothesis represents an attractive property in the 789 (i.e., their full marginal distribution) and more



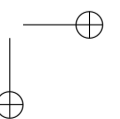
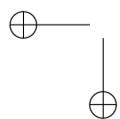


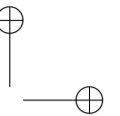
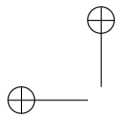
790 importantly, it implicitly assumes exchangeability 822 population allele frequencies (Ω), the S_B statistic
791 of SNPs both across the populations and along the 823 relies on a covariance matrix called F (Refoyo-
792 genome. Such an assumption, which pertains to 824 Martinez *et al.*, 2019) that specifies an a priori
793 the null hypothesis of neutral differentiation only 825 inferred admixture graph summarizing the history
794 (and consequently of no association with binary 826 of the sampled populations. The covariance
795 population-specific covariable), might actually be 827 matrix F thus represents a simplified version of
796 viewed as conservative even in the presence of 828 Ω that may only partially capture the covariance
797 background LD across the populations, providing 829 structure of the population allele frequencies. In
798 that a reasonably large number of SNPs is 830 addition, to compute S_B , the graph root allele
799 analyzed. Interestingly, the almost absence of 831 frequencies are estimated as the average of allele
800 clustering of associated SNPs we observed in 832 frequencies across the sampled population, which
801 the *D. sukukii* genome suggested a very limited 833 might result in biased estimates, particularly when
802 extent of across-population LD, presumably 834 the graph is unbalanced. Deriving the matrix
803 resulting from large effective population sizes. 835 F from Ω (e.g., Pickrell and Pritchard, 2012)
804 This conversely led to a high mapping resolution. 836 might actually allow interpreting C_2 as a Bayesian
805 In practice, when dealing with large data sets, 837 counterpart of the S_B statistic, thereby benefiting
806 a sub-sampling strategy consisting in analyzing 838 from the aforementioned advantages regarding the
807 data sets thinned by marker position also 839 estimation of the parameters Ω and π and allowing
808 allows further reduction of across-population LD 840 proper analysis of Pool-Seq data.
809 (Gautier *et al.*, 2018). Finally, it should be
810 noticed that information from LD might be at 841
811 least partially recovered by combining C_2 or XtX 842
812 derived p-values into local scores (Fariello *et al.*, 843
813 2017). 844

814 Other less computationally intensive (but 845
815 less flexible and versatile) approaches may be 846
816 considered to estimate the C_2 statistic. For 847
817 instance, the C_2 statistic is closely related to the 848
818 S_B statistic recently proposed by Refoyo-Martinez 849
819 *et al.* (2019) to identify footprints of selection in 850
820 admixture graphs. However, while the C_2 statistic 851
821 relies on the full scaled covariance matrix of 852

Conclusion and perspectives

Our genome-wide association approach allowed identifying genomic regions and genes most likely involved in adaptive processes underlying the invasion success of *D. sukukii*. The approach can be transposed to any other invasive species, and more generally to any species models for which binary traits of interest can be defined at the population level. The major advantage of our approach is that it does not require a preliminary, often extremely laborious, phenotypic characterization of the populations





853 considered (for example using common garden
854 experiments) in order to inform candidate traits
855 for which genomic associations are sought. As
856 a matter of fact, in our association study the
857 populations analyzed are simply classified into two
858 categories: invasive or native.

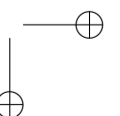
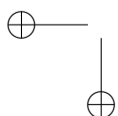
859 The functional and phenotypic interpretation
860 of the signals obtained by our genome scan
861 methods remains challenging. Such interpretation
862 requires a good functional characterization of
863 the genome of the studied species or, failing
864 that, of a closely related species (i.e. *D.*
865 *melanogaster* in our study). Following a strategy
866 sometimes referred to as “reverse ecology”
867 since it goes from gene(s) to phenotype(s) (Li
868 *et al.*, 2008), it is then necessary to test and
869 validate via quantitative genetic experiments
870 whether the inferred candidate traits show
871 significant differences between native and invasive
872 populations. The functional interpretation of the
873 statistical association results can also benefit
874 from experimental validation approaches based
875 on techniques using RNA interference (RNA-
876 silencing, e.g. Janitz *et al.*, 2006) and/or more
877 genome editing approaches (e.g., Karageorgi
878 *et al.*, 2017) targeting the identified candidate
879 variants. Hopefully, such a combination of
880 statistical, molecular and quantitative approaches
881 will provide useful insights into the genomic and
882 phenotypic responses to invasion, and by the
883 same, will help better predict the conditions under
884 which invasiveness can be enhanced or suppressed.

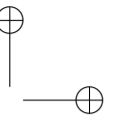
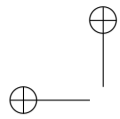
885 **Materials and Methods**

886 Simulation study

887 We used computer simulations to evaluate the
888 performance of the novel statistical framework
889 described in the section *New Approach*. Simulated
890 data sets were generated under the SIMUPOP
891 environment (Peng and Kimmel, 2005) using
892 individual-based forward-in-time simulations
893 implemented on a modified version of the code
894 developed by de Villemereuil *et al.* (2014) for
895 the so-called *HsIMM-C* demographic scenario.
896 This corresponded to an highly structured
897 isolation with migration demographic model
898 (Figure 1A) that was divided in two successive
899 periods: (i) a neutral divergence phase leading
900 to the differentiation of an ancestral population
901 into 16 populations after four successive fission
902 events (at generations $t=50$, $t=150$, $t=200$ and
903 $t=300$); and (ii) an adaptive phase (lasting 200
904 generations) during which individuals of the 16
905 populations were subjected to selective pressures
906 exerted by two environmental constraints (*ec1*
907 and *ec2*), each constraint having two possible
908 modalities (*a* or *b*). We thus had a total of four
909 possible environments in our simulation setting
910 (Figure 1A).

911 All the simulated populations consisted of 500
912 diploid individuals reproducing under random-
913 mating with non-overlapping generations. From
914 generation $t=150$ (with four populations), the
915 migration rate $m_{jj'}$ between two populations j
916 and j' was set to $m_{jj'} = \frac{m}{2^{p-1}}$ where p is the





917 number of populations in the path connecting
918 k to k' in the population tree. The migration
919 rate between the two ancestral populations from
920 generation $t=50$ to $t=150$ was set to $m=0.005$.
921 For illustration purposes, some of the migration
922 edges were displayed in Figure 1A.

923 Following de Villemereuil *et al.* (2014), a
924 simulated genotyping data set consisted of
925 320 individuals (20 per populations) that were
926 genotyped for 5,000 bi-allelic SNPs regularly
927 spread along ten chromosomes of one Morgan
928 length and with a frequency of 0.5 for the reference
929 allele (randomly chosen) in the root population.
930 Polygenic selection acting during the adaptive
931 phase was simulated by choosing 50 randomly
932 distributed SNPs (among the previous 5,000 ones)
933 that influenced individual fitness according to
934 either the *ec1* or *ec2* environmental constraints
935 (with 25 SNPs for *ec1* and 25 SNPs for *ec2*).

936 The fitness of each individual, given its
937 genotype, can be defined at each generation.
938 let $p(o)=j$ ($j=1,\dots,16$) denote the population
939 of origin of individual o ($o=1,\dots,16\times 500$),
940 and $e_k(j)=1$ (respectively $e_k(j)=-1$) if the
941 environmental constraint eck ($k=1,2$) of
942 population j is of type a (respectively b).
943 Let further denote $s_i(k)$ the local selective
944 coefficient of SNP i such that $s_i(k)=0$ if the SNP
945 is neutral with respect to eck and $s_i(k)=0.01$
946 otherwise. The fitness of each individual o (at each
947 generation) given its genotypes at all the SNPs
948 is then defined using a cumulative multiplicative

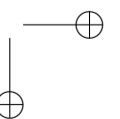
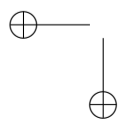
949 fitness function as:

$$w(o) = \prod_{i=1}^I \prod_{k=1}^2 (1 + e_k(p(o))(1 - g_i(o))s_i(k)) \quad (6)$$

950 where $g_i(o)$ is the genotype of individual o at
951 marker i coded as the number of the reference
952 allele (0, 1 or 2).

