1   **Title (100 char):**
2   Global invasion history of the world's most abundant pest butterfly: a citizen science population genomics study
3
4   **Running title:**
5   Global invasion history of *Pieris rapae*
6

7   **Authors:**
8   Sean F. Ryan,[1,2]
9   Eric Lombaert[3],
10  Anne Espeset[4],
11  Roger Vila[5],
12  Gerard Talavera[5,6],
13  Vlad Dincă[7],
14  Mark A. Renshaw[8],
15  Matthew W. Eng[9],
16  Meredith M. Doellman[9],
17  Emily A. Hornett[10],
18  Yiyuan Li[9],
19  Michael E. Pfrender[9,11],
20  DeWayne Shoemaker[1]
21

22  **Affiliations:**
23  [1]Department of Entomology and Plant Pathology, University of Tennessee, Knoxville, TN 37996 USA
24  [2]Department of Applied Ecology, North Carolina State University, NC 27695 USA
25  [3]INRA, Université Côte d'Azur, CNRS, ISA, Sophia-Antipolis, France[4]Program in Ecology, Evolution, and
26  Conservation Biology
27  [4]Department of Biology, University of Nevada, Reno, NV 89557 USA
28  [5]Institut de Biologia Evolutiva (CSIC-UPF), Passeig Marítim de la Barceloneta 37, Barcelona 08003, Spain
29  [6]Department of Organismic and Evolutionary Biology and Museum of Comparative Zoology, Harvard University,
30  26 Oxford Street, Cambridge, MA 02138, USA
31  [7]Department of Ecology and Genetics, PO Box 3000, 90014 University of Oulu, Finland
32  [8]Oceanic Institute of Hawai'i Pacific University, Waimanalo, HI 96795
33  [9] Department of Biological Sciences, University of Notre Dame, South Bend, IN 46556
34  [10] Department of Evolution, Ecology and Behaviour, University of Liverpool, Liverpool L69 3BX, United Kingdom
35  [11] Environmental Change Initiative, University of Notre Dame, South Bend, IN 46556
36

37  **Corresponding author:**
38  Sean F. Ryan,
39  Email: citscisean@gmail.com
40

41  **Keywords:**
42  Invasive Species, Invasion History, Genomics, Agricultural Pest, Citizen Science, Approximate Bayesian
43  Computation

44    **Abstract**

45          A major goal of invasion and climate change biology research is to understand the

46    ecological and evolutionary responses of organisms to anthropogenic disturbance, especially

47    over large spatial and temporal scales. One significant, and sometimes unattainable, challenge of

48    these studies is garnering sufficient numbers of relevant specimens, especially for species spread

49    across multiple continents. We developed a citizen science project, "Pieris Project", to

50    successfully amass thousands of specimens of the invasive agricultural pest *Pieris rapae*, the

51    small cabbage white butterfly, from 32 countries worldwide. We then generated and analyzed

52    genomic (ddRAD) and mitochondrial DNA sequence data for these samples to reconstruct and

53    compare different global invasion history scenarios. Our results bolster historical accounts of the

54    global spread and timing of *P. rapae* introductions. The spread of *P. rapae* over the last ~160

55    years followed a linear series of at least four founding events, with each introduced population

56    serving as the source for the next. We provide the first molecular evidence supporting the

57    hypothesis that the ongoing divergence of the European and Asian subspecies of *P. rapae*

58    (~1,200 yrBP) coincides with the domestication of brassicaceous crops. Finally, the international

59    success of the Pieris Project allowed us to nearly double the geographic scope of our sampling

60    (i.e., add >1,000 specimens from 13 countries), demonstrating the power of the public to aid

61    scientists in collections-based research addressing important questions in ecology and

62    evolutionary biology.

63

64    **Non-technical summary:** We provide genetic evidence that the success of the small cabbage

65    white butterfly—its rise to one of the most widespread and abundant butterflies on the planet—

66    was largely facilitated by human activities, through the domestication of its food plants and the

67    accidental movement of the butterfly by means of trade and human movement (migration).

68    Through an international citizen science project—Pieris Project—people from around the world

69    helped to unravel the global invasion history of this agricultural pest butterfly by collecting

70    samples for DNA analysis. The success of this citizen science project demonstrates the power of

71    the public to aid in collections-based research that address important questions related to ecology

72    and evolutionary biology.

73

## Introduction

Invasive species—species spread to places beyond their natural range, where they generate a negative impact (e.g., extirpate or displace native fauna, spread disease, destroy agricultural crops[1])—continue to increase in number, with no signs of saturation[2]. The spread of invasive species often is driven by (human) migration, global trade and transportation networks[3], and, in some cases, domestication of wild plants and animals[4]. A critical, often first step to mitigating the spread and impacts of invasive species is to understand their invasion history, including assessing source populations, routes of spread, number of independent invasions, and the effects of genetic bottlenecks, among other factors. Such detailed knowledge is crucial from an applied perspective (e.g., developing an effective biological control program) as well as for addressing basic questions associated with the invasion process (e.g., genetic changes and adaptation to novel environments)[5].

Unraveling a species' invasion history often requires sampling across large spatial and temporal scales, which can be challenging and expensive, particularly for many invasive species found on multiple continents. Citizen science—research in which non-scientists play a role in project development, data collection or discovery and is subject to the same system of peer review as conventional science[6]—is a potentially powerful means to overcome some of these challenges. A major strength of citizen science is that it can greatly enhance the scale and scope of science and its impact on society[7]. Consequently, there are now thousands of citizen science projects worldwide. Yet, still very few involve agricultural pests[6] and nearly all rely on observations (e.g., sightings or photographs), limiting their capacity to address fundamental questions in ecology and evolution.

97    *Pieris rapae*, the small cabbage white butterfly, is the world's most widespread and

98    abundant pest butterfly. Caterpillars of this species are a serious agricultural pest of crops in the

99    Brassicaceae family (e.g., cabbage, kale, broccoli, brussels sprouts)[8]. This butterfly is believed to

100   have originated in Europe and subsequently undergone a range expansion into Asia several

101   thousand years ago as a result of domestication and trade of its host plants[9,10]. The Europe and

102   Asia populations recognized today are believed to represent separate subspecies—*P. rapae rapae*

103   and *P. rapae crucivora*, respectively.

104       The small cabbage white butterfly has been introduced to many other parts of the world

105   over the last ~160 years. These invasions are unique in that there is a wealth of historical records

106   (observations and collections) documenting the putative dates of first introduction (North

107   America in the 1860s[11], New Zealand in 1930[12], and Australia in 1937[13]). Detailed accounts and

108   observations from what was essentially a 19th century citizen science project led by the

109   entomologist Samuel Scudder provide a chronology of the spread of *P. rapae* across North

110   America and suggest that there were multiple independent introductions[11]. While the small

111   cabbage white butterfly ranks as one of the most successful and abundant invasive species, a

112   detailed analysis of its invasion history has never been undertaken[9,14]. In addition, the

113   consequences of this rapid invasion on the population genetic structure and diversity also are

114   unknown.

115       Here, we employ a collection-based citizen science approach to obtain range-wide, long-

116   term, population-level sampling of this globally distributed invasive agricultural pest. Molecular

117   genomics tools are then applied to this global collection of specimens to demonstrate that a

118   citizen science approach can be used to address a multitude of important questions in ecology

119    and evolutionary biology, including the reconstruction of the global invasion history of *P. rapae*

120    and assessment of historical and contemporary patterns of genetic structure and diversity.

121

## 122    Results

### 123    *Citizen scientist assisted sampling*

124            The international citizen science project—Pieris Project—recruited more than 100

125    participants that collected >1,000 butterflies from 13 countries in less than three years. The

126    majority of our participants (citizen scientists) were recruited through entomological and

127    lepidopterist societies and other organizations related to nature and science. These citizen

128    scientist collections were supplemented with collections from researchers bringing the total to

129    >3,000 *P. rapae* from the period of 2002-2017 (median collection year: 2014, Fig S1). Nearly

130    half (338/794) of the specimens used to generate mtDNA or ddRADseq data were from citizen

131    scientists. Of the 32 countries represented in our collection, five countries (Portugal, Czech

132    Republic, Gibraltar, Turkey, and South Korea) were made up of specimens sampled entirely by

133    citizen scientists, three countries (Russia, Australia, and New Zealand) had the majority of

134    specimens coming from citizen scientists, and three countries (United States, Canada, and Spain)

135    had nearly half the specimens coming from citizen scientists (Table S1). These samples

136    collectively cover nearly the entire native and invaded ranges, consisting of 293 localities

137    spanning 32 countries (Fig 1, up-to-date collections map); note, we do not have collections from

138    South America because there are no (known) populations of *Pieris rapae* on the continent.

139            A total of 22,059 autosomal (ddRADseq) Single Nucleotide Polymorphisms (SNPs) for

140    559 individuals (average depth: $74X \pm 28$ sd; average missingness: $2.9\% \pm 4.3$ sd) passed quality

141    filtering (Fig 1a). We also sequenced a 502 bp region of the mitochondrial gene cytochrome c

142   oxidase subunit 1 (*COI*) from 751 individuals (632 of these individuals were also used to

143   generate ddRADseq data) and supplemented these sequences with 251 additional sequences from

144   various online databases (total individuals with *COI* sequence = 1,002; Fig 1b).

145

146   *Global patterns of autosomal genetic differentiation and diversity*

147   We filtered the ddRADseq data for autosomal markers and found evidence for at least

148   seven genetically distinct clusters (ADMIXTURE lowest cross-validation error: 0.25 for K = 7)

149   (Fig 2a). These genetic clusters largely correspond to the continental regions sampled and we

150   refer to them henceforth as populations, named based on their sampling region: Europe, North

151   Africa, Asia, Russia (east), North America (east), North America (west), Australia/New Zealand

152   (Fig 2e). The greatest genetic differentiation was between Asia (including Russia (east)) and all

153   other populations; average $F_{ST}$ = 0.26±0.03sd (Fig 2c). Visual inspection of ancestry assignments

154   (at higher values of K) suggests additional hierarchical levels of structure, primarily in Asia, but

155   also within North America, and between Australia and New Zealand (Fig S2a,b). Surprisingly,

156   we were unable to detect (geographically coherent) structure within Europe (except for Malta

157   being distinct from the rest of Europe) or within Australia.

158   Almost all recently introduced populations (i.e., North America, Australia and New

159   Zealand) exhibit lower observed heterozygosity and nucleotide diversity compared with

160   populations in the native range (i.e., Europe and Asia), consistent with population bottlenecks

161   associated with these introductions (Fig 3). North America (east) was a notable exception among

162   the introduced populations, with observed heterozygosity higher than populations found in the

163   native range. All estimates of Tajima's *D* fell within the range of -1 to 1, suggesting most

164   populations are near equilibrium. However, there is a negative relationship between estimates of

165    Tajima's *D* and time since introduction—i.e., more recent introductions have higher (positive)

166    estimates of Tajima's *D*, suggesting that North America (west), New Zealand and Australia are

167    still recovering from repeated population bottlenecks.

168

169    ***Global Invasion History***

170            We compared a number of alternative invasion history scenarios for both the native and

171    introduced populations using ddRADseq autosomal data within an approximate Bayesian

172    computation random forest (ABC-RF) framework. We used an iterative process for selecting

173    each bifurcation event, starting with native populations (Europe and Asia), then Russia (east) and

174    North Africa, followed by the recently introduced populations (North America (east and west)),

175    New Zealand and Australia, and then simulated a full model that incorporated all the best

176    supported scenarios to get final parameter estimates (Fig 2e, Table 1). Based on this full final

177    scenario (Fig 2d), posterior model checking revealed that the observed values of only six

178    summary statistics out of 928 (i.e. 0.6%) fall in the tail of the probability distribution of statistics

179    calculated from the posterior simulation (i.e. $p < 0.05$ or $p > 0.95$), which indicates that the

180    chosen model fitted well the observed genetic data. From parameter estimation (Table 1), we

181    found the greatest support for a scenario with an ancestral population undergoing a demographic

182    expansion *ca.* 20,000 (32,000–4,900) yrBP (Years Before Present) (Table S2). In evaluating the

183    source for the Europe and Asia populations, we found the strongest support for the scenario of an

184    ancestral population giving rise to both the Europe and Asia populations (~85% posterior

185    probability), *ca.* 1,200 (2,900–300) yrBP, over scenarios with Europe as the source for Asia, or

186    Asia as the source for Europe (Fig S3a; Table 1). We evaluated multiple scenarios to determine

187    the source for the Russia (east) and North Africa populations and found the strongest support for

188   a scenario with Asia giving rise to the Russia (east) *ca.* 300 (800–200) yrBP, and the Europe

189   population giving rise to the North Africa population ca. 200 (600–200) yrBP (Fig S3b; Table 1).

190       We found strong support (total of 996 random forest votes out of 1000) for Europe being

191   the source of introduction to North America (east) (Fig S3c; Table 1). The scenario of a single

192   introduction had only slightly better support than the scenario with multiple (two) introductions,

193   and both have a similar number of random forest votes (576 and 418, respectively, out of 1,000,

194   for dataset 1). Thus, we cannot clearly distinguish between these two scenarios, and prior error

195   rate was consequently relatively high (~33%). However, subsequent analyses performed by

196   considering multiple introductions for the formation of North America (east) does not

197   qualitatively change any of the following results (results not shown).

198       We found the strongest support for North America (east) serving as the source for the

199   genetically distinct North America (west) population when compared to alternative scenarios

200   with Asia or Europe (for both the scenario with ~ 400 or 200 yrBP prior estimate for date of

201   introduction) as the source (Fig S3d; Table 1). This introduction was estimated to have occurred

202   *ca.* 137 yrBP. For the introduction into New Zealand, we found strong support for North

203   America (west) being the source, when compared to Europe, Asia, or North America (east) as the

204   source (Fig S3e; Table 1). The New Zealand population was found to have the greatest support

205   as being the source for the introduction to Australia (Fig S3f; Table 1). All of these results were

206   obtained with dataset 1 but were qualitatively confirmed by the analyses of datasets 2 and 3

207   (Table 1).

208       Demographic parameter estimates from ABC-RF analyses suggest each introduced

209   population underwent a severe bottleneck, but the intensity (duration and number of founders

210   with respect to the effective size of the source population) differed among populations (Table

211    S2). Specifically, New Zealand and, more importantly, North America (west) were estimated to

212    have undergone the most intense bottlenecks, whereas North America (east) and, to a lesser

213    extent, Australia suffered less intense bottlenecks.

214

215    *Global patterns of mtDNA haplotype diversity and distribution*

216         A total of 88 COI haplotypes were identified from 1,002 individuals, and 85% of these

217    individuals harbored one of the eleven most common haplotypes (Fig 4; Fig S4, Table S3). The

218    geographic distribution of mtDNA haplotypes is consistent with the invasion routes identified

219    from ABC-RF analyses—haplotypes in introduced populations are largely a subset of those from

220    putative source populations or differ by only one to two mutations from haplotypes in high

221    frequency in the putative source populations (Fig 4a; see interactive figure to plot haplotypes

222    individually— http://www.pierisproject.org/ResultsInvasionHistory.html).

223         Estimates of mtDNA haplotype diversity (richness) were highest in Asia and Europe and

224    had large confidence intervals (based on rarefaction curve analysis), indicating these populations

225    were likely undersampled (Fig 4c). All introduced populations had substantially lower estimates

226    of mtDNA haplotype diversity. New Zealand, Australia, North America (west), North Africa,

227    and Russia (east) were estimated to have less than a dozen mtDNA haplotypes, whereas North

228    America (east) had substantially (~3X) more mtDNA haplotypes and was significantly greater

229    than all other introduced populations (based on non-overlapping confidence intervals). Estimates

230    of mtDNA nucleotide diversity are similar to those observed for autosomal markers using

231    ddRADseq data, with the exception of Australia, which had higher nucleotide diversity than New

232    Zealand and North America (west) (Fig 4d). From a global perspective, there appears to be a

233    general trend of decreasing nucleotide diversity with increasing distance from southern Europe

234    and the eastern Mediterranean region (Fig S5).

235        Considering the mtDNA haplotypes found in North America and their frequencies in sub-

236    populations of Europe, we estimate that the minimum number of individuals that would need to

237    have been sampled from the sub-population of England (2.3) is $23 \pm 12$ sd individuals or $123 \pm$

238    88 sd individuals for Spain/southern France (2.4) to account for all haplotypes in North America.

239

240    **Discussion**

241    *Citizen science greatly expands range-wide sampling of butterflies*

242        We show for the first time how the public can contribute to our understanding of species

243    invasions through a collection-based citizen science project. Collections made by citizen

244    scientists dramatically expanded the geographic scope of our study by nearly doubling the

245    number of countries we were able to include in our analyses. Samples from many of these

246    countries either solely or primarily came from citizen scientists, including nearly the entire

247    region of Australasia. Moreover, these contributions allowed us to increase substantially the total

248    number of individuals in each of the populations studied, with nearly half of all specimens

249    sequenced in our study coming from citizen scientists. We estimate that the use of citizen science

250    to aid in the collection of *P. rapae* from across its near-global range resulted in tens of thousands

251    of dollars (USD) in cost-savings that would have been required to cover salary and travel costs.

252    We believe our citizen science approach can be applied to other systems, particularly to

253    organisms that are easily identifiable (e.g., spotted lantern fly or Giant African snail) and easy to

254    transport (e.g., dead invertebrates), to address questions in invasion biology as well as a broad

255    range of questions in ecology and evolutionary biology. As examples, we currently are

256    leveraging our large collection to address questions concerning the effects of climate and land

257    use changes on wing pigmentation of this butterfly and to identify genomic regions underlying

258    ecological selection.

259         The development, implementation, and maintenance of this project was not trivial, as is

260    the case with many citizen science projects, and required considerable time and effort engaging

261    the public (e.g., contacting organizations, using social media, responding to emails) and

262    processing samples. Collections-based citizen science projects that focus on less charismatic

263    species or incorporate non-lethal forms of sampling (e.g., eDNA) and are easy to collect (slow

264    moving or sessile) may have the greatest success. We suggest that those interested in applying a

265    collections-based citizen science approach seek advice from, or build collaborations with,

266    individuals with experience in the field of citizen science.

267

268    ***Geographic spread and divergence of <u>P. rapae</u> driven by host plant domestication***

269         The Europe and Asia populations of *P. rapae* are believed to represent distinct

270    subspecies—*P. rapae rapae* and *P. rapae crucivora*, respectively—based on phenotypic

271    differences[9] and evidence for reproductive isolation[15]. Our study provides further support for

272    this, revealing the two main genetic lineages recovered by ddRADseq for *P. rapae* worldwide

273    correspond to the Europe (*Pieris rapae rapae*) and Asia (*Pieris rapae crucivora*) populations. It

274    has long been hypothesized that the domestication of brassicaceous crops, which are the primary

275    host plants for this butterfly, aided the spread and/or divergence of the Europe and Asia

276    subspecies[10]. Our data support this hypothesis. We estimate the divergence of both the Europe

277    and Asia populations (subspecies) occurring within the last ~3000 years. This estimate overlaps

278    with estimates for the domestication of brassicaceous crops that occurred during this time period

279    (i.e., *Brassica oleracea* and *Brassica rapa*)[16–18]. However, our results do not support the

280    hypothesis that Europe was the source of the Asia population[9]. Instead, our results suggest both

281    the Europe and Asia populations independently diverged from an ancestral population, most

282    likely *ca.* 1,200 yrBP. This time period corresponds with the intensification in the cultivation of

283    *B. oleracea* varieties, such as cabbage and brussels sprouts[19].

284        It remains unclear whether *P. rapae* spread across and occupied Europe and Asia during

285    this expansion event, and then diverged *in situ* in response to the domestication of brassicaceous

286    crops, or whether *P. rapae* was more restricted in distribution (e.g., confined to Europe or the

287    eastern Mediterranean region) and diverged in association with the domestication of

288    brassicaceous crops across Eurasia. Consistent with the hypothesis that Europe and Asia *P.*

289    *rapae* populations diverged as they spread out of the eastern Mediterranean, the range boundaries

290    of the Europe and Asia populations (subspecies) abut in the eastern Mediterranean region and

291    genetic diversity generally decreases with increasing distance from this region. In further support

292    of this hypothesis, there is growing evidence that the domestication of *Brassica oleracea* and

293    *Brassica rapa* originated in the Mediterranean region[18]. However, additional sampling in western

294    Asia is needed to further evaluate this hypothesis.

295        Interestingly, the putative ancestral population that gave rise to the Europe and Asia

296    populations appears to have undergone a rapid increase in effective population size *ca.* 7,000–

297    28,000 yrBP. This time period overlaps with early human development of agriculture. However,

298    our median estimate for the date of this expansion was *ca.* 20,000 yrBP, placing it at the end of

299    the last glacial maximum *ca.* 23,000–19,000 yrBP[20]. Changes in the distribution and demography

300    of species in response to glacial–interglacial cycles is well documented[21,22], and may be more

301    likely to have facilitated a major demographic shift in *P. rapae*, as the earliest domestication of

302    brassicaceous crops was relatively recent (earliest evidence being *ca.* 7,000 yrBP)[17].

303

304    ***Recent invasion history largely reflects historical records, but with a few unexpected findings***

305            Although historical records of species invasions can be misleading[23], our molecular

306    genomics-based reconstruction of the *P. rapae* global invasion history is largely consistent with

307    that historically documented through observations. As expected, we found Europe to be the most

308    likely source of this butterfly's introduction into North America. However, we unexpectedly

309    found that there was no discernable nuclear genetic structure in Europe (even when K = 30),

310    making it impossible to narrow down with confidence the source population to a specific locality

311    or country (e.g., England vs Spain). However, mtDNA haplotype distributions and frequencies in

312    European countries suggest England as the most likely source—i.e., fewer individuals would be

313    required to produce the mtDNA haplotypes found in North America, if England was the source

314    rather than Spain and southern France. We do not know what specific factors account for the

315    lack of genetic structure in Europe. One possibility is long-distance dispersal of this species[24]

316    coupled with historic and/or ongoing human assisted dispersal has led to high levels of gene

317    flow. Another interesting possibility supported by some evidence[24–26] is that *P. rapae* is

318    migratory or undergoes migratory-like events.

319            Historical records obtained by Samuel Scudder pointed to possible multiple introductions

320    of *P. rapae* into North America, which occurred during or shortly after its initial invasion[11].

321    Confirming multiple introductions from the same source population early in an invasion,

322    particularly one that quickly underwent a rapid expansion, is extremely difficult. The best-fit

323    model to our data suggests a scenario with a single introduction, but it seems reasonable there

324    were multiple introductions for a couple of reasons. First, both competing scenarios—one vs

325    multiple introductions from Europe—had a similar number of random forest votes (576 vs 418

326    out of 1,000) and the selected scenario (i.e. one introduction from Europe) had a low posterior

327    probability estimate (~0.5; in contrast, all scenarios that were chosen in the other analyses had

328    posterior probability estimates >0.70). Second, we estimated a rather low bottleneck intensity,

329    with a founding population size of ~50-100 individuals; this estimate is much higher than a

330    previous estimate of one to four individuals[27]. It seems unlikely that North America was founded

331    from a single introduction given this rather large estimated founding population size and the

332    reasonable assumption that no more than a few dozen (unrelated) butterflies would be

333    transported on any one ship. Third, multiple introduction events (from Europe) would help

334    explain the higher heterozygosity found in North America (east) than in the native range—i.e.,

335    multiple introductions aided in the rebound in genetic diversity as the introduced population

336    spread across eastern North America.

337         Our second unexpected finding was the identification of a genetically distinct population

338    within North America that is restricted to the western USA. We found evidence of admixture in

339    areas where the two North America (east and west) populations come into contact, suggesting

340    that these genetically distinct populations are neither geographically or reproductively isolated

341    from each other. The geographic extent of this admixture zone is not clear from our sampling,

342    nor are the consequences of gene flow between these populations. We initially hypothesized this

343    western population represents an early introduction brought by Spaniards during the 1600s, but

344    our data instead indicate that it most likely results from a secondary founder event from the

345    North America (east) population brought during the ~1860-1880s as a result of the rapid

346    development of railroad lines[28], specifically from the eastern USA to central California (Video

347    S1).

348        Our results confirm previous speculation that North America (west), likely San

349    Francisco, California, was the source of introduction to the Hawaii islands, based on individuals

350    from Hawaii being assigned to the North America (west) cluster but not being reported in the

351    Hawaiian islands until 1987[29] (after the arrival of *P. rapae* to central California). Also,

352    unexpectedly, our results suggest that the introduction of *P. rapae* to New Zealand came from

353    North America (west), and not from Europe, as was believed to be the most likely case given the

354    United Kingdom was the largest exporter into New Zealand at the time[30]. Lastly, previous

355    speculation that New Zealand is the immediate source of *P. rapae* in Australia[13] is supported by

356    our data.

357

358    ***Implications for invasion biology***

359        Our study also sheds light more broadly on invasion biology. Growing evidence shows

360    many invasive populations are able to flourish and adapt to new environments despite substantial

361    loss of genetic diversity—a phenomenon termed the genetic paradox of invasions[31]. *Pieris rapae*

362    is a remarkable example of this paradox. This butterfly has rapidly expanded its range following

363    each new founding event, despite repeated population bottlenecks—at least four separate times,

364    with each new founding population the product of a previously bottlenecked population (i.e.,

365    multiple serial founding events). Whether introduced populations maintained high genetic

366    variation in ecologically relevant traits following each founding event remains unclear. Evidence

367    of local adaption for thermal tolerances among populations in North America[32] suggests such

368    variation exists. However, resolving this paradox and the persistent puzzle of how this butterfly

369    has been an extremely successful invader into new environments will require future studies to

370    assess the relative contributions of factors such as adaptative evolution, phenotypic plasticity,

371    natural enemy escape, and domestication of its host plants.