953 Sampling of *D. suzukii* populations and DNA
954 extraction

955 Adult *D. suzukii* flies were sampled in the field
956 at a total of 22 localities (hereafter termed
957 sample sites) distributed throughout most of
958 the native and invasive range of the species
959 (Fig 1 and Table S1). Samples were collected
960 between 2013 and 2016 using baited traps (with
961 a vinegar-alcohol-sugar mixture) and sweep nets,
962 and stored in ethanol. Only four of the 22
963 samples were composed of flies which directly
964 emerged in the lab from fruits collected in the
965 field (Table S1). Native Asian samples consisted
966 of a total of six sample sites including four
967 Chinese and two Japanese localities. Samples from
968 the invasive range were collected in Hawaii (1
969 sample site), Continental US (6 sites), Brazil (1
970 site), Europe (7 sites) and the French island of
971 La Réunion (1 site). The continental US (plus
972 Brazil) and European (plus La Réunion Island)
973 populations are representative of two separate
974 invasion routes (the American and European
975 routes, respectively), with different native source
976 populations and multiple introduction events in
977 both invaded areas (Fraimout *et al.* 2017; see
978 Table S1).



979 Pool sequencing 1011 removed. Filtered reads were then mapped onto
980 For each of the 22 sampling sites, the thoraxes 1012 the newly released WT3-2.0 *D. suzukii* genome
981 of 50 to 100 representative adult flies (Table 1013 assembly (Paris *et al.*, 2019), using default options
982 S1) were pooled for DNA extraction using the 1014 of the *mem* program from the *bwa* 0.7.17 software
983 EZ-10 spin column genomic DNA minipreps kit 1015 (Li, 2013; Li and Durbin, 2009). Read alignments
984 (Bio basic inc.). Barcoded DNA PE libraries 1016 with a mapping quality Phred-score <20 or
985 with insert size of ca. 550 bp were further 1017 PCR duplicates were further removed using the
986 prepared using the Illumina Truseq DNA Library 1018 *view* (option -q 20) and *markdup* programs from
987 Preparation Kit following manufacturer protocols 1019 the *SAMtools* 1.9 software (Li *et al.*, 2009),
988 using the 22 DNA pools samples. The DNA 1020 respectively.
989 libraries were then validated on a DNA1000 1021 Variant calling was then performed on the
990 chip on a Fragment Analyzer (Agilent) to 1022 resulting *mpileup* file using *VarScan mpileup2cns*
991 determine size and quantified by qPCR using 1023 v2.3.4 (Koboldt *et al.*, 2012) (options *-min-*
992 the Kapa library quantification kit to determine 1024 *coverage* 50 *-min-avg-qual* 20 *-min-var-freq* 0.001
993 concentration. The cluster generation process was 1025 *-variants-output-vcf* 1). The resulting *vcf* file was
994 performed on cBot (Illumina) using the Paired- 1026 processed with the *vcf2pooldata* function from the
995 End Clustering kit (Illumina). Each pool DNA 1027 R package *poolfstats* v1.1 (Hivert *et al.*, 2018)
996 library was further paired-end sequenced on a 1028 retaining only bi-allelic SNPs covered by >4
997 HiSeq 2500 (Illumina) using the Sequence by 1029 reads, <99.9th overall coverage percentile in each
998 Synthesis technique (providing 2x125 bp reads, 1030 pool and with an overall MAF>0.01 (computed
999 respectively) with base calling achieved by the 1031 from read counts). In total, n=11,564,472 SNPs
1000 RTA software (Illumina). The Pool-Seq data were 1032 (respectively n=1,966,184 SNPs) SNPs mapping
1001 deposited in the Sequence Read Archive (SRA) 1033 to the autosomal contigs (respectively X-
1002 repository under the BioProject accession number 1034 chromosome contigs) were used for genome-wide
1003 PRJNA576997. 1035 association analysis. The median coverage per
1004 Raw paired-end reads were filtered using *fastp* 1036 pool ranged from 58X to 88X and from 34X to 84X
1005 0.19.4 (Chen *et al.*, 2018) run with default options 1037 for autosomal and X chromosomes, respectively
1006 to remove contaminant adapter sequences and 1038 (Table S2). As previously described (Gautier
1007 trim for poor quality bases (i.e., with a phred- 1039 *et al.*, 2018), the autosomal and X-chromosome
1008 quality score <15). Read pairs with either one 1040 data sets were divided into sub-data sets of ca.
1009 read with a proportion of low quality bases over 1041 75,000 SNPs each (by taking one SNP every 154
1010 40% or containing more than 5 N bases were 1042 SNPs and one SNPs every 26 SNPs along the

underlying autosomal and X-chromosome contigs, respectively).
Genome scan analyses
All genome-wide scans were performed using an upgraded version (2.2) of BAYPASS (Gautier, 2015) (available from <http://www1.montpellier.inra.fr/CBGP/software/baypass/>), that includes the new C_2 and XtX statistics estimated as described in the above section *New Approach*. We always used the BAYPASS core model with default options for the MCMC algorithm to obtain estimates of four statistical items: (i) the scaled covariance matrix (Ω); (ii) the SNP-specific XtX overall differentiation statistic in the form of both \widehat{XtX} , the posterior mean of XtX (Gautier, 2015) and \widehat{XtX}^* , our newly described calibrated estimator; (iii) our novel C_2 statistic in the form of the calibrated estimator described above; and (iv) Bayes Factor reported in deciban units (db) as a measure of support for association with contrasts of each SNP based on a linear regression model (Coop *et al.*, 2010; Gautier, 2015). For BF, a value >15 db (respectively >20 db) provides very strong (respectively decisive) evidence in favor of association according to the Jeffreys' rule (Jeffreys, 1961).

For the *D. sukukii* data sets, we specified the pool haploid sample sizes, for either autosomes or the X-chromosome (Table S1), to activate the Pool-Seq mode of BAYPASS. The C_2^{WW} statistic for the contrast of the six native and

16 worldwide invasive populations was estimated jointly with the C_2^{EU} and C_2^{AM} statistics for the contrast of the six native and eight invasive populations of the European and American invasion routes, respectively. For these two latter estimates, this simply amounted to setting $c_j=0$ (see eq. 3) for all population j not considered in the corresponding contrast analysis. Finally, two additional independent runs (using the option -seed) were performed to assess reproducibility of the MCMC estimates. We found a fairly high correlation across the different independent runs (Pearson's $r > 0.92$ for autosomal and $r > 0.87$ and X-chromosome data) for the different estimators and thus only presented results from the first run. Similarly and for each chromosome type (i.e., autosomes or the X chromosome), a near perfect correlation of the posterior means of the estimated Ω matrix elements was observed across independent runs as well as within each run across SNP sub-samples, with the corresponding FMD distances (Gautier, 2015) being always smaller than 0.4. We thus only reported results regarding the Ω matrix that were obtained from a single randomly chosen sub-data set analysed in the first run.

Acknowledgments

AE, MG and LO acknowledge financial support from the National Research Fund ANR (France) through the project ANR-16-CE02-0015-01 (SWING), the Languedoc-Roussillon region

(France) through the European Union program
FEFER FSE IEJ 2014-2020 (project CPADROL)
and the INRA scientific department SPE (AAP-
SPE 2016 and 2018). MGX acknowledges financial
support from France Génomique National
infrastructure, funded as part of “Investissement
d’avenir” program managed by Agence Nationale
pour la Recherche (contract ANR-10-INBS-09).
We are grateful to the genotoul bioinformatics
platform Toulouse Midi-Pyrenees for providing
computing resources, Nicolas Rode for useful
discussions and comments on a previous version
of the manuscript and Nicolas Ris, Jon Koch,
Masahito Kimura, Simon Fellous, Vincent
Debat, Marta Pascual, Ruth Hufbauer, Marindia
Depra, Isabel Martinez, Pierre Girod and Maxi
Richmond for help in collecting some of the *D.*
suzukii samples.

References

Adrion, J. R., Kousathanas, A., Pascual, M., Burrack, H. J.,
Haddad, N. M., Bergland, A. O., Machado, H., Sackton,
T. B., Schlenke, T. A., Watada, M., Wegmann, D.,
and Singh, N. D. 2014. *Drosophila suzukii*: the genetic
footprint of a recent, worldwide invasion. *Mol. Biol. Evol.*, 31(12): 3148–63.