372

373    **Methods**

374    *Specimens collection and DNA extractions*

375        The *Pieris rapae* specimens were collected as part of an international citizen science

376    project—Pieris Project (pierisproject.org)—that was launched in June 2014 and through

377    collections by researchers. A website was created in 2014 for Pieris Project that included a

378    description of the research goals and collection protocol—specimens were to be individually

379    placed in hand-made or glassine envelopes, labeled with location and date collected, and placed

380    in a freezer overnight, then air-dried for at least two days and shipped using standard mail. The

381    project was advertised through social media—Twitter (@PierisProject) and Facebook

382    (https://www.facebook.com/pierisproject/), and through listservs, social media, and blogs of

383    Entomological and Lepidopterists societies and nature/science/citizen science related

384    organizations (e.g., YourWildLife, eButterfly, National Geographic, and SciStarter). Once

385    received, specimens were stored in 95% ethanol and kept at -20 °C; depending on the collector,

386    specimens were air-dried for a few days to years prior to being placed in ethanol. Genomic DNA

387    was isolated from tissue from prothorax or (2-3) legs using DNeasy Blood and Tissue Kit spin-

388    columns (Qiagen, Hilden; Cat No./ID: 6950).

389        To estimate the contributions by scientists, we binned the collector of each specimen into

390    one of two categories: 1) researcher, and 2) citizen scientist. Collectors whose identity was

391    known (>90% of participants) were not considered as citizen scientists if they had a college

392    degree in biology. This makes our estimated contribution by citizen scientists a conservative

393    estimate, as some of these participants may consider themselves citizen scientists. There is a

394    great deal of debate as to what does or does not constitute being a citizen scientist. Our threshold

395    is based on our previously stated definition of citizen science that is accepted by many within the

396    field of citizen science.

397

398    *ddRADseq sequencing and filtering*

399          Nine reduced-complexity libraries were generated using a double-digest restriction-site-

400    associated DNA fragment procedure (ddRAD)[33] following Ryan et al., 2018[34]. Briefly, genomic

401    DNA (~400 ng) was digested with the restriction enzymes EcoR1 and Mse1 and a universal

402    Mse1 and barcoded EcoR1 adapter ligated to the digested DNA. Ligated products were diluted

403    10 times with 0.1X TE buffer prior to PCR enrichment. Amplified products with unique

404    barcodes were pooled into a single mixture prior to purification. The library was purified three

405    times with 0.8X volume of Agencourt Ampure XP beads (Beckman Coulter, A63881). At the

406    end of each round of purification, elution volume was reduced to 0.25 – 0.5X of the beginning

407    sample volume. After three rounds of purification, each library (1.0 µg) was size selected for 400

408    to 600 bp fragment length using 1.5% DF Cassette and BluePippin System (Sage Science).

409    Libraries were evaluated by Bioanalyzer 2100 system and sequenced across one lane Illumina

410    MiSeq (University of Notre Dame, Genomics Core Facility), 14 lanes of Illumina HiSeq 4000

411    (12 at University of Illinois and 2 at Beijing Genomics Institute); most samples were sequenced

412    on 2 (some 3) independent lanes.

413          Raw reads were demultiplexed and barcodes/cutsite removed using a custom Python

414    script. Reads were further trimmed and cleaned with the program Trimmomatic (v0.32)[35] using

415    default settings. The first 5 bp and any after 80 bps were then trimmed from all reads and only

416    reads at least 76 bp in length were retained, resulting in all reads being exactly 76 bp.

417        Reads were then aligned to the *Pieris rapae* genome v1 [27] using BWA-aln (v0.7.15)[36].

418    Variant calling was performed using GATK's Haplotypecaller (v3.8)[37,38] with the default

419    settings. Filters were applied in the following order: kept only biallelic SNPs, applied GATK's

420    "hard filtering" (QD < 2.0 || MQ < 40.0 || MQRankSum < -12.5 || ReadPosRankSum < -8.0),

421    SNPs with a genotype quality (GQ) < 20 were converted to missing data, removed SNPs with

422    minor allele frequency less than 0.01, kept SNPs with min of 1X coverage for 50% of

423    individuals, removed SNPs with coverage > 95th percentile (112.8 X coverage), removed

424    individuals with > 75% missing data, kept only SNPs with a minimum of 10X coverage in 90%

425    of individuals, and removed individuals with >25% missing data. Finally, SNPs with

426    heterozygosity > 0.6 were considered potential paralogs and were discarded.

427        As there is no linkage map for *P. rapae* and the genome is not assembled into

428    chromosomes, we applied a simple heterozygosity method to determine whether SNPs were

429    autosomal or sex (Z) linked. To do this we used the expectation that females should be

430    homozygous at all SNPs on the Z-chromosome—females are the heterogametic sex (ZW) in

431    Lepidoptera. Using 231 females, we calculated the percentage that were heterozygous or

432    homozygous at each site (SNP) using the is_het function from the R package vcfR (v1.6.0)[39] and

433    custom scripts in R. SNPs with greater than 25% missing data were removed. A scaffold (and all

434    SNPs within) was considered putatively Z-linked if > 60% of the SNPs fell below the threshold

435    of having less than 1% of the females being heterozygous (average number of SNPs for each

436    scaffold was 84 ± 101; mode = 4).

437    To complement the heterozygosity method, we also inferred the chromosome assignment

438    of each *P. rapae* scaffold using the approach by Ryan et al. (2017)[40]. Briefly, we used blastx

439    (ncbi-blast-2.2.30+)[41,42] to blast all peptide sequences within each scaffold of the *P. rapae*

440    assembly against the *Bombyx mori* genome (silkdb v2.0)[43]. The *B. mori* scaffold with the most

441    significant blast hits (based on Bit scores) was retained and used to determine the putative

442    chromosome of each *P. rapae* scaffold. All, except one scaffold, of those we identified as Z-

443    linked using the heterozygosity method, mapped to chromosome 1 (Z) and 2 (W) of *B. mori* (Fig

444    S6). That we found some regions of *P. rapae* mapping to the *B. mori* chromosome 2 (W)

445    suggests they are not completely syntenic—the *P. rapae* genome was assembled from males and

446    thus there should be no scaffold mapping to this chromosome. Some *P. rapae* scaffolds mapping

447    to chromosome 2 (W) of *B. mori* were recovered as actually being on the chromosome 1 (Z)

448    based on the heterozygosity method.

449    Using a subset of the putative Z-linked markers—SNPs where < 1% of females had a

450    heterozygous call (i.e., SNPs with a high likelihood of being Z-linked)—we validated the sex of

451    each individual. Females and males with >20% or <20% of these SNPs being heterozygous were

452    considered possibly mislabeled males and females respectively. These individuals were flagged,

453    and the specimens double-checked visually; in all cases, visual identification confirmed that

454    these individuals were mislabeled.

455

456    ***Inference of Population Structure and Diversity***

457    Population structure was investigated with the model-based clustering algorithm

458    ADMIXTURE[44] using default settings and a cross-validation = 10 for K 1-30 using a modified

459    SNP dataset, i.e., pruned for LD ($r^2 > 0.2$; calculated using VCFtools geno-r2) (17,917 SNPs).

460    The optimal K was that with the lowest cross validation error. The Bayesian program

461    fastSTRUCTURE[45] as well as a non-model-based multivariate approach—Discriminant Analysis

462    of Principal Components (DAPC; Fig S2c)[46]—were also used to confirm the results from

463    ADMIXTURE (see Supplementary Materials for more details) using the R package adegenet

464    (v2.1.1)[47]. Genetic assignments were plotted using custom scripts and the R package pophelper

465    (v2.2.3)[48]. A Neighbor-Joining tree based on genetic distance was constructed in the poppr

466    (v2.8.0)[49] and ape (v5.1)[50] packages in R, that included the species *Pieris napi*, *Pieris brassicae*,

467    and *Pieris canidia* as outgroups, using only sites with at least 15X coverage in 90% individuals

468    from this new dataset. Trees were visualized using FigTree (v1.4.4)[51]. Population differentiation

469    was estimated between all populations using the Weir & Cockerham's estimator of $F_{ST}$[52]

470    implemented in VCFtools (v0.1.15) using 10 kb windows and a window step size of 5 kb.

471         All measures of genetic diversity ($H_{obs}$, $\pi$, and Tajima's $D$) were calculated using SNPs

472    restricted to scaffolds longer than 100kb (22,059 SNPs). In an attempt to minimize the Wahlund

473    effect (i.e., reduction of heterozygosity caused by subpopulation structure), individuals were split

474    into spatially contiguous subpopulations from within the seven identified by ADMIXTURE

475    (N=34; one subpopulation from Mexico was not included because it contained only three

476    individuals); these were the same subpopulations used for the ABC-RF analyses. To control for

477    differences in sample size, we computed each statistic 1,000 times using a random subset

478    (without replacement) of seven individuals (size of smallest population). Heterozygosity was

479    estimated using the R package adegenet v2.1.1. Calculations for $\pi$, and Tajima's $D$ were

480    estimated using a dataset containing invariant sites (i.e., vcf files were created using gatk-4.0.4.0

481    with the flag -allSites true and the same filters as described above were then applied) with

482    VCFtools (v0.1.15) using 10 kb windows (and a window step size of 5 kb used for estimating $\pi$).

483

### *ABC-RF-based inferences of global invasion history*

485   An approximate Bayesian computation analysis (ABC)[53] was carried out to infer the

486   global invasion history of *Pieris rapae*. The eight populations considered in the ABC analysis

487   corresponded to the seven identified by ADMIXTURE, with an additional separation of New

488   Zealand and Australia for geographical reasons. Each population was represented in the analysis

489   by a single sub-population (individuals sampled within the same subregion and within a three-

490   year period) (dataset 1). ABC is a model-based Bayesian method allowing posterior probabilities

491   of historical scenarios to be computed, based on historical data and massive simulations of

492   genetic data. We used historical information (e.g., dates of first observation of invasive

493   populations) to define 6 sets of competing introduction scenarios that were analyzed sequentially

494   (Table 1 and Fig S3). Step by step, subsequent analyses used the results obtained from the

495   previous analyses, until the most recent invasive populations were considered. The first set of

496   competing scenarios (three scenarios) considered the evolutionary relationship between the

497   Asian and European populations. In the second analysis (four scenarios), we explored the links

498   between Asia, Europe, North Africa and Russia (east). In the third analysis (four scenarios), we

499   set North America (east) as the target and determined whether it originated from Asia or Europe,

500   either through one or two introductions. In the fourth analysis (5 scenarios), North America

501   (west) could be originating either from Europe, Asia or North America (east), and the

502   introduction could be ancient (400 yrBP) in the case of Europe. In the fifth analysis (four

503   scenarios), New Zealand could be originating either from Europe, Asia, North America (east) or

504   North America (west). Finally, the sixth analysis (five scenarios) aimed at deciphering the origin

505   of the Australian population by testing as source population New Zealand, North America

506    (west), and admixtures between New Zealand and either Europe, Asia or North America (west).

507    All scenarios of all analyses are detailed in Table 1 and Fig S3.

508    In our ABC analysis, historical and demographic parameter values for simulations were

509    drawn from prior distributions defined from historical data and demographic parameter values

510    available from empirical studies on *Pieris rapae*[11–13,29], as described in Table S5. Simulated and

511    observed datasets were summarized using the whole set of summary statistics proposed by

512    DIYABC[54] for SNP markers, describing genetic variation per population (e.g., mean gene

513    diversity across loci), per pair (e.g., mean across loci of $F_{ST}$ distances), or per triplet (e.g., mean

514    across loci of admixture estimates) of populations (see the DIYABC v2.1.0 for details about

515    statistics), plus the linear discriminant analysis axes[55] as additional summary statistics (Table

516    S6). The total number of summary statistics ranged from 18 to 388 depending on the analysis

517    (Table 1).

518    To compare the scenarios, we used a random forest process[56,57]. Random forest is a

519    machine-learning algorithm which uses hundreds of bootstrapped decision trees to perform

520    classification using a set of predictor variables, here the summary statistics. Some simulations

521    are not used in tree building at each bootstrap (i.e., the out-of-bag simulations) and can thus be

522    used to compute the "prior error rate", which provides a direct method for cross-validation. We

523    simulated a 10,000 SNPs datasets for each competing scenario using the standard Hudson's

524    algorithm for minor allele frequency (i.e., only polymorphic SNPs over the entire dataset are

525    considered), so the number of used markers ranged between 13,974 and 17,116 depending on the

526    analysis (Table 1). We then grew a classification forest of 1,000 trees based on the simulated

527    datasets. The random forest computation applied to the observed dataset provides a classification

528    vote (i.e., the number of times a model is selected among the 1,000 decision trees). The scenario

529    with the highest classification vote was selected as the most likely scenario, and we then

530    estimated its posterior probability by way of a second random forest procedure of 1,000 trees[57].

531    To evaluate the global performance of our ABC-RF scenario choice, we computed the prior error

532    rate based on the available out-of-bag simulations and conducted the complete scenario selection

533    analysis with two additional datasets with different sub-populations (dataset 2 and dataset 3)

534    representative of the same populations as dataset 1[58].

535    We then performed a posterior model checking analysis on a full final scenario including

536    all 8 populations (dataset 1), to determine whether this scenario matches well with the observed

537    genetic data. Briefly, if a model fits the observed data correctly, then data simulated under this

538    model with parameters drawn from their posterior distribution should be close to the observed

539    data. The lack of fit of the model to the data with respect to the posterior predictive distribution

540    can be measured by determining the frequency at which the observed summary statistics are

541    extreme with respect to the simulated summary statistics distribution (hence defining a tail-area

542    probability or $p$-value, for each summary statistic). We simulated 100,000 data sets under the full

543    final scenario (17,609 SNP and 928 summary statistics), and then obtained a 'posterior sample'

544    of 5,000 values of the posterior distributions of parameters through a rejection step based on

545    Euclidean distances and a linear regression post-treatment[53]. We simulated 1,000 new datasets

546    with parameter values drawn from this "posterior sample", and each observed summary statistic

547    was compared with the distribution of the 1,000 simulated test statistics, and its $p$-value,

548    corrected for multiple comparisons with the false discovery rate procedure[59], was computed.

549    Finally, 10,000 simulated datasets of the full final scenario were used to infer posterior

550    distribution values of all parameters, and some relevant composite parameters under a regression

551    by random forest methodology[60], with classification forests of 1,000 trees. The simulation steps,

552   the computation of summary statistics, as well as the model checking analysis were performed

553   using DIYABC v2.1.0. All scenario comparisons and parameter estimations were carried out in

554   R using the package abcrf (v1.7.1)[57].

555

556   *mtDNA sequencing and analysis*

557        A 1,600 bp region of COI was amplified using primers optimized to work with multiple

558   species within the genera *Pieris* (Pieridae_COI_F 5-AAATTTACAATYTATCGCTTA-3,

559   Pieridae_COI_R 5-TGGGGTTTAAATCCATTACATATW-3). When these primers failed we

560   amplified a 658 bp region of COI using previously published primers[61]. PCR amplicons were

561   purified using magnetic beads and amplified using standard fluorescent cycle sequencing PCR

562   reactions (ABI Prism Big Dye terminator chemistry, Applied Biosystems). Sequencing reactions

563   were purified using Agencourt CleanSeq magnetic beads (Beckman Coulter) and run on an ABI-

564   3730XL-96 capillary sequencer (Applied Biosystems) at the University of Florida biotechnology

565   facility (ICBR) or Macrogen (Macrogen Inc). Individuals with both forward and reverse reads

566   were assembled in Geneious 11.0.4 using the De Novo Assemble tool with default settings. The

567   find heterozygotes tool (peak similarity set to 50%) was used to find and discard any sequences

568   found to be heterozygous. Reads were trimmed to 502 bp and aligned (error probability limit of

569   0.001) with sequences from GenBank and Barcode of Life databases using MUSCLE Alignment

570   in Geneious with default settings.

571        To evaluate whether we were adequately sampling mtDNA haplotype diversity, we

572   plotted rarefaction curves (estimates of haplotype richness by sampling effort) for each

573   population using iNEXT[62] and predicted the total haplotypes for each population assuming 1,000

574    sampled individuals. A median-joining haplotype network was created using POPART[63] for all

575    populations and for each population separately.

576        In an effort to further pinpoint whether the introductions in North America came from

577    western (i.e., United Kingdom) or southwestern (i.e., Spain and France) Europe, we estimated

578    the minimum number of individuals that would need to be sampled from each of these native

579    populations to generate the mtDNA diversity found in North America. Specifically, for each

580    native subpopulation, we randomly sampled (with replacement) a haplotype from each

581    subpopulation based on their haplotype frequencies, until all haplotypes represented in North

582    America were sampled and simulated this procedure 10,000 times for each subpopulation. This

583    approach assumes that the true source population will be the most parsimonious—i.e., require

584    sampling of fewer individuals to create the diversity found in North America.

585

596    2016-29083 and by the National Geographic Society (grant WW1-300R-18). E.A.H. was

597    supported by a Marie Curie Actions IO Fellowship no. 330136.

598

599    **Competing Interests:** the authors declare no competing interests.

600

601    **Data Availability:**

602        Demultiplexed ddRADseq reads generated in this study are available through NCBI's

603    Sequence Read Archive associated with Bioproject (<ID>, SRA: <ID>). All new *COI* sequences

604    were deposited to the Barcode of Life Database (BOLD; <ID>). All metadata and scripts

605    associated with analyses in this study have been deposited on DRYAD (<link>).

606

607    **Contributions:**

608        Author contributions: S.F.R. designed the research project. S.F.R conceived of, created,

609    implemented, and runs the citizen science project—Pieris Project; S.F.R. performed all

610    molecular work, with assistance from M.M.D. with prepping ddRAD libraries and M.A.R. with

611    developing mtDNA primers; S.F.R. performed genomic diversity and structure analyses; E.L.

612    conducted ABC-RF analyses with assistance from S.F.R.; S.F.R., R.V., G.T., V.D., A.E., E.A.H.

613    contributed specimens and/or mtDNA sequences; S.F.R., M.W.E. and A.E. designed educational

614    material related to the citizen science project; and S.F.R. wrote the manuscript with contributions

615    from E.L., A.E., R.V., G.T., V.D., M.A.R., M.M.D., M.W.E., E.A.H., Y.Y., M.E.P., and D.D.S.

616    All authors read and approved the final manuscript.

617

618    **Competing Financial Interests**

619    The authors declare no competing financial interests.

**References**:

1.  Meyerson, L. A. & Mooney, H. A. Invasive alien species in an era of globalization. *Front. Ecol. Environ.* **5**, 199–208 (2007).

2.  Seebens, H. *et al.* No saturation in the accumulation of alien species worldwide. *Nat. Commun.* **8**, 14435 (2017).

3.  Westphal, M. I., Browne, M., MacKinnon, K. & Noble, I. The link between international trade and the global distribution of invasive alien species. *Biol. Invasions* **10**, 391–398 (2008).

4.  Chen, Y. H. Crop domestication, global human-mediated migration, and the unresolved role of geography in pest control. *Elem Sci Anth* **4**, 000106 (2016).

5.  Estoup, A. & Guillemaud, T. Reconstructing routes of invasion using genetic data: why, how and so what? *Mol. Ecol.* **19**, 4113–4130 (2010).

6.  Ryan, S. F. *et al.* The role of citizen science in addressing grand challenges in food and agriculture research. *Proc. Biol. Sci.* **285**, (2018).

7.  McKinley, D. C. *et al.* Citizen science can improve conservation science, natural resource management, and environmental protection. *Biol. Conserv.* **208**, 15–28 (2017).

8.  Hely, P. C., Gellatley, J. G., Pasfield, G. & Agriculture, N. S. W. D. of. *Insect pests of fruit and vegetables in NSW*. (Clayton, Vic. : Inkata Press, 1982).

9.  Fukano, Y., Satoh, T., Hirota, T., Nishide, Y. & Obara, Y. Geographic expansion of the cabbage butterfly (Pieris rapae) and the evolution of highly UV-reflecting females. *Insect Sci.* **19**, 239–246 (2012).

10. Hiura, I. Monshirochou-zoku no Rekishi. *Konchu Shizen* **3**, 9–15 (1968).

11. Scudder, S. H. & History, B. S. of N. *The introduction and spread of Pieris rapae in North America, 1860-1885 [i.e. 1886]*. (Boston Society of Natural History, 1887).

12. Ashby, J. W. & Pottinger, R. P. Natural regulation of Pieris rapae Linnaeus (Lepidoptera : Pieridae) in Canterbury, New Zealand. *N. Z. J. Agric. Res.* **17**, 229–239 (1974).

13. Braby, M. F., Upton, M. S., Collection, A. N. I. & Entomology, C. *The butterflies of Australia : their identification, biology and distribution*. (Melbourne : CSIRO Publishing, 2000).

14. Seiter, S. & Kingsolver, J. Environmental determinants of population divergence in life-history traits for an invasive species: climate, seasonality and natural enemies. *J. Evol. Biol.* **26**, 1634–1645 (2013).

15. McQueen, E. W. & Morehouse, N. I. Rapid Divergence of Wing Volatile Profiles Between Subspecies of the Butterfly Pieris rapae (Lepidoptera: Pieridae). *J. Insect Sci.* **18**, (2018).

652   16.   Qi, X. *et al.* Genomic inferences of domestication events are corroborated by written records in
653          Brassica rapa. *Mol. Ecol.* **26**, 3373–3388 (2017).

654   17.   Prakash, S., Wu, X.-M. & Bhat, S. R. History, Evolution, and Domestication of Brassica Crops. in
655          *Plant Breeding Reviews* 19–84 (Wiley-Blackwell, 2011). doi:10.1002/9781118100509.ch2

656   18.   Maggioni, L., von Bothmer, R., Poulsen, G. & Lipman, E. Domestication, diversity and use of
657          Brassica oleracea L., based on ancient Greek and Latin texts. *Genet. Resour. Crop Evol.* **65**, 137–159
658          (2018).

659   19.   Maggioni, L. Domestication of Brassica oleracea L. (2015). Available at:
660          https://pub.epsilon.slu.se/12424/. (Accessed: 15th November 2018)

661   20.   Clark, P. U. *et al.* The Last Glacial Maximum. *Science* **325**, 710–714 (2009).

662   21.   Hewitt, G. M. Genetic consequences of climatic oscillations in the Quaternary. *Philos. Trans. R. Soc.*
663          *B Biol. Sci.* **359**, 183–195 (2004).

664   22.   Hewitt, G. The genetic legacy of the Quaternary ice ages. *Nature* **405**, 907–913 (2000).

665   23.   Fischer, M. L. *et al.* Historical Invasion Records Can Be Misleading: Genetic Evidence for Multiple
666          Introductions of Invasive Raccoons (Procyon lotor) in Germany. *PloS One* **10**, e0125441 (2015).

667   24.   Jones, R. E., Gilbert, N., Guppy, M. & Nealis, V. Long-Distance Movement of Pieris rapae. *J. Anim.*
668          *Ecol.* **49**, 629–642 (1980).

669   25.   Williams, C. *The migration of butterflies*. (Oliver & Boyd, 1930).

670   26.   John, E., Cottle, N., McArthur, A. & Markis, C. Eastern Mediterranean migrations of Pieris rapae
671          (Linnaeus, 1758) (Lepidoptera: Pieridae): observations in Cyprus, 2001 and 2007. *Entomol. Gaz.* **59**,
672          71–78 (2008).

673   27.   Shen, J. *et al.* Complete genome of Pieris rapae, a resilient alien, a cabbage pest, and a source of
674          anti-cancer proteins. *F1000Research* **5**, 2631 (2016).

675   28.   Atack, J. Historical Geographic Information Systems (GIS) database of U.S. Railroads for 1830-1972.
676          (2016).

677   29.   Opler, P. A. & Krizek, G. O. *Butterflies East of the Great Plain: an illustrated natural history.* (Johns
678          Hopkins University Press, 1984).

679   30.   *New Zealand Official Yearbook (NZOYB)*. (1930).

680   31.   Allendorf, F. W. & Lundquist, L. L. Introduction: Population Biology, Evolution, and Control of
681          Invasive Species. *Conserv. Biol.* **17**, 24–30 (2003).

682   32.   Kingsolver, J. G., Massie, K. R., Ragland, G. J. & Smith, M. H. Rapid population divergence in
683          thermal reaction norms for an invading species: breaking the temperature-size rule. *J. Evol. Biol.*
684          **20**, 892–900 (2007).

685 33. Peterson, B. K., Weber, J. N., Kay, E. H., Fisher, H. S. & Hoekstra, H. E. Double Digest RADseq: An
686      Inexpensive Method for De Novo SNP Discovery and Genotyping in Model and Non-Model Species.
687      *PLoS ONE* **7**, e37135 (2012).

688 34. Ryan, S. F. *et al.* Climate-mediated hybrid zone movement revealed with genomics, museum
689      collection, and simulation modeling. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E2284–E2291 (2018).