Alberto, F. J., Boyer, F., Orozco-terWengel, P., Streeter, I.,
Servin, B., de Villemereuil, P., Benjelloun, B., Librado,
P., Biscarini, F., Colli, L., Barbato, M., Zamani,
W., Alberti, A., Engelen, S., Stella, A., Joost, S.,
Ajmone-Marsan, P., Negrini, R., Orlando, L., Rezaei,
H. R., Naderi, S., Clarke, L., Flicek, P., Wincker,
P., Coissac, E., Kijas, J., Tosser-Klopp, G., Chikhi,
A., Bruford, M. W., Taberlet, P., and Pompanon, F.
2018. Convergent genomic signatures of domestication
in sheep and goats. *Nat Commun*, 9(1): 813.

Asplen, M. K., Anfora, G., Biondi, A., Choi, D.-S., Chu, D.,
Daane, K. M., Gibert, P., Gutierrez, A. P., Hoelmer,
K. A., Hutchison, W. D., Isaacs, R., Jiang, Z.-L.,
Krpti, Z., Kimura, M. T., Pascual, M., Philips, C. R.,
Plantamp, C., Ponti, L., Vtek, G., Vogt, H., Walton,
V. M., Yu, Y., Zappal, L., and Desneux, N. 2015.
Invasion biology of spotted wing drosophila (*drosophila*
suzukii): a global perspective and future priorities.
Journal of Pest Science, 88(3): 469–494.

Balanya, J., Oller, J. M., Huey, R. B., Gilchrist, G. W.,
and Serra, L. 2006. Global genetic change tracks global
climate warming in *drosophila subobscura*. *Science*,
313(5794): 1773–1775.

Barrett, S. C. H. 2015. Foundations of invasion genetics:
the baker and stebbins legacy. *Molecular Ecology*, 24(9):
1927–1941.

Bock, D. G., Caseys, C., Cousens, R. D., Hahn, M. A.,
Heredia, S. M., Hubner, S., Turner, K. G., Whitney,
K. D., and Rieseberg, L. H. 2015. What we still don’t
know about invasion genetics. *Molecular ecology*, 24(9):
2277–2297.

Bonhomme, M., Chevalet, C., Servin, B., Boitard, S.,
Abdallah, J., Blott, S., and Sancristobal, M. 2010.
Detecting selection in population trees: the lewontin and
krakauer test extended. *Genetics*, 186(1): 241–262.

Charlesworth, B. 1994. *Evolution in Age-Structured*
Populations. Cambridge Studies in Mathematical
Biology. Cambridge University Press, 2 edition.

Chen, S., Zhou, Y., Chen, Y., and Gu, J. 2018. fastp: an
ultra-fast all-in-one fastq preprocessor. *Bioinformatics*,
34(17): i884–i890.

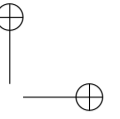
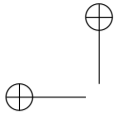
Chiu, J. C., Jiang, X., Zhao, L., Hamm, C. A., Cridland,
J. M., Saelao, P., Hamby, K. A., Lee, E. K., Kwok, R. S.,
Zhang, G., Zalom, F. G., Walton, V. M., and Begun,
D. J. 2013. Genome of *drosophila suzukii*, the spotted
wing drosophila. *G3 (Bethesda)*, 3(12): 2257–71.

Cini, A., Ioriatti, C., and Anfora, G. 2012. A review
of the invasion of *drosophila suzukii* in europe and a

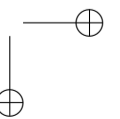
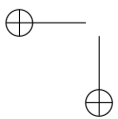
- 1179 draft research agenda for integrated pest management. 1218 4700–4711.
- 1180 *Bulletin of Insectology*, 65: 149–160. 1219
- 1181 Clemente, F., Gautier, M., and Vitalis, R. 2018. Inferring 1220
1182 sex-specific demographic history from snp data. *PLoS* 1221
1183 *Genet.*, 14(1): e1007191. 1222
- 1184 Colautti, R. I. and Barrett, S. C. H. 2013. Rapid adaptation 1223
1185 to climate facilitates range expansion of an invasive 1224
1186 plant. *Science*, 342(6156): 364–6. 1225
- 1187 Colautti, R. I. and Lau, J. A. 2015. Contemporary 1226
1188 evolution during invasion: evidence for differentiation, 1227
1189 natural selection, and local adaptation. *Molecular* 1228
1190 *Ecology*, 24(9): 1999–2017. 1229
- 1191 Coop, G., Witonsky, D., Rienzo, A. D., and Pritchard, J. K. 1230
1192 2010. Using environmental correlations to identify loci 1231
1193 underlying local adaptation. *Genetics*, 185(4): 1411– 1232
1194 1423. 1233
- 1195 de Villemereuil, P. and Gaggiotti, O. E. 2015. A new 1234
1196 FST-based method to uncover local adaptation using 1235
1197 environmental variables. *Methods in Ecology and* 1236
1198 *Evolution*, 6(11): 1248–1258. 1237
- 1199 de Villemereuil, P., Frichot, E., Bazin, E., Francois, O., and 1238
1200 Gaggiotti, O. E. 2014. Genome scan methods against 1239
1201 more complex models: when and how much should we 1240
1202 trust them? *Mol Ecol*, 23(8): 2006–2019. 1241
- 1203 Dlugosch, K. M., Anderson, S. R., Braasch, J., Cang, 1242
1204 F. A., and Gillette, H. D. 2015. The devil is in the 1243
1205 details: genetic variation in introduced populations and 1244
1206 its contributions to invasion. *Molecular Ecology*, 24(9): 1245
1207 2095–2111. 1246
- 1208 Ellstrand, N. C. and Schierenbeck, K. A. 2000. 1247
1209 Hybridization as a stimulus for the evolution of 1248
1210 invasiveness in plants? *Proceedings of the National* 1249
1211 *Academy of Sciences*, 97(13): 7043–7050. 1250
- 1212 Eoche-Bosy, D., Gautier, M., Esquibet, M., Legeai, F., 1251
1213 Bretaudeau, A., Bouchez, O., Fournet, S., Grenier, 1252
1214 E., and Montarry, J. 2017. Genome scans on 1253
1215 experimentally evolved populations reveal candidate 1254
1216 regions for adaptation to plant resistance in the potato 1255
1217 cyst nematode *globodera pallida*. *Mol. Ecol.*, 26(18): 1256
- 1219 Estoup, A., Ravigne, V., Hufbauer, R., Vitalis, R., Gautier,
1220 M., and Facon, B. 2016. Is there a genetic paradox
1221 of biological invasion? *Annual Review of Ecology,*
1222 *Evolution, and Systematics*, 47(1): 51–72.
- 1223 Facon, B., Hufbauer, R. A., Tayeh, A., Loiseau, A.,
1224 Lombaert, E., Vitalis, R., Guillemaud, T., Lundgren,
1225 J. G., and Estoup, A. 2011. Inbreeding depression is
1226 purged in the invasive insect *harmonia axyridis*. *Curr.*
1227 *Biol.*, 21(5): 424–7.
- 1228 Fariello, M. I., Boitard, S., Mercier, S., Robelin, D., Faraut,
1229 T., Arnould, C., Recoquillay, J., Bouchez, O., Salin,
1230 G., Dehais, P., Gourichon, D., Leroux, S., Pitel, F.,
1231 Leterrier, C., and SanCristobal, M. 2017. Accounting
1232 for linkage disequilibrium in genome scans for selection
1233 without individual genotypes: The local score approach.
1234 *Mol. Ecol.*, 26(14): 3700–3714.
- 1235 Foll, M., Gaggiotti, O. E., Daub, J. T., Vatsiou, A., and
1236 Excoffier, L. 2014. Widespread signals of convergent
1237 adaptation to high altitude in asia and america. *Am. J.*
1238 *Hum. Genet.*, 95(4): 394–407.
- 1239 Fraimout, A., Debat, V., Fellous, S., Hufbauer, R. A.,
1240 Foucaud, J., Pudlo, P., Marin, J.-M., Price, D. K.,
1241 Cattel, J., Chen, X., Depra, M., Duyck, P. F., Guedot,
1242 C., Kenis, M., Kimura, M. T., Loeb, G., Loiseau, A.,
1243 Martinez-Sanudo, I., Pascual, M., Richmond, M. P.,
1244 Shearer, P., Singh, N., Tamura, K., Xureb, A., Zhang, J.,
1245 and Estoup, A. 2017. Deciphering the routes of invasion
1246 of *drosophila suzukii* by means of abc random forest.
1247 *Mol. Biol. Evol.*, 34(4): 980–996.
- 1248 Francois, O., Martins, H., Caye, K., and Schoville, S. D.
1249 2016. Controlling false discoveries in genome scans for
1250 selection. *Mol. Ecol.*, 25(2): 454–69.
- 1251 Frichot, E., Schoville, S. D., Bouchard, G., and Francois,
1252 O. 2013. Testing for associations between loci and
1253 environmental gradients using latent factor mixed
1254 models. *Mol Biol Evol*, 30(7): 1687–1699.
- 1255 Frichot, E., Schoville, S. D., de Villemereuil, P., Gaggiotti,
1256 O. E., and Francois, O. 2015. Detecting adaptive