690 35. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data.
691      *Bioinformatics* **30**, 2114–2120 (2014).

692 36. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
693      *Bioinforma. Oxf. Engl.* **25**, 1754–1760 (2009).

694 37. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation
695      DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).

696 38. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
697      generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).

698 39. Knaus, B. J. & Grünwald, N. J. vcfr: a package to manipulate and visualize variant call format data in
699      R. *Mol. Ecol. Resour.* **17**, 44–53 (2017).

700 40. Ryan, S. F. *et al.* Patterns of divergence across the geographic and genomic landscape of a butterfly
701      hybrid zone associated with a climatic gradient. *Mol. Ecol.* **26**, 4725–4742 (2017).

702 41. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J.
703      Mol. Biol.* **215**, 403–410 (1990).

704 42. Shiryev, S. A., Papadopoulos, J. S., Schäffer, A. A. & Agarwala, R. Improved BLAST searches using
705      longer words for protein seeding. *Bioinforma. Oxf. Engl.* **23**, 2949–2951 (2007).

706 43. Wang, J. *et al.* SilkDB: a knowledgebase for silkworm biology and genomics. *Nucleic Acids Res.* **33**,
707      D399-402 (2005).

708 44. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
709      individuals. *Genome Res.* **19**, 1655–1664 (2009).

710 45. Raj, A., Stephens, M. & Pritchard, J. K. fastSTRUCTURE: variational inference of population
711      structure in large SNP data sets. *Genetics* **197**, 573–589 (2014).

712 46. Jombart, T., Devillard, S. & Balloux, F. Discriminant analysis of principal components: a new
713      method for the analysis of genetically structured populations. *BMC Genet.* **11**, 94 (2010).

714 47. Jombart, T. adegenet: a R package for the multivariate analysis of genetic markers. *Bioinforma.
715      Oxf. Engl.* **24**, 1403–1405 (2008).

716 48. Francis, R. M. pophelper: an R package and web app to analyse and visualize population structure.
717      *Mol. Ecol. Resour.* **17**, 27–32 (2017).

718    49.    Kamvar, Z. N., Tabima, J. F. & Grünwald, N. J. Poppr: an R package for genetic analysis of
719           populations with clonal, partially clonal, and/or sexual reproduction. *PeerJ* **2**, e281 (2014).

720    50.    Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R language.
721           *Bioinformatics* **20**, 289–290 (2004).

722    51.    Rambaut, A. *FigTree v1.4: Tree Figure Drawing Tool*. (2009).

723    52.    Weir, B. S. & Cockerham, C. C. Estimating F-Statistics for the Analysis of Population Structure.
724           *Evolution* **38**, 1358–1370 (1984).

725    53.    Beaumont, M. A., Zhang, W. & Balding, D. J. Approximate Bayesian Computation in Population
726           Genetics. *Genetics* **162**, 2025–2035 (2002).

727    54.    Cornuet, J.-M. *et al.* DIYABC v2.0: a software to make approximate Bayesian computation
728           inferences about population history using single nucleotide polymorphism, DNA sequence and
729           microsatellite data. *Bioinformatics* **30**, 1187–1189 (2014).

730    55.    Estoup, A. *et al.* Estimation of demo-genetic model probabilities with Approximate Bayesian
731           Computation using linear discriminant analysis on summary statistics. *Mol. Ecol. Resour.* **12**, 846–
732           855 (2012).

733    56.    Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).

734    57.    Pudlo, P. *et al.* Reliable ABC model choice via random forests. *Bioinformatics* **32**, 859–866 (2016).

735    58.    Lombaert, E. *et al.* Complementarity of statistical treatments to reconstruct worldwide routes of
736           invasion: the case of the Asian ladybird Harmonia axyridis. *Mol. Ecol.* **23**, 5979–5997 (2014).

737    59.    Cornuet, J.-M., Ravigné, V. & Estoup, A. Inference on population history and model checking using
738           DNA sequence and microsatellite data with the software DIYABC (v1.0). *BMC Bioinformatics* **11**,
739           401 (2010).

740    60.    Raynal, L. *et al.* Raynal L, Marin J-M, Pudlo P, Ribatet M, Robert CP, Estoup A. ABC random forests
741           for Bayesian parameter inference. *arXiv* **1605.05537v4**, (2017).

742    61.    Hebert, P. D. N., Stoeckle, M. Y., Zemlak, T. S. & Francis, C. M. Identification of Birds through DNA
743           Barcodes. *PLoS Biol.* **2**, (2004).

744    62.    Hsieh, T. C., Ma, K. H. & Chao, A. iNEXT: an R package for rarefaction and extrapolation of species
745           diversity (Hill numbers). *Methods Ecol. Evol.* **7**, 1451–1456 (2016).

746    63.    Leigh, J. W. & Bryant, D. popart: full-feature software for haplotype network construction.
747           *Methods Ecol. Evol.* **6**, 1110–1116 (2015).

748

749

750 **Figure Legends:**

751 **Fig 1. Sample size by location. a**, ddRADseq (N = 559). **b**, mtDNA (N = 1,002). Size of points
752 corresponds to sample size. Explore these data further through interactive data visualizations.

753 **Fig 2. Global invasion history and patterns of genetic structure and diversity of *Pieris***
754 ***rapae*. a**, Genetic ancestry assignments based on the program Admixture. **b**, Rooted neighbor-
755 joining tree based on Nei's genetic distance. **c**, Among population genetic differentiation based
756 on Weir and Cockerham's $F_{ST}$; New Zealand and Australia are treated separately. **d**, Graphical
757 illustration of divergence scenario chosen in ABC-RF analysis (Table 1), **e**, Geographic
758 representation of divergence scenario with the highest likelihood based on ABC-RF analysis;
759 points are colored based on their population assignment using Admixture (Fig 2a) and dates
760 represent median estimates from ABC-RF analysis. All analysis based on 558 individuals
761 genotyped for 17,917 ddRADseq SNPs. Explore these data further through interactive data
762 visualizations.

763 **Fig 3**. **Patterns of autosomal genetic diversity—observed heterozygosity, pairwise**
764 **nucleotide diversity, and Tajima's *D*—by population.**

765 **Fig 4**. **Global patterns of mitochondrial haplotype diversity. a**, Geographic distribution of all
766 88 mtDNA haplotypes discovered (unique color for each haplotype; see interactive data
767 visualizations to explore individual haplotypes); note points jittered to avoid overlapping
768 (hidden) points, thus coordinates are approximate, and the color used for haplotypes are
769 unrelated to those used in other panels. **b**, Haplotype network inferred using median-joining
770 algorithm and colored by population. Hash marks between haplotypes represent base changes
771 (mutations). **c**, Number of unique mtDNA haplotypes by population as well as subpopulation
772 estimated using a rarefaction approach (see Methods) and plotted by geographic location. **d**,
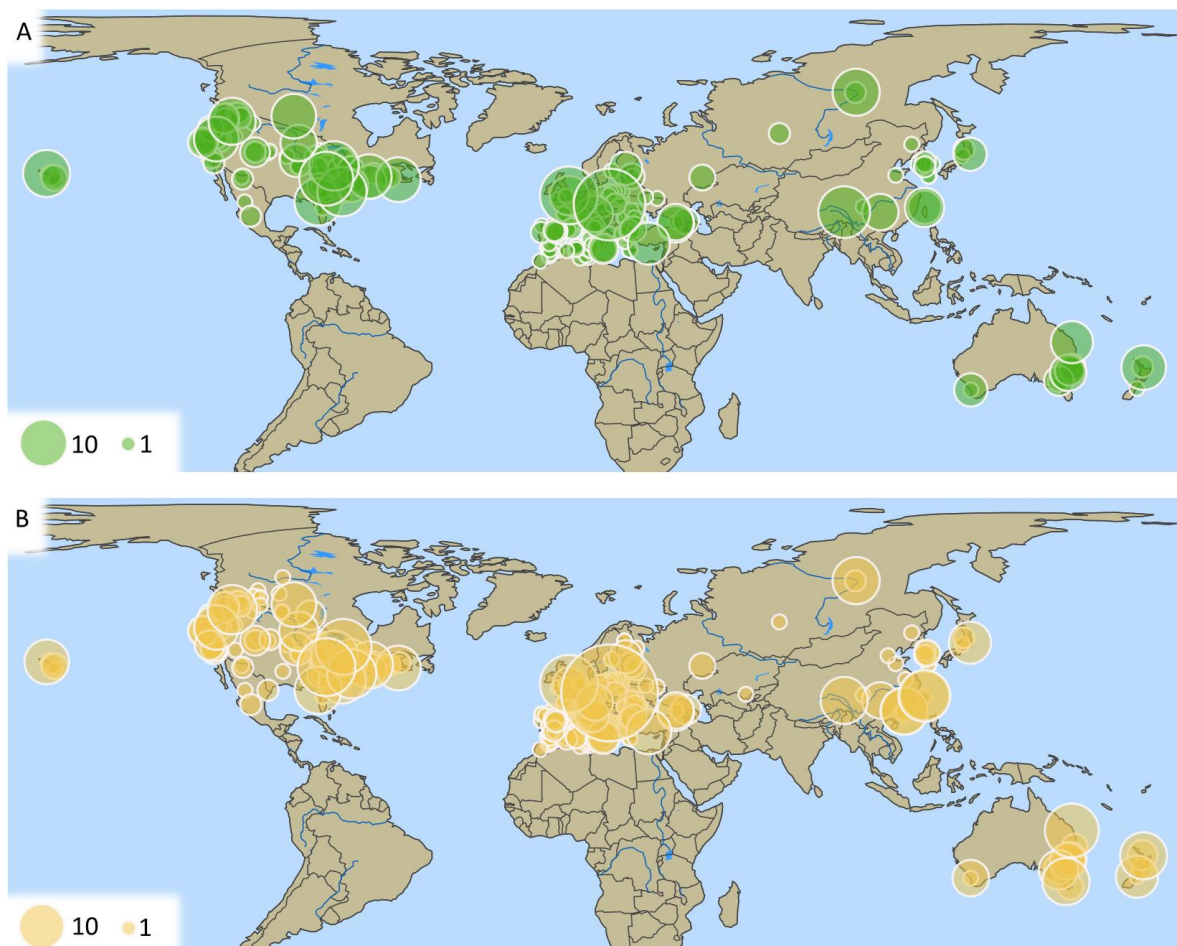773 Pairwise nucleotide diversity by population.

## Tables

**Table 1.** Description of the competing scenarios and results of the six successive ABC analyses to infer the invasion history of *Pieris rapae*.

| Step Scenario | Prior error rate | | | Random forest votes | | | Posterior probability | | |
|---|---|---|---|---|---|---|---|---|---|
| | Data 1 | Data 2 | Data 3 | Data 1 | Data 2 | Data 3 | Data 1 | Data 2 | Data 3 |
| *Analysis 1 - Native area - 18 summary statistics; 13,974 SNPs* | 13.82% | 14.49% | 14.29% | | | | | | |
| S1: Asia is the source of Europe | | | | 57 | 188 | 207 | - | - | - |
| S2: Europe is the source of Asia | | | | 132 | 132 | 43 | - | - | - |
| **S3: Asia and Europe both derived from an ancestral population** | | | | **811** | **680** | **750** | **0.8479** | **0.8173** | **0.8353** |
| *Analysis 2 - Russia and North Africa - 115 summary statistics; 15,533 SNPs* | 16.26% | 17.22% | 17.07% | | | | | | |
| **S1: Asia and Europe are respectively the sources of Russia and Africa** | | | | **602** | **676** | **521** | **0.7020** | **0.7549** | **0.6352** |
| S2: Asia and Africa are respectively the sources of Russia and Europe | | | | 162 | 180 | 146 | - | - | - |
| S3: Africa and Russia are respectively the sources of Europe and Asia | | | | 56 | 30 | 76 | - | - | - |
| S4: Europe and Russia are respectively the sources of Africa and Asia | | | | 180 | 114 | 257 | - | - | - |
| *Analysis 3 - North America east (NAE) - 51 summary statistics; 16,753 SNPs* | 32.82% | 31.83% | 31.82% | | | | | | |
| S1: Asia is the source of NAE, 1 introduction | | | | 0 | 9 | 2 | - | - | - |
| **S2: Europe is the source of NAE, 1 introduction** | | | | **576** | **553** | **559** | **0.5010** | **0.6136** | **0.5064** |
| S3: Asia is the source of NAE, 2 introductions | | | | 6 | 4 | 2 | - | - | - |
| S4: Europe is the source of NAE, 2 introductions | | | | 418 | 434 | 437 | - | - | - |
| *Analysis 4 – North America west (NAW) - 116 summary statistics; 17,049 SNPs* | 11.44% | 11.54% | 10.95% | | | | | | |
| S1: Asia is the source of NAW | | | | 19 | 35 | 13 | - | - | - |
| S2: Europe is the source of NAW | | | | 144 | 121 | 41 | - | - | - |
| **S3: NAE is the source of NAW** | | | | **721** | **720** | **933** | **0.8518** | **0.9288** | **0.9524** |
| S4: Europe is the source of NAW ~ 1600 CE | | | | 85 | 73 | 7 | - | - | - |
| S5: Europe is the source of NAW ~ 1600 CE; NAW is the source of NAE | | | | 31 | 51 | 6 | - | - | - |
| *Analysis 5 – New Zealand - 223 summary statistics; 17,100 SNPs* | 2.18% | 2.30% | 2.18% | | | | | | |
| S1: Asia is the source of New Zealand | | | | 2 | 6 | 5 | - | - | - |
| S2: Europe is the source of New Zealand | | | | 14 | 14 | 28 | - | - | - |
| S3: NAE is the source of New Zealand | | | | 16 | 52 | 130 | - | - | - |
| **S4: NAW is the source of New Zealand** | | | | **968** | **928** | **837** | **0.9739** | **0.9760** | **0.9802** |
| *Analysis 6 - Australia - 388 summary statistics; 17,116 SNPs* | 14.94% | 15.00% | 14.81% | | | | | | |
| **S1: New Zealand is the source of Australia** | | | | **631** | **733** | **613** | **0.7797** | **0.8420** | **0.8124** |
| S2: NAW is the source of Australia | | | | 63 | 33 | 62 | - | - | - |
| S3: New Zealand and Europe are the source of Australia (admixture) | | | | 15 | 10 | 18 | - | - | - |
| S4: New Zealand and Asia are the source of Australia (admixture) | | | | 15 | 7 | 10 | - | - | - |
| S5: New Zealand and NAW are the source of Australia (admixture) | | | | 276 | 217 | 297 | - | - | - |

Results are provided for all three datasets. For each ABC analysis a forest of 1,000 trees was grown. The lines in bold characters corresponds to the selected (most likely) scenarios.

778    **Figures**



779

780    **Fig 1. Sample size by location. a**, ddRADseq (N = 559). **b**, mtDNA (N = 1,002). Size of points
781    corresponds to sample size. Explore these data further through interactive data visualizations.
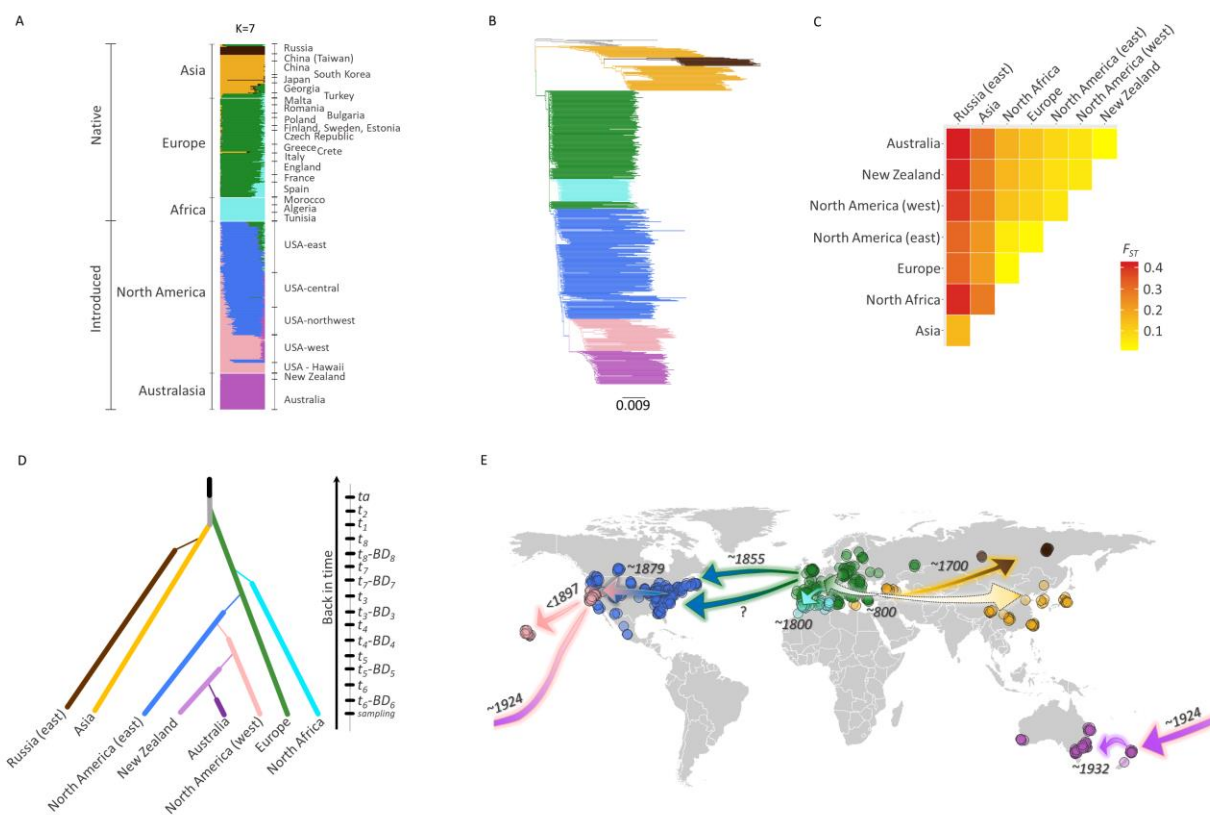
782



783

**Fig 2. Global invasion history and patterns of genetic structure and diversity of *Pieris rapae*. a**, Genetic ancestry assignments based on the program Admixture. **b**, Rooted neighbor-joining tree based on Nei's genetic distance. **c**, Among population genetic differentiation based on Weir and Cockerham's $F_{ST}$; New Zealand and Australia are treated separately. **d**, Graphical illustration of divergence scenario chosen in ABC-RF analysis (Table 1), **e**, Geographic representation of divergence scenario with the highest likelihood based on ABC-RF analysis; points are colored based on their population assignment using Admixture (Fig 2a) and dates (Common Era) represent median estimates from ABC-RF analysis. All analysis based on 558 individuals genotyped for 17,917 ddRADseq SNPs. Explore these data further through interactive data visualizations.

794

795

796

**Fig 3**. **Patterns of autosomal genetic diversity—observed heterozygosity, pairwise nucleotide diversity, and Tajima's *D*—by population.**

799

**Fig 4**. **Global patterns of mitochondrial haplotype diversity. a**, Geographic distribution of all
88 mtDNA haplotypes discovered (unique color for each haplotype; see interactive data
visualizations to explore individual haplotypes); note points jittered to avoid overlapping
(hidden) points, thus coordinates are approximate, and the color used for haplotypes are
unrelated to those used in other panels. **b**, Haplotype network inferred using median-joining

805    algorithm and colored by population. Hash marks between haplotypes represent base changes
806    (mutations). **c**, Number of unique mtDNA haplotypes by population as well as subpopulation
807    estimated using a rarefaction approach (see Methods) and plotted by geographic location. **d**,
808    Pairwise nucleotide diversity by population. Explore these data further through interactive data
809    visualizations.

810 **Supplementary Tables**

811 **Table S1.** Contribution of specimens made by citizen scientists.

812

| Country | Specimens collected by citizen scientists | Total specimens collected | % of specimens collected by citizen scientists |
|---|---|---|---|
| Czech Republic | 43 | 43 | 100% |
| Portugal | 2 | 2 | 100% |
| Gibraltar | 4 | 4 | 100% |
| Turkey | 10 | 10 | 100% |
| South Korea | 11 | 11 | 100% |
| Russia | 16 | 18 | 89% |
| Australia | 62 | 76 | 82% |
| New Zealand | 17 | 25 | 68% |
| USA | 149 | 305 | 49% |
| Spain | 13 | 27 | 48% |
| Canada | 11 | 24 | 46% |
| Romania | 4 | 12 | 33% |
| Bulgaria | 2 | 9 | 22% |
| Algeria | 0 | 12 | 0% |
| Austria | 0 | 2 | 0% |
| China | 0 | 22 | 0% |
| England | 0 | 26 | 0% |
| Estonia | 0 | 4 | 0% |
| Finland | 0 | 7 | 0% |
| France | 0 | 12 | 0% |
| Georgia | 0 | 22 | 0% |
| Greece | 0 | 12 | 0% |
| Italy | 0 | 14 | 0% |
| Japan | 0 | 11 | 0% |
| Malta | 0 | 10 | 0% |
| Mexico | 0 | 6 | 0% |
| Morocco | 0 | 12 | 0% |
| Poland | 0 | 12 | 0% |
| Sweden | 0 | 2 | 0% |
| Taiwan | 0 | 24 | 0% |
| Tunisia | 0 | 12 | 0% |
| Ukraine | 0 | 2 | 0% |

813

814

815 **Table S2.** Prior and posterior distributions of all parameters and several composite parameters of
816 the full final complete scenario (Fig 2d) performed with dataset 1.

| Parameters | Prior distributions | | | | Posterior distributions | | | |
|---|---|---|---|---|---|---|---|---|
| | Q 5% | median | mean | Q 95% | Q 5% | median | mean | Q 95% |
| Raw parameters | | | | | | | | |
| $N_1$ | 159 | 10,110 | 107,800 | 629,178 | 2,630 | 13,662 | 80,180 | 592,324 |
| $N_2$ | 158 | 9,872 | 108,500 | 636,286 | 8,217 | 52,803 | 177,076 | 710,940 |
| $N_3$ | 156 | 9,942 | 107,600 | 626,685 | 10,994 | 184,625 | 300,675 | 881,015 |
| $N_4$ | 159 | 10,010 | 108,200 | 630,976 | 2,305 | 88,066 | 245,186 | 914,933 |
| $N_5$ | 158 | 10,050 | 108,900 | 632,660 | 3,477 | 226,668 | 303,680 | 860,610 |
| $N_6$ | 160 | 10,350 | 108,100 | 628,855 | 1,863 | 119,285 | 294,245 | 884,571 |
| $N_7$ | 158 | 10,140 | 108,900 | 630,040 | 2,033 | 98,892 | 239,822 | 872,549 |
| $N_8$ | 160 | 9,863 | 108,100 | 629,561 | 1,038 | 39,466 | 187,263 | 794,756 |
| $N_A$ | 799 | 68,110 | 193,500 | 792,653 | 26,142 | 167,928 | 247,860 | 829,410 |
| $N_D$ | 126 | 1,469 | 23,160 | 123,273 | 269 | 14,173 | 22,626 | 75,263 |
| $NF_3$ | 3 | 20 | 43 | 159 | 13 | 90 | 98 | 191 |
| $NF_4$ | 3 | 20 | 43 | 159 | 11 | 52 | 63 | 152 |
| $NF_5$ | 3 | 20 | 43 | 159 | 10 | 68 | 77 | 164 |
| $NF_6$ | 3 | 20 | 43 | 158 | 11 | 88 | 91 | 179 |
| $NF_7$ | 3 | 20 | 43 | 159 | 13 | 87 | 89 | 179 |
| $NF_8$ | 3 | 20 | 43 | 159 | 6 | 64 | 71 | 172 |
| $DB_3$ | 2 | 15 | 16 | 29 | 1 | 5 | 6 | 15 |
| $DB_4$ | 2 | 15 | 15 | 29 | 5 | 15 | 16 | 28 |
| $DB_5$ | 2 | 16 | 16 | 29 | 2 | 15 | 15 | 29 |
| $DB_6$ | 2 | 15 | 15 | 29 | 1 | 7 | 9 | 21 |
| $DB_7$ | 2 | 16 | 16 | 29 | 2 | 9 | 11 | 26 |
| $DB_8$ | 2 | 16 | 16 | 29 | 3 | 15 | 15 | 29 |
| $t_1$ | 974 | 4,162 | 4,566 | 9,260 | 908 | 3,576 | 4,159 | 8,849 |
| $t_2$ | 973 | 4,165 | 4,562 | 9,265 | 850 | 3,011 | 3,656 | 8,510 |
| $t_3$ | 467 | 480 | 480 | 494 | 467 | 481 | 480 | 494 |
| $t_4$ | 398 | 411 | 411 | 425 | 397 | 408 | 409 | 424 |
| $t_5$ | 260 | 273 | 273 | 287 | 261 | 274 | 274 | 287 |
| $t_6$ | 235 | 249 | 249 | 262 | 235 | 248 | 248 | 262 |
| $t_7$ | 540 | 1,205 | 1,788 | 5,121 | 511 | 674 | 880 | 1,814 |
| $t_8$ | 539 | 1,200 | 1,783 | 5,119 | 529 | 859 | 1,057 | 2,312 |
| $t_a$ | 14,547 | 55,320 | 55,170 | 95,499 | 14,751 | 60,482 | 58,546 | 96,484 |
| Composite parameters | | | | | | | | |
| $BNsev_3$ | 42 | 6,208 | 180,900 | 900,381 | 203 | 1,062 | 11,425 | 51,050 |
| $BNsev_4$ | 41 | 6,155 | 181,600 | 907,241 | 2,513 | 57,450 | 153,718 | 635,627 |
| $BNsev_5$ | 41 | 6,186 | 181,700 | 913,085 | 618 | 23,343 | 99,199 | 459,562 |
| $BNsev_6$ | 42 | 6,290 | 182,500 | 924,448 | 151 | 12,860 | 40,460 | 134,639 |
| $BNsev_7$ | 41 | 6,053 | 183,900 | 930,965 | 539 | 6,793 | 34,863 | 167,428 |
| $BNsev_8$ | 41 | 6,162 | 179,500 | 886,500 | 309 | 5,580 | 13,628 | 46,122 |

Note: $BNsev_i$ = bottleneck severity of population $i$ computed as $[BD_i \times N_{parental\ population\ of\ population\ i}) / NF_i]$,
with parental populations being populations 2, 3, 4, 5, 2 and 1 for populations 3, 4, 5, 6, 7 and 8
817 respectively.