- 1257 evolution based on association with ecological gradients: 1296
1258 Orientation matters! *Heredity (Edinb)*. 1297
1259 Gautier, M. 2015. Genome-wide scan for adaptive 1298
1260 divergence and association with population-specific 1299
1261 covariates. *Genetics*, 201(4): 1555–79. 1300
1262 Gautier, M., Foucaud, J., Gharbi, K., Cezard, T., Galan, 1301
1263 M., Loiseau, A., Thomson, M., Pudlo, P., Kerdelhue, 1302
1264 C., and Estoup, A. 2013. Estimation of population allele 1303
1265 frequencies from next-generation sequencing data: pool- 1304
1266 versus individual-based genotyping. *Mol Ecol*, 22(14): 1305
1267 3766–3779. 1306
1268 Gautier, M., Yamaguchi, J., Foucaud, J., Loiseau, A., 1307
1269 Ausset, A., Facon, B., Gschloessl, B., Lagnel, J., 1308
1270 Loire, E., Parrinello, H., Severac, D., Lopez-Roques, 1309
1271 C., Donnadiou, C., Manno, M., Berges, H., Gharbi, K., 1310
1272 Lawson-Handley, L., Zang, L.-S., Vogel, H., Estoup, A., 1311
1273 and Prud'homme, B. 2018. The genomic basis of color 1312
1274 pattern polymorphism in the harlequin ladybird. *Curr.* 1313
1275 *Biol.*, 28(20): 3296–3302.e7. 1314
1276 Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. 1315
1277 2003. *Bayesian Data Analysis, Second Edition*. CRC 1316
1278 Press. 1317
1279 Grau, J., Grosse, I., and Keilwagen, J. 2015. Proc: 1318
1280 computing and visualizing precision-recall and receiver 1319
1281 operating characteristic curves in r. *Bioinformatics*, 1320
1282 31(15): 2595–7. 1321
1283 Gunther, T. and Coop, G. 2013. Robust identification 1322
1284 of local adaptation from allele frequencies. *Genetics*, 1323
1285 195(1): 205–220. 1324
1286 Hivert, V., Leblois, R., Petit, E. J., Gautier, M., and Vitalis, 1325
1287 R. 2018. Measuring genetic differentiation from pool-seq 1326
1288 data. *Genetics*, 210(1): 315–330. 1327
1289 Hoskins, R. A., Carlson, J. W., Wan, K. H., Park, S., 1328
1290 Mendez, I., Galle, S. E., Booth, B. W., Pfeiffer, B. D., 1329
1291 George, R. A., Svirskas, R., Krzywinski, M., Schein, J., 1330
1292 Accardo, M. C., Damia, E., Messina, G., Mndez-Lago, 1331
1293 M., de Pablos, B., Demakova, O. V., Andreyeva, E. N., 1332
1294 Boldyreva, L. V., Marra, M., Carvalho, A. B., Dimitri, 1333
1295 P., Villasante, A., Zhimulev, I. F., Rubin, G. M., 1334
Karpen, G. H., and Celniker, S. E. 2015. The release
6 reference sequence of the drosophila melanogaster
genome. *Genome Res.*, 25(3): 445–58.
Janitz, M., Vanhecke, D., and Lehrach, H. 2006. *High-
Throughput RNA Interference in Functional Genomics*,
pages 97–104. Springer Berlin Heidelberg, Berlin,
Heidelberg.
Jeffreys, H. 1961. *Theory of Probability*. Oxford University
Press, 3rd edition.
Kankare, M., Salminen, T., Laiho, A., Vesala, L., and
Hoikkala, A. 2010. Changes in gene expression linked
with adult reproductive diapause in a northern malt fly
species: a candidate gene microarray study. *BMC Ecol.*,
10: 3.
Karageorgi, M., Bracker, L. B., Lebreton, S., Minervino, C.,
Cavey, M., Siju, K. P., Kadow, I. C. G., Gompel, N., and
Prud'homme, B. 2017. Evolution of multiple sensory
systems drives novel egg-laying behavior in the fruit pest
drosophila suzukii. *Curr. Biol.*, 27(6): 847–853.
Koboldt, D. C., Zhang, Q., Larson, D. E., Shen, D.,
McLellan, M. D., Lin, L., Miller, C. A., Mardis, E. R.,
Ding, L., and Wilson, R. K. 2012. VarScan 2: somatic
mutation and copy number alteration discovery in
cancer by exome sequencing. *Genome Res.*, 22(3):
568–76.
Lee, C. E. 2002. Evolutionary genetics of invasive species.
Trends in ecology and evolution, 17(8): 386–391.
Lee, C. E. and Gelembiuk, G. W. 2008. Evolutionary origins
of invasive populations. *Evolutionary Applications*, 1(3):
427–448.
Li, H. 2013. Aligning sequence reads, clone sequences and
assembly contigs with bwa-mem. *arXiv*, 1303.3997.
Li, H. and Durbin, R. 2009. Fast and accurate
short read alignment with burrows-wheeler transform.
Bioinformatics, 25(14): 1754–60.
Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan,
J., Homer, N., Marth, G., Abecasis, G., and Durbin,
R. 2009. The sequence alignment/map format and
samtools. *Bioinformatics*, 25(16): 2078–9.

- 1335 Li, Y. F., Costello, J. C., Holloway, A. K., and Hahn, M. W. 1374
1336 2008. "reverse ecology" and the power of population 1375
1337 genomics. *Evolution*, 62(12): 2984–2994. 1376
- 1338 Mueller, L. D. and Ayala, F. J. 1981. Trade-off between 1377
1339 r-selection and k-selection in drosophila populations. 1378
1340 *Proc. Natl. Acad. Sci. U.S.A.*, 78(2): 1303–5. 1379
- 1341 Ochocki, B. M. and Miller, T. E. X. 2017. Rapid evolution 1380
1342 of dispersal ability makes biological invasions faster and 1381
1343 more variable. *Nat Commun*, 8: 14315. 1382
- 1344 Ometto, L., Cestaro, A., Ramasamy, S., Grassi, A., Revadi, 1383
1345 S., Siozios, S., Moretto, M., Fontana, P., Varotto, C., 1384
1346 Pisani, D., Dekker, T., Wrobel, N., Viola, R., Pertot, 1385
1347 I., Cavalieri, D., Blaxter, M., Anfora, G., and Rota- 1386
1348 Stabelli, O. 2013. Linking genomics and ecology 1387
1349 to investigate the complex evolution of an invasive 1388
1350 drosophila pest. *Genome Biol Evol*, 5(4): 745–57. 1389
- 1351 Paris, M., Boyer, R., Jaenichen, R., Wolf, J., Karageorgi, 1390
1352 M., Green, J., Cagnon, M., Parinello, H., Estoup, A., 1391
1353 Gautier, M., Gompel, N., and Prud'homme, B. 2019. 1392
1354 Near-chromosome level genome assembly of the fruit 1393
1355 pest drosophila *suzukii* using long-read sequencing. 1394
1356 *submitted*. 1395
- 1357 Pavlidis, P., Jensen, J. D., Stephan, W., and Stamatakis, A. 1396
1358 2012. A critical assessment of storytelling: gene ontology 1397
1359 categories and the importance of validating genomic 1398
1360 scans. *Mol. Biol. Evol.*, 29(10): 3237–48. 1399
- 1361 Peng, B. and Kimmel, M. 2005. simupop: a forward- 1400
1362 time population genetics simulation environment. 1401
1363 *Bioinformatics*, 21(18): 3686–3687. 1402
- 1364 Pickrell, J. K. and Pritchard, J. K. 2012. Inference of 1403
1365 population splits and mixtures from genome-wide allele 1404
1366 frequency data. *PLoS Genet.*, 8(11): e1002967. 1405
- 1367 Puzey, J. and Vallejo-Marin, M. 2014. Genomics 1406
1368 of invasion: diversity and selection in introduced 1407
1369 populations of monkeyflowers (*mimulus guttatus*). *Mol.* 1408
1370 *Ecol.*, 23(18): 4472–85. 1409
- 1371 Refoyo-Martinez, A., da Fonseca, R. R., Halldrsdttir, 1410
1372 K., Arnason, E., Mailund, T., and Racimo, F. 2019. 1411
1373 Identifying loci under positive selection in complex 1412
population histories. *Genome Res.*, 29(9): 1506–1520.
- Reznick, D. N., Losos, J., and Travis, J. 2019. From low
to high gear: there has been a paradigm shift in our
understanding of evolution. *Ecol. Lett.*, 22(2): 233–244.
- Roesti, M., Kueng, B., Moser, D., and Berner, D. 2015.
The genomics of ecological vicariance in threespine
stickleback fish. *Nat Commun*, 6: 8767.
- Roughgarden, J. 1971. Density-dependent natural
selection. *Ecology*, 52(3): 453–468.
- Schlotterer, C., Tobler, R., Kofler, R., and Nolte, V. 2014.
Sequencing pools of individuals - mining genome-wide
polymorphism data without big funding. *Nat Rev
Genet*, 15(11): 749–763.
- Schmidt, P. S. and Paaby, A. B. 2008. Reproductive
diapause and life-history clines in north american
populations of drosophila melanogaster. *Evolution*,
62(5): 1204–15.
- Schmidt, P. S., Matzkin, L., Ippolito, M., and Eanes, W. F.
2005. Geographic variation in diapause incidence, life-
history traits, and climatic adaptation in drosophila
melanogaster. *Evolution*, 59(8): 1721–32.
- Schmidt, P. S., Zhu, C.-T., Das, J., Batavia, M., Yang, L.,
and Eanes, W. F. 2008. An amino acid polymorphism
in the couch potato gene forms the basis for climatic
adaptation in drosophila melanogaster. *Proc. Natl.
Acad. Sci. U.S.A.*, 105(42): 16207–11.
- Storey, J. D. and Tibshirani, R. 2003. Statistical
significance for genomewide studies. *Proc. Natl. Acad.
Sci. U.S.A.*, 100(16): 9440–5.
- Wang, M. and Kong, L. 2019. pblat: a multithread blat
algorithm speeding up aligning sequences to genomes.
BMC Bioinformatics, 20(1): 28.
- Weir, B. S. and Cockerham, C. C. 1984. Estimating
f-statistics for the analysis of population structure.
Evolution, 38(6): 1358–1370.
- Welles, S. and Dlugosch, K. 2018. *Population Genomics of
Colonization and Invasion*, pages 1–29.
- Westram, A. M., Galindo, J., Rosenblad, M. A., Grahame,
J. W., Panova, M., and Butlin, R. K. 2014. Do the



1413 same genes underlie parallel phenotypic divergence in
1414 different *littorina saxatilis* populations? *Mol. Ecol.*,
1415 23(18): 4603–16.
1416 Williams, J. L., Kendall, B. E., and Levine, J. M.
1417 2016. Rapid evolution accelerates plant population
1418 spread in fragmented experimental landscapes. *Science*,
1419 353(6298): 482–485.
1420 Wu, N., Zhang, S., Li, X., Cao, Y., Liu, X., Wang, Q., Liu,
1421 Q., Liu, H., Hu, X., Zhou, X. J., James, A. A., Zhang, Z.,
1422 Huang, Y., and Zhan, S. 2019. Fall webworm genomes
1423 yield insights into rapid adaptation of invasive species.
1424 *Nat Ecol Evol*, 3(1): 105–115.



Sample name	Sampling site (Lat.; Long.)	Status	Sampling date	Sampling method	Auto. (X) haploid sample size	Auto. (X) median coverage
CN-Bei	Beijing, China (40.00;116.4)	Native	June 2014	Baited trap	100 (89)	88 (84)
CN-Lia	Liaoyuan, China (42.96;125.1)	Native	Aug. 2014	Baited trap	100 (42)	63 (41)
CN-Nin	Ningbo, China (30.02;121.5)	Native	July 2014 & May 2016	Baited trap	100 (86)	59 (56)
CN-Shi	Shiping county, China (23.7;102.5)	Native	June 2014 & May 2016	Baited trap	100 (53)	61 (34)
JP-Sap	Sapporo, Japan (43.05;141.4)	Native	July 2014	Mullberry	100 (54)	77 (40)
JP-Tok	Tokyo, Japan (35.64;139.4)	Native	June 2016	Mullberry, plum	100 (90)	62 (56)
DE-Dos	Dossenheim, Germany (49.45;8.660)	Invasive (EU)	Aug. 2015	Baited trap	100 (58)	73 (49)
DE-Jen	Jena, Germany (50.93;11.56)	Invasive (EU)	Sept. 2016	Baited trap	200 (150)	70 (56)
ES-Bar	Barcelona, Spain (41.36;1.964)	Invasive (EU)	July 2014	Baited trap	100 (50)	71 (37)
FR-Cor	Corsica, France (42.35;9.003)	Invasive (EU)	Aug. 2016	Baited trap	100 (75)	66 (58)
FR-Lez	Montpellier, France (43.70;3.834)	Invasive (EU)	July 2014	Baited trap	200 (150)	82 (65)
FR-Par	Paris, France (48.84;2.361)	Invasive (EU)	Nov. 2016	Baited trap	200 (150)	65 (54)
FR-Run	La Réunion, France (-21.15;55.64)	Invasive (EU)	Sept. 2016	Cattley guava	200 (150)	86 (68)
IT-Tre	Trento, Italy (46.04;11.15)	Invasive (EU)	Sept. 2014	Baited trap	200 (140)	63 (48)
BR-Pal	Porto Alegre, Brazil (-27.72;-52.17)	Invasive (AM)	July 2014	Baited trap	100 (67)	68 (53)
US-Col	Fort Collins, USA (40.57;-105.1)	Invasive (AM)	Sept. 2015	Baited trap	100 (74)	72 (59)
US-Haw	Hawaii (Hilo), USA (19.67;-155.5)	Invasive (AM)	June 2016	Baited trap	100 (75)	87 (71)
US-Nca	Raleigh, USA (35.70;-80.62)	Invasive (AM)	Oct. 2016	Raspberries,Blackberries	200 (150)	67 (54)
US-Sdi	San-Diego, USA (32.72;-117.2)	Invasive (AM)	May 2014	Baited trap	100 (68)	82 (61)
US-Sok	Dayton, USA (45.22;-123.1)	Invasive (AM)	Oct. 2014	Baited trap, sweep net	150 (95)	58 (38)
US-Wat	Watsonville, USA (36.90;-121.8)	Invasive (AM)	Oct. 2014	Sweep net	100 (54)	65 (37)
US-Wis	Barneveld, USA (42.97;-89.69)	Invasive (AM)	Nov. 2016	Baited trap	150 (120)	70 (58)

TAB. S1 : Description of the 22 *D. suzukii* population samples. The populations representative of the European and American invasion routes are denoted Invasive (EU) and Invasive (AM), respectively (column 3). For each population sample, the thoraxes of 50 to 100 adult flies were pooled; hence the haploid sample sizes of autosomal loci ranging from 100 to 200 (column 7). Pool-samples included both females and male adults, with different proportions of the two sexes depending on the sample; hence the variable number of haploid sample sizes for the X chromosome (column 7).