818

819    **Table S3.** Metadata for specimens used in this study.

820    Included as supplementary file (too large)

821

822    **Table S4.** Populations used for ABC-RF analyses.

823    Included as supplementary file (too large)

824

825 **Table S5.** Prior distributions of demographic and historical parameters used in ABC analyses

826 processed to retrace the worldwide invasion routes of Pieris rapae.

| Parameters | Distribution | Quantile 5% | Median | Mean | Quantile 95% |
|---|---|---|---|---|---|
| $N_D$ | Log-Uniform [100 – 1,000,000] | 126 | 1,482 | 23,610 | 128,053 |
| $N_A$ | Log-Uniform [100 – 1,000,000] | 774 | 67,560 | 192,800 | 790,045 |
| $N_j, N_i, N_{ia}, N_{ib}$ | Log-Uniform [100 – 1,000,000] | 159 | 10,280 | 108,600 | 629,089 |
| $NF_j, NF_i, NF_{ia}, NF_{ib}$ | Log-Uniform [2 – 200] | 3 | 20 | 43 | 159 |
| $BD_j, BD_i, BD_{ia}, BD_{ib}$ | Uniform [1 – 30] | 2 | 16 | 15 | 29 |
| $ta$ | Uniform [10,000 – 100,000] | 14,547 | 54,930 | 54,990 | 95,453 |
| $t_j$ | Log-Uniform [500 – 10,000] | 585 | 2,265 | 3,197 | 8,617 |
| $t_i, t_{ia}$ | Uniform [$x_i - x_i+30$] | DV | DV | DV | DV |
| $t_{mix}, t_{ib}$ | Uniform [$165 - x_i+30$] | DV | DV | DV | DV |
| $t_{4old}$ | Uniform [1245 – 1275] | 1,246 | 1,260 | 1,260 | 1,274 |
| $ar_i$ | Uniform [0.1 – 0.9] | 0.14 | 0.50 | 0.50 | 0.86 |

Notes: Index $i$ stands for the number of the invasive population, i.e. 3, 4, 5 or 6 for North America (east), North America (west), New Zealand or Australia respectively. Index $j$ stands for the number of the ancient putative native population, i.e. 1, 2, 7 or 8 for Asia, Europe, Africa or Russia (east) respectively. $N_D$ and $N_A$ = stable effective population size (number of diploid individuals) of the ancestral native population respectively before and after a demographic expansion event ($N_G < N_A$); $N_j, N_i$ = stable effective population size (number of diploid individuals) of the putative native and invasive populations; $NF_i$ = effective number of founders during a bottleneck lasting $BD_i$ generation(s) for population $i$; $ta$ = time of the demographic expansion in the ancestral native population; $t_j$ = merging time of the putative native populations into the ancestral one; $t_i$ = introduction time of invasive populations $i$ with bounds $x_i$ fixed from dates of first observation of established population; $t_{4old}$ corresponds to the particular case of an old introduction hypothesis of the North American (west) population in ABC analysis 4; $N_{ia}, N_{ib}, NF_{ia}, NF_{ib}, BD_{ia}$ and $BD_{ib}, t_{ia}, t_{ib}$ and $t_{mix}$ are the parameters associated to an admixture event leading to the formation of invasive population $i$; $ar_i$ = admixture rate. Depending on the scenarios considered, various conditions were applied to times so that coalescent times fit with each scenario's topology. All times are expressed in number of generations assuming 3 generations per year, and running back in time from time 0 which corresponds to year 2015. All prior quantities presented were computed from $10^5$ values. DV = different values were possible. See Figure S3 for a graphical representation of the evolutionary scenarios with

827 associated historical and demographic parameters considered in the ABC analyses.

828

829     **Table S6.** Summary statistics used in all DIYABC simulations (Cornuet et al., 2014).

| DIYABC abbreviation | Description |
|---|---|
| *Single sample statistics for each sampled population* | |
| HP0 | Proportion of loci with zero gene diversity |
| HM1 | Mean gene diversity across polymorphic loci (Nei, 1987) |
| HV1 | Variance of gene diversity across polymorphic loci |
| HMO | Mean gene diversity across all loci |
| *Two sample statistics for each pairwise sample combination* | |
| FP0 | Proportion of loci with zero F ST distance (Weir & Cockerham, 1984) |
| FM1 | Mean across loci of non-zero F ST distances |
| FV1 | Variance across loci of non-zero F ST distances |
| FMO | Mean across loci of F ST distances |
| NP0 | Proportion of loci with zero Nei's distance (Nei, 1972) |
| NM1 | Mean across loci of non-zero Nei's distances |
| NV1 | Variance across loci of non-zero Nei's distances |
| NMO | Mean across loci of Nei's distances |
| *Admixture statistics (Choisy et al., 2004) for each combination of parental and admixed populations* | |
| AP0 | Proportion of loci with zero admixture estimates |
| AM1 | Mean across loci of non-zero admixture estimate |
| AV1 | Variance across loci of non-zero admixture estimated |
| AMO | Mean across all locus admixture estimates |

830

831

832    **Supplementary Figures**

833



834

835    **Fig S1.** Sample sizes of *Pieris rapae* specimens by year collected.
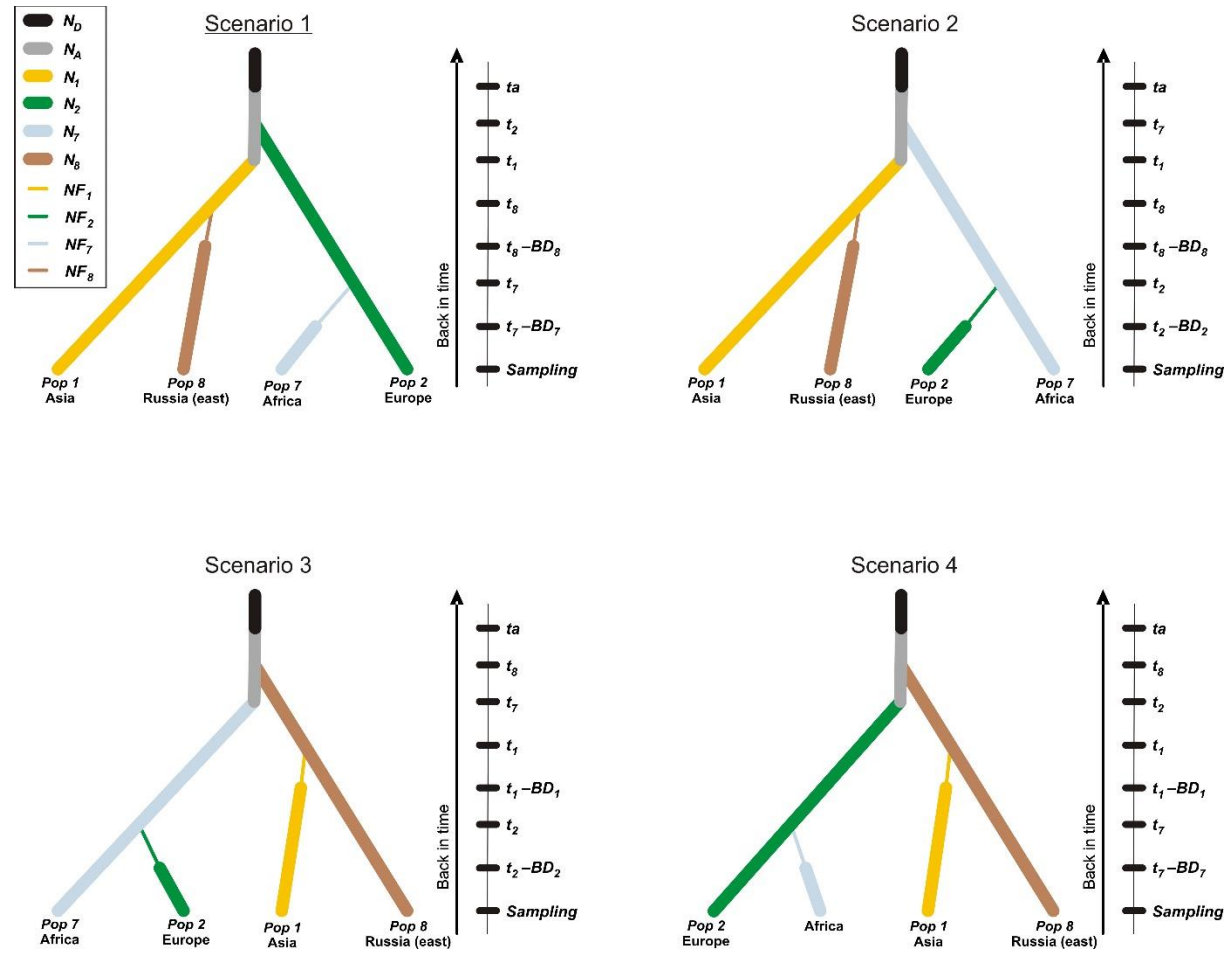
836