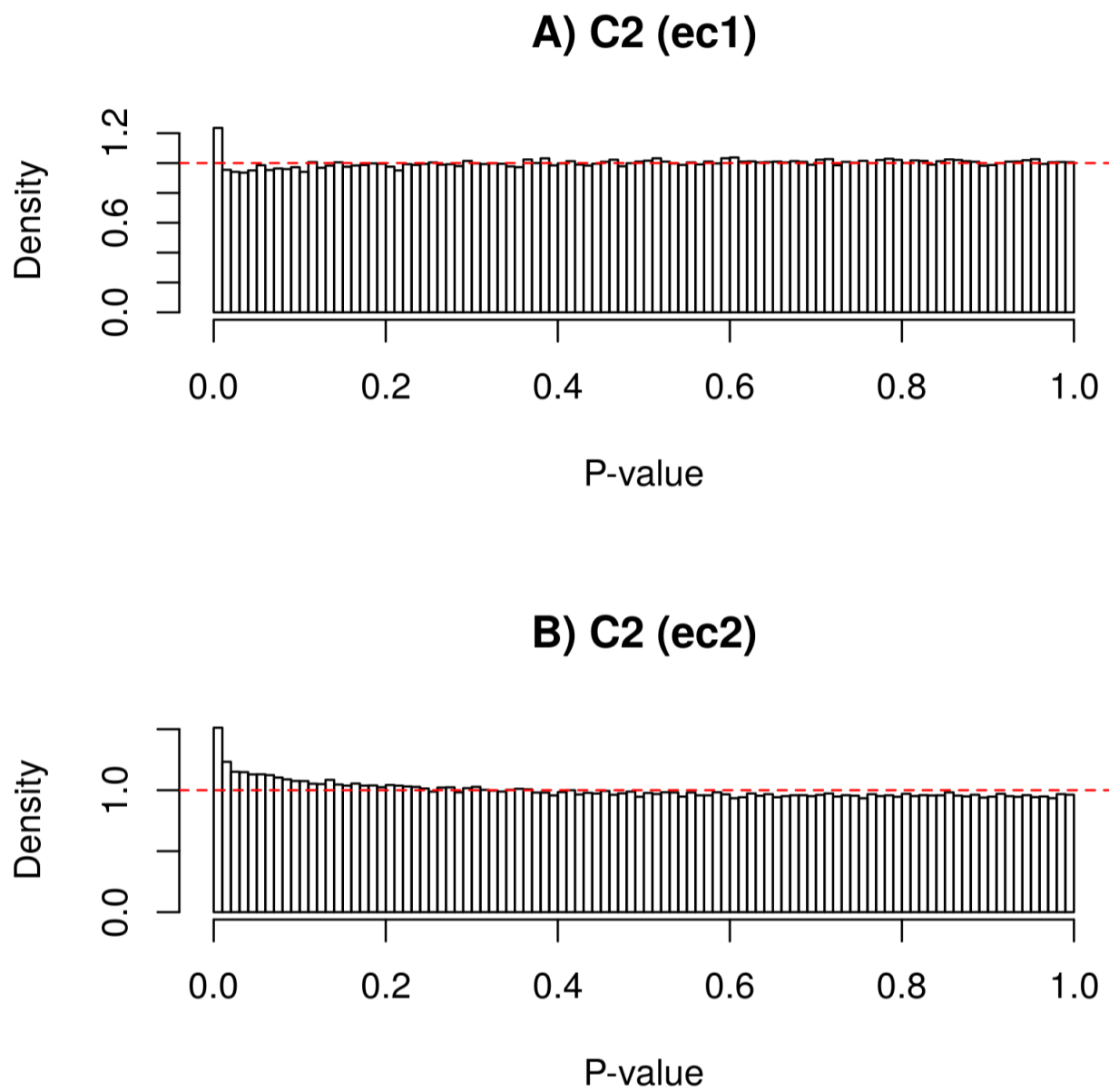


FIG. S1. Distribution of the p-values computed on the simulated data ($n=500,000$ SNPs) and derived from the C_2 statistics for the environmental contrasts *ec1* A) and *ec1* B), assuming a χ^2 null distribution (with one degree of freedom). The red horizontal dashed line represents the uniform distribution.

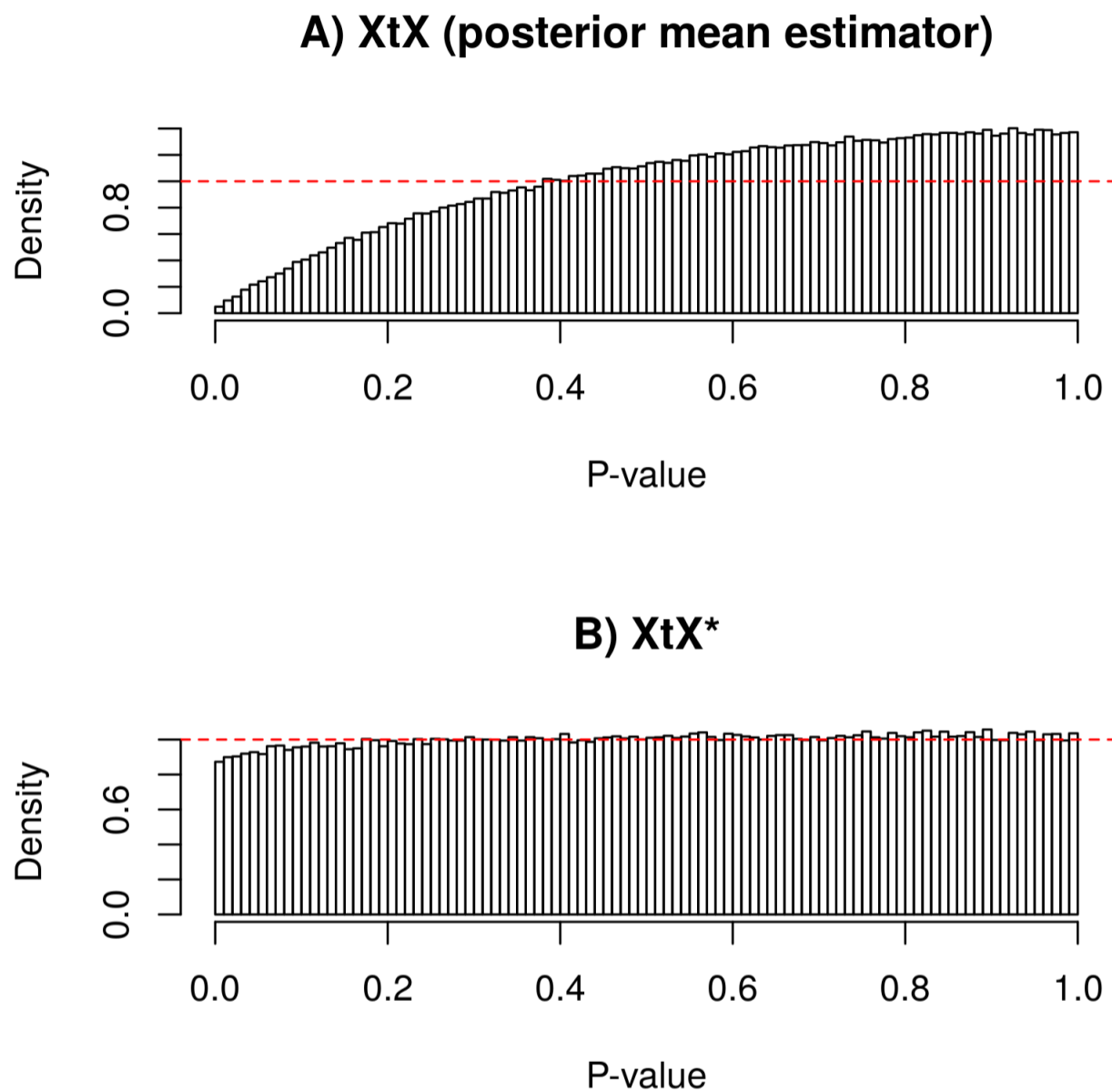


FIG. S2. Distribution of the p-values computed on the simulated data sets ($n=500,000$ SNPs) and derived from the SNP differentiation estimator XtX (posterior mean estimator) A) and the new estimator \widehat{XtX}^* B), assuming a χ^2 null distribution (with $J=16$, the number of population, degrees of freedom). To account for the bilateral nature of the underlying test (SNPs might be over or under-differentiated if under directional or balancing selection), p-value were computed as $p=1-2|\Phi_{\chi^2(J)}(\widehat{XtX})-0.5|$, where $\Phi_{\chi^2(J)}$ represents the cumulative density function of the χ^2 distribution with J degrees of freedom. The red horizontal dashed line represents the uniform distribution.

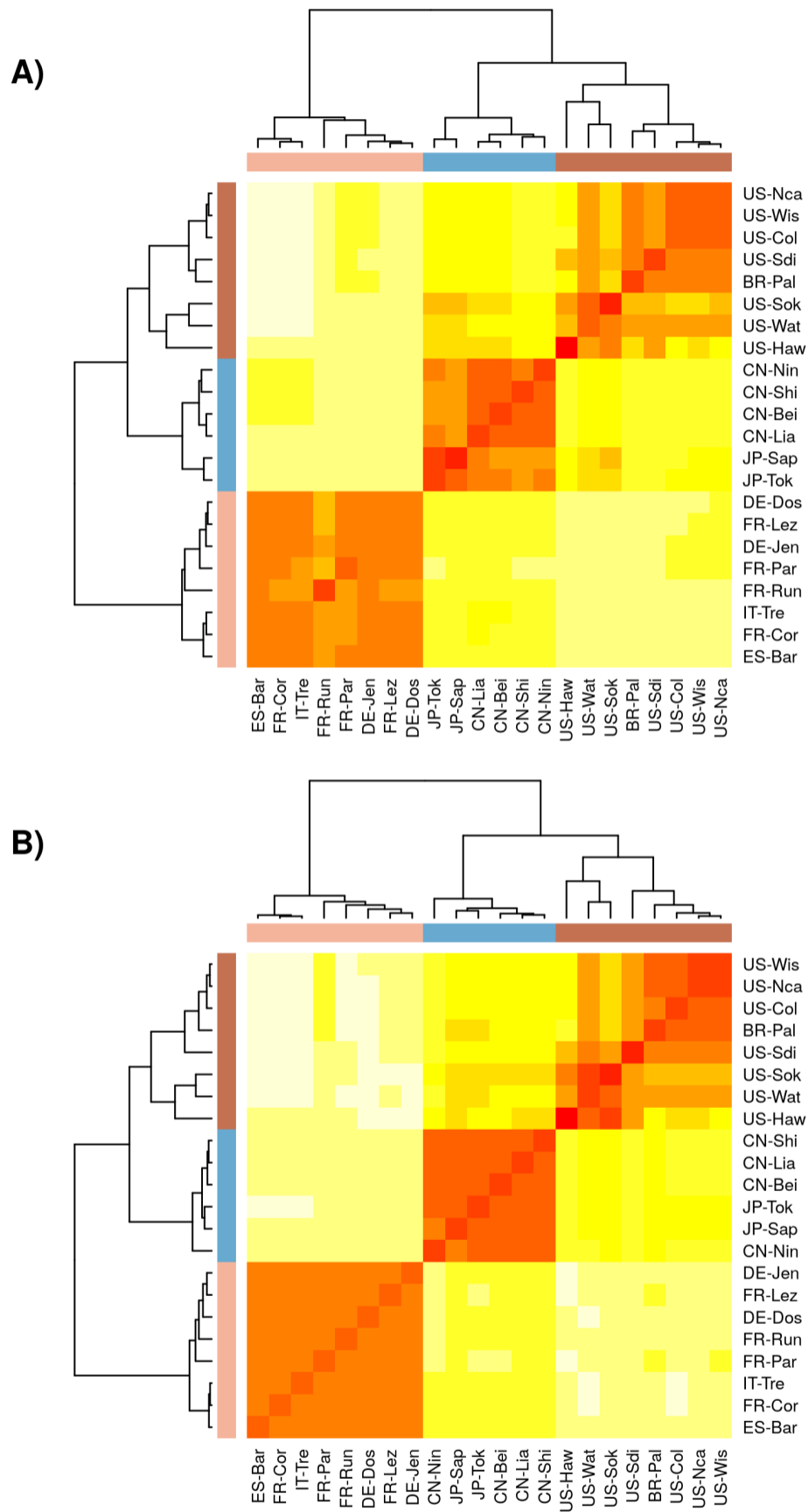


FIG. S3. Correlation plot representation of the scaled covariance matrices of population allele frequencies (Ω) among all 22 *D. suzukii* populations based on autosomal (A) and X-linked (B) SNPs.

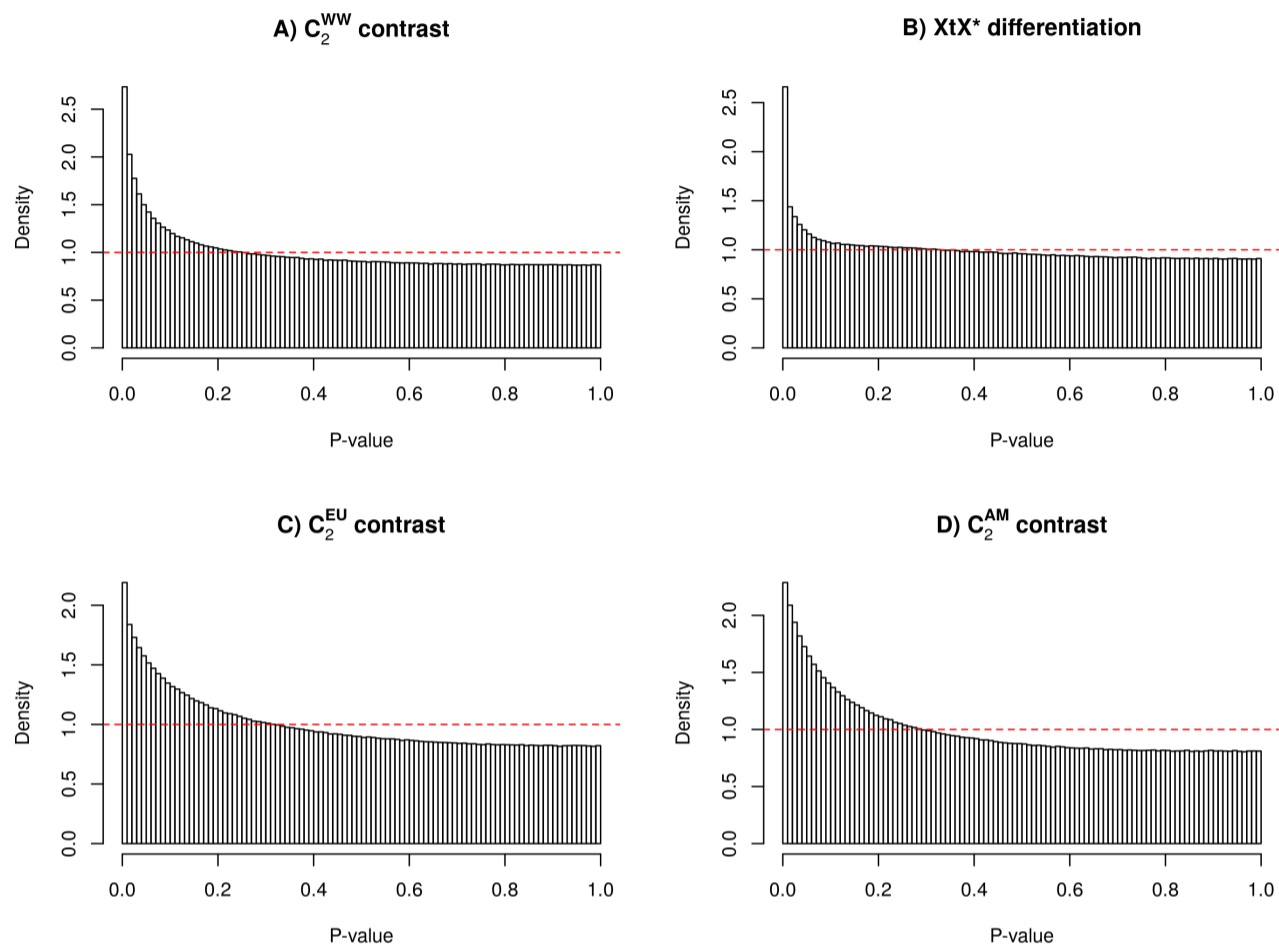


FIG. S4. Distribution of the p-values derived from the C_2 contrast statistics estimated from different analyses [A), C) and D)], and from the $\widehat{XtX^*}$ statistic for genetic differentiation among all 22 populations [B)]. The C_2 contrast statistics were estimated for 6 native *vs.* 16 worldwide invasive populations (C_2^{WW}) [A)]; 6 native *vs.* 8 invasive populations of the European invasion route (C_2^{EU}) [C)]; and 6 native *vs.* 8 invasive populations of the American invasion route (C_2^{AM}) [D)]. The distribution of the (two-sided) p-values derived from the XtX^* differentiation statistics (among all 22 *D. sukuzii* populations) is given in B). The red dashed line correspond to the uniform distribution.