837

**Fig S2.** Population ancestry assignment plots for K:2-30, using **a**, ADMIXTURE, **b**,
fastSTRUCTURE, and **c**, Discriminant Analysis of Principal Components (DAPC). For each
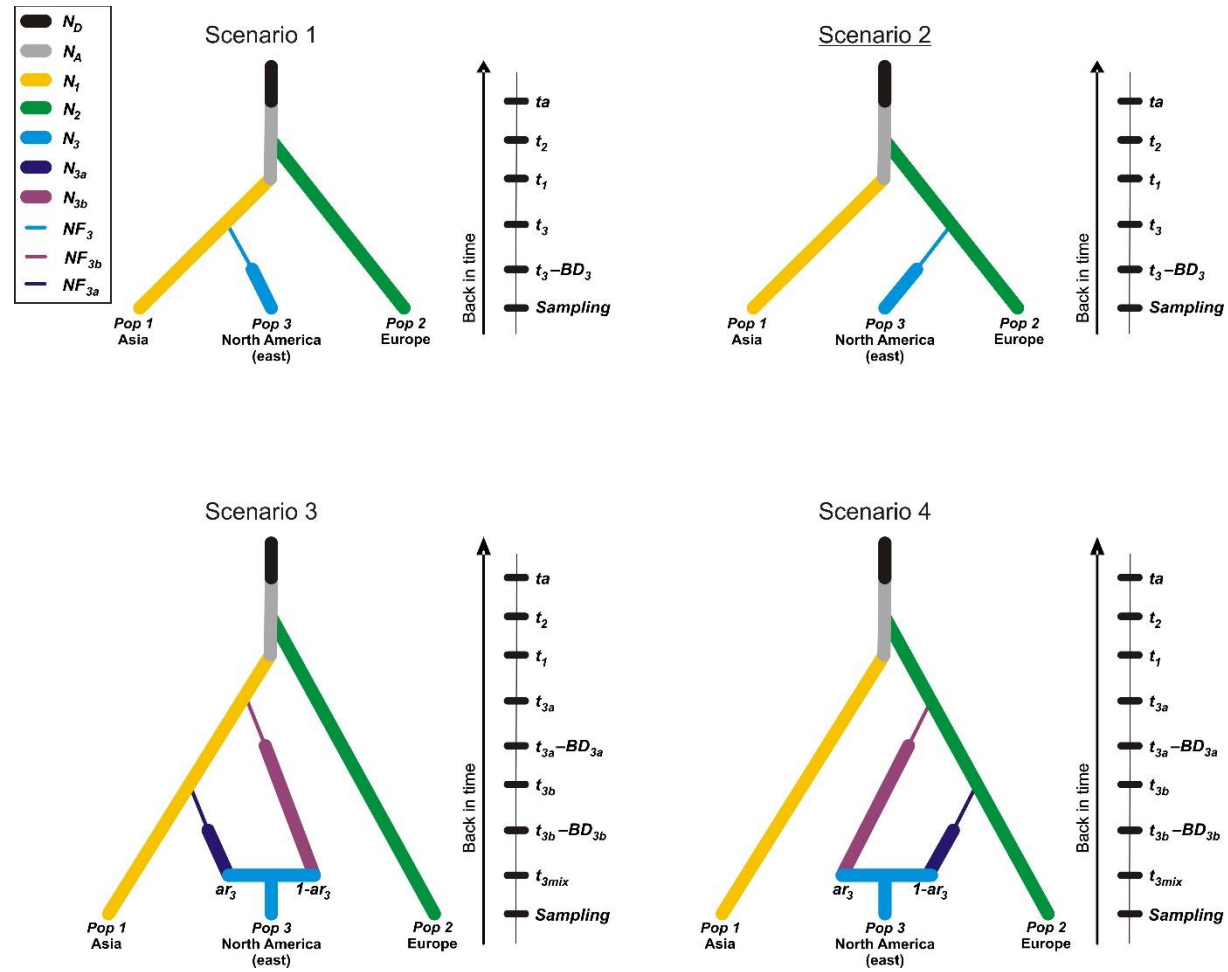analysis the evaluation for optimal K is included.
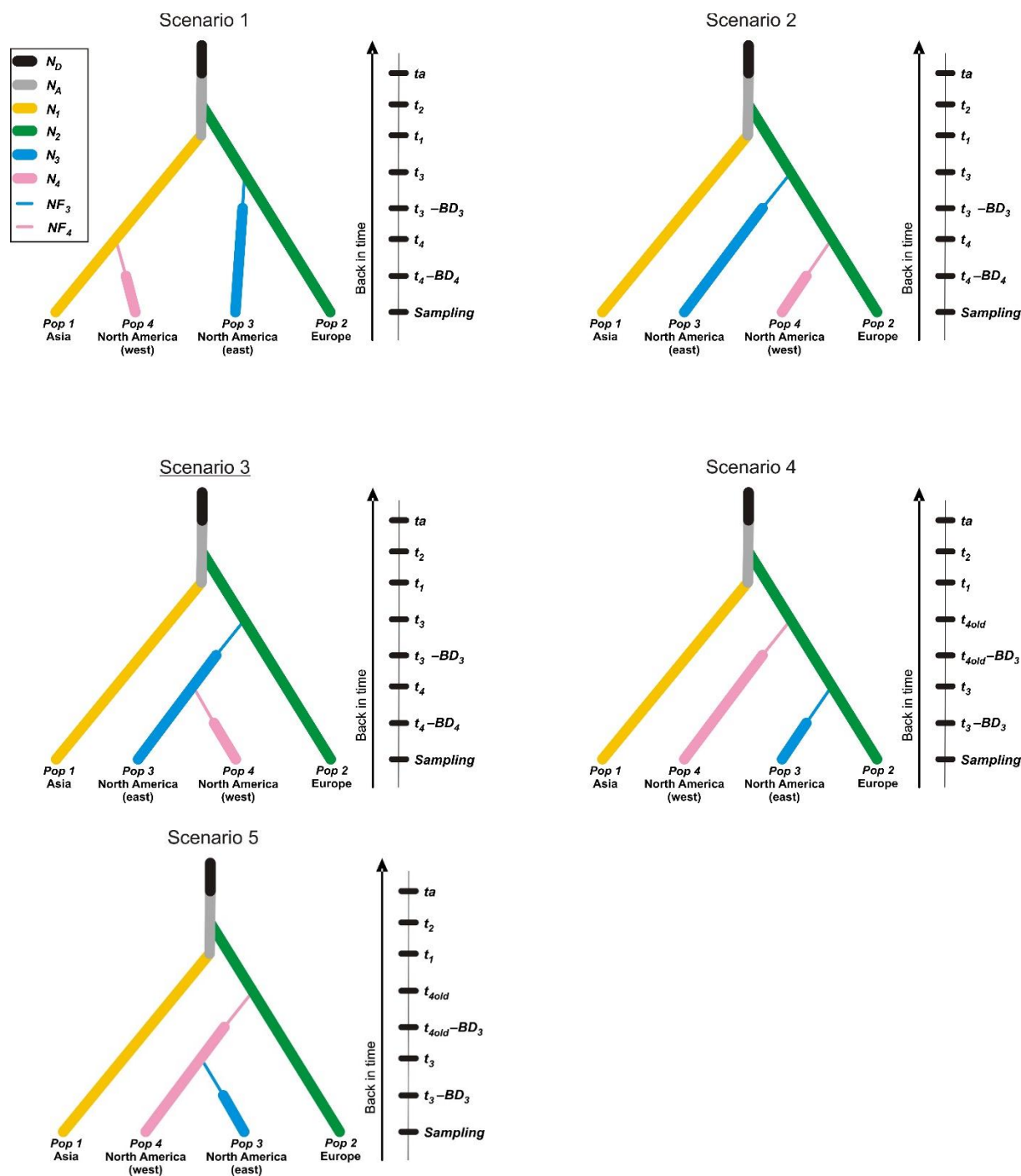
841

842

## A. Analysis 1 – Native area

843

844



845

846

847

848    B. Analysis 2 – Russia (east) and North Africa

849

850



851

852

853

854    C. Analysis 3 – North America (east)

855

856



857

858

859

860 D. Analysis 4 – North America (west)



861



862

863

864

865 | E. Analysis 5 – New Zealand



866

867

868

869

## F. Analysis 6 – Australia

870



871

872

873 **Fig S3.** Schematic representation of each set of scenarios used in the ABC analyses to decipher
874 the worldwide invasion routes of *Pieris rapae* (see also Table 1). Population numbers are as
875 follows: 1 for Asia; 2 for Europe; 3 for North America (east); 4 for North America (west); 5 for
876 New Zealand; 6 for Australia; 7 for North Africa; 8 for Russia (east). For each analysis, the name
877 of the most likely scenario is underlined. Thin lines indicate bottlenecks. For parameters
878 descriptions and priors see Table S5. Time is not to scale.

879
**Fig S4.** Median-joining haplotype networks for each population. Hash marks between haplotypes
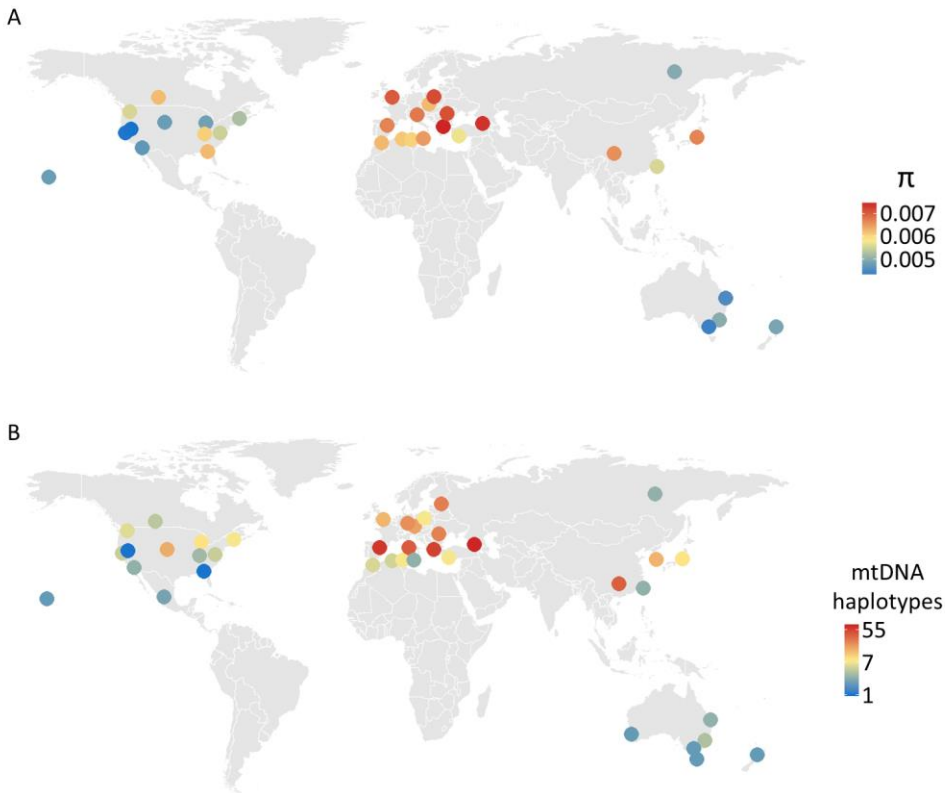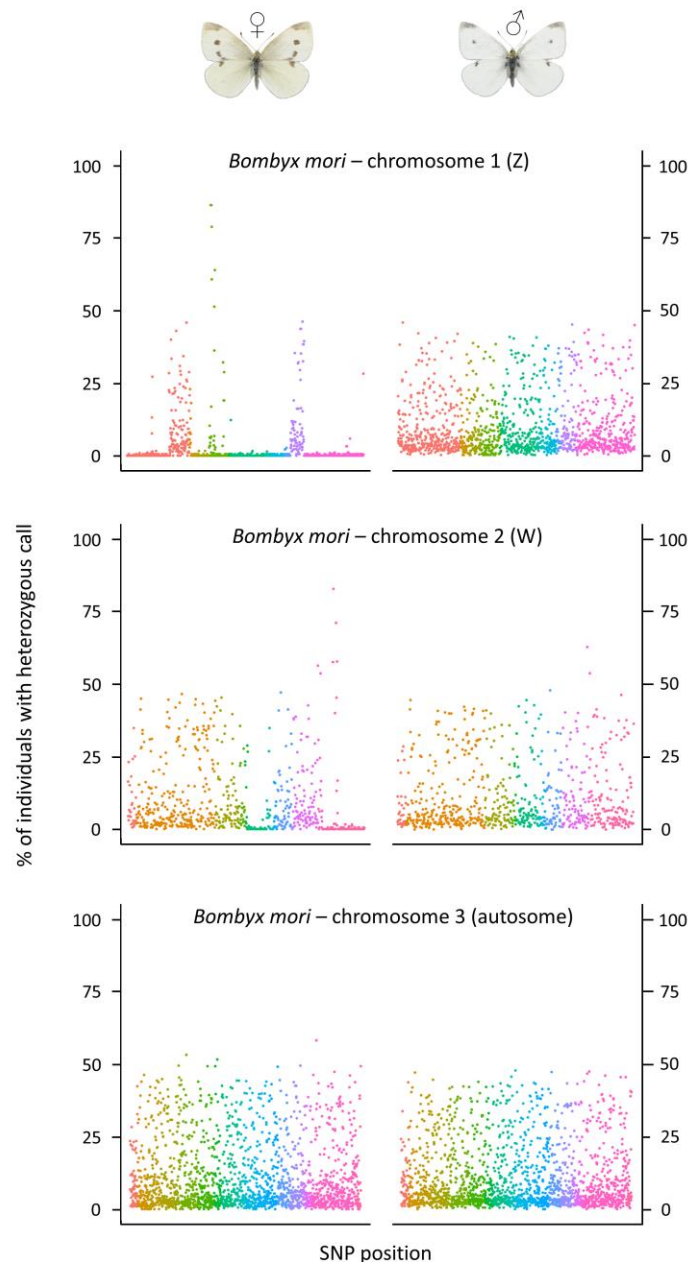represent base changes (mutations).

882

883



884

**Fig S5. Global patterns of genetic diversity. a,** Estimate of pairwise nucleotide diversity for each subpopulation based on autosomal ddRADseq data. **b**, mtDNA haplotype diversity estimated from rarefaction curves (note, colors are based on a log scale).

888

**Fig S6. Percentage of individuals with heterozygous calls for each locus, plotted separately females and males.** The location of each locus is based on its position within each *P. rapae* scaffold, with each *P. rapae* scaffold then ordered in each *B. mori* chromosome based on its homology to each *B. mori* scaffold (see Methods). Loci are colored by the *B. mori* scaffold to which they are associated. An autosome (chromosome 3) is plotted for reference and the pattern reflects those observed in other autosomes—no discernable difference in heterozygosity between males and females. Note, the W chromosome was not sequenced or assembled in the reference genome used in this study and is thus likely to be made up of portions of other chromosomes, including the Z (regions with no heterozygosity in females).

898    *Video included as a supplementary file*
899    **Video S1.** Development of railroad lines in the United States from 1830-1972. Railroad line data
900    were obtained from Atack, 2016[28] and plotted by their date of operation. Note the competition of
901    railroad lines connecting eastern and western US in 1872, a few years prior (1879) to when a
902    small population originating from North America (east) was believed to be introduced to that
903    exact region—North America (west) (i.e., central California).