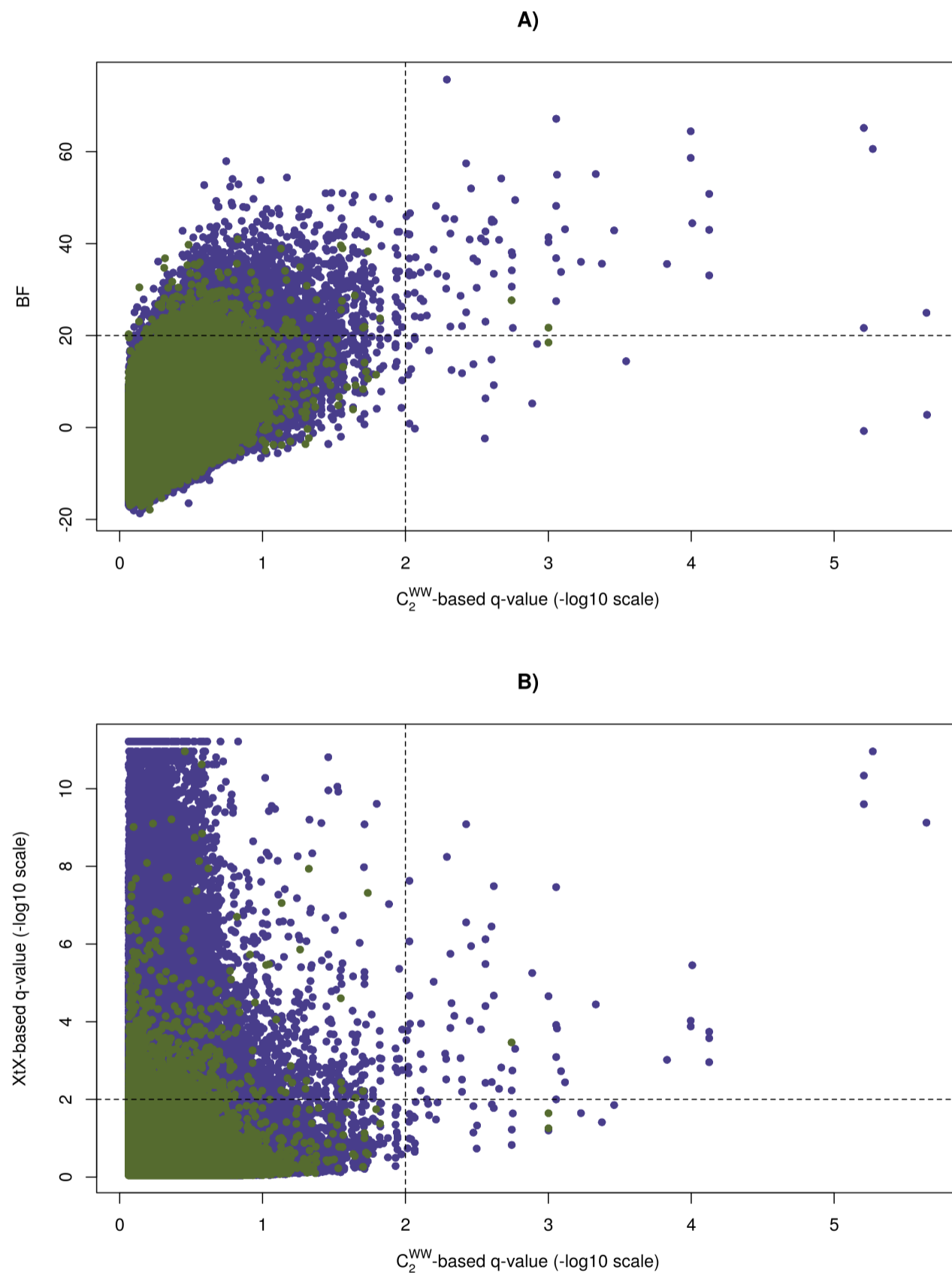


FIG. S5. Comparison of the C_2 statistics for the native vs. invasive status of the 22 *D. sukikii* populations (C_2^{WW}) with

BF for association A) and with XtX* overall differentiation estimates B). In A) the dashed horizontal line indicates the BF=20 db threshold of decisive evidence according to the Jeffreys' rule (Jeffreys, 1961) and the dashed vertical line to the 0.1% q-value threshold for the C_2 derived q-values. In B) the horizontal and vertical dashed lines indicate the 0.1% q-value threshold for the XtX* and C_2 derived q-values, respectively. Values for the autosomal (X-linked) SNPs are plotted in purple (green).

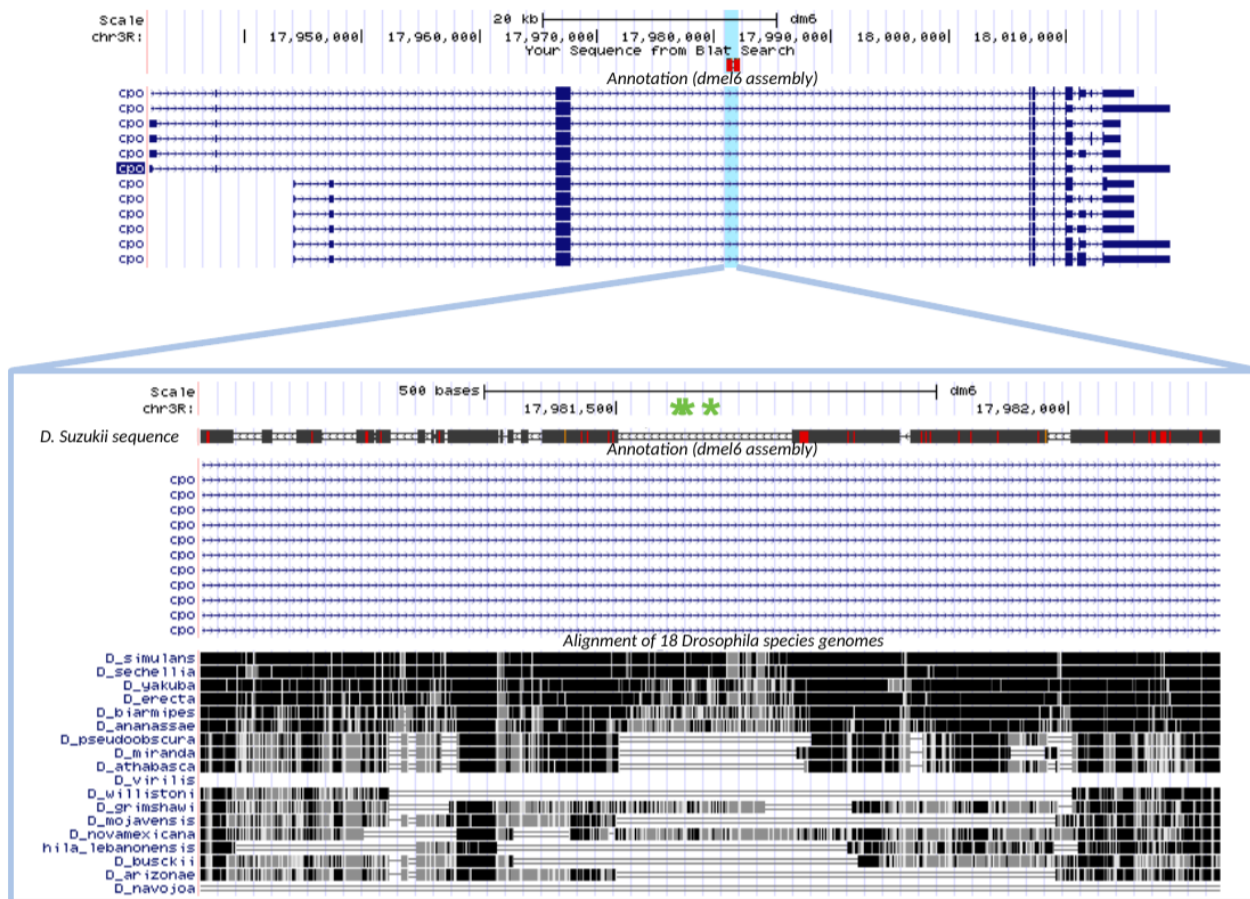


FIG. S6. Mapping of the three significant SNPs within the *cpo* gene onto the *dmel6* reference genome of *D. melanogaster* and alignment with genomes from other *drosophila* species. The aligned *D. suzukii* sequence consisted of a 1,193 bp sequence spanning the three significant SNPs (separated by 193 bp) indicated by a green star in the lower panel (the two first SNPs being those significant for the three contrast analyses). The plots were generated using the UCSC genome browser (<https://genome.ucsc.edu/>).