

23 **Abstract**

24

25 Population genetics studies on non-model organisms typically involve sampling few
26 markers from multiple individuals. Next-generation sequencing approaches open up the
27 possibility of sampling many more markers from fewer individuals to address the same
28 questions. Here, we applied a target gene capture method to deep sequence ~1000
29 independent autosomal regions of a non-model organism, the blacktip reef shark
30 (*Carcharhinus melanopterus*). We devised a sampling scheme based on the predictions of
31 theoretical studies of metapopulations to show that sampling few individuals, but many
32 loci, can be extremely informative to reconstruct the evolutionary history of species. We
33 collected data from a single deme (SID) from Northern Australia and from a scattered
34 sampling representing various locations throughout the Indian Ocean (SCD). We explored
35 the genealogical signature of population dynamics detected from both sampling schemes
36 using an ABC algorithm. We then contrasted these results with those obtained by fitting the
37 data to a non-equilibrium finite island model. Both approaches supported an Nm value ~40,
38 consistent with philopatry in this species. Finally, we demonstrate through simulation that
39 metapopulations exhibit greater resilience to recent changes in effective size compared to
40 unstructured populations. We propose an empirical approach to detect recent bottlenecks
41 based on our sampling scheme.

42

43

44

45 **Introduction**

46

47 Obtaining sufficient sequence information is no longer a limiting factor for
48 deciphering the demographic and selective history of a species. Many studies now use
49 population genomic approaches to shed light on evolutionary questions that also have
50 practical applications for species conservation and management. While next-generation
51 sequencing (NGS) techniques have made such approaches possible, studying non-model
52 organisms without a reference genome remains a challenge. Thus far, whole-genome re-
53 sequencing¹, transcriptome analysis² and *de novo* restriction site associated DNA (RAD)
54 sequencing³ have been used to investigate the population dynamics of non-model
55 organisms. All of these approaches have limitations: re-sequencing whole genomes of many
56 individuals remains expensive and superfluous when reliable results can be obtained using
57 a subset of the genome⁴. Generating transcriptome data sets suitable for investigating
58 adaptive change can be complicated by the requirement for high quality RNA and sampling
59 equivalent life stages and tissue types for all individuals⁵. Short sequence lengths, high rates
60 of allelic dropout⁶ and missing data⁷, as well as the vagaries of the bioinformatics pipeline
61 that is used⁸, can influence estimates of genetic diversity obtained with RAD-seq
62 approaches. Target gene capture approaches offer a promising alternative to these methods.
63 They permit sequencing of pre-selected regions of the genome (see the review from⁹) that
64 have been determined *a priori* to be informative for the research question and are generally
65 effective across a range of sample types and qualities. This allows the generation of
66 relatively unbiased, more complete datasets than can be obtained using RAD technologies,
67 and at an affordable cost compared with whole-genome re-sequencing.

68 The population structure of the study organism is another factor that may complicate
69 our ability to make robust demographic inferences from population genomic data. Previous
70 work has shown that ignoring metapopulation structure can sometimes mislead
71 interpretations of the demographic and evolutionary history of a species¹⁰⁻¹³. This is a
72 problem because populations are rarely entirely isolated in nature but instead belong to a
73 network of demes that exchange migrants at different rates. Despite these concerns, few
74 studies have attempted to estimate demographic parameters under complex metapopulation
75 models (see ¹⁴ and ¹⁵ for some exceptions). The predicted coalescent patterns that are
76 associated with metapopulations sampled at different levels are of special interest. The gene
77 genealogy of a metapopulation is characterized by the product of the effective population
78 size N and the migration rate m ¹⁶⁻¹⁹. The coalescent history of a sample of lineages can be
79 divided into two phases: the *scattering* phase within each single deme and the *collecting*
80 phase when each of the lineages comes from a different deme^{16,20}. This separation of time
81 scales holds true for several metapopulation models, including range expansions¹⁹.
82 Contrasting the shape of the genealogy at different sampling levels (e.g. “single deme” vs.
83 “scatter”, where each lineage comes from a different deme) indirectly provides information
84 on the long term Nm ¹⁸. It is therefore possible to make inferences about the evolutionary
85 history of the metapopulation using relatively simple demographic models that investigate
86 changes in effective population size across specific sampling schemes. A few individuals
87 judiciously sampled from the range distribution of a species can be sufficient to uncover its
88 demography, when enough loci are available for statistical inferences. Such an approach is
89 now tenable with the large datasets that are being generated using NGS and has important

90 applications for the study of endangered animals, or of any organism for which specimens
91 are not easily obtained.

92 Here, we demonstrate the use of such an approach to study the demographic history
93 of a non-model organism, the blacktip reef shark (*Carcharhinus melanopterus*). These animals
94 are highly philopatric and are often found on remote coral islands and atolls^{21,22}. Their
95 proclivity for philopatry and strong habitat preference for coral reefs predisposes them to
96 exhibit a disjunct meta-population structure over their widespread distribution throughout
97 the Indo-Pacific. We devised a two-layered sampling scheme based on the above theoretical
98 considerations associated with metapopulations²⁰: we collected individuals from a single
99 deme from Northern Australia and from a scattered sampling from various locations
100 throughout the Indo-Pacific. We then used a newly developed target gene capture
101 approach²³ to generate sequence data for ~1000 pre-specified independent orthologous
102 autosomal regions. This approach allowed us to minimise the sampling of individuals, while
103 ensuring a comprehensive sampling of loci for a species for which no closely related
104 reference genome exists. Demographic parameters were estimated in two ways: i) indirectly
105 by contrasting the gene genealogies at different sampling levels in the metapopulation
106 (“single deme” vs. “scatter”); ii) directly by applying a non-equilibrium finite island model
107 to both sampling schemes. We developed an Approximate Bayesian Computation (ABC)
108 framework^{24,25} with recombination to estimate parameters and to compare demographic
109 models. The site frequency spectrum (SFS) computed on unphased data was used as
110 summary statistic, thus avoiding the level of uncertainty that is introduced by phasing and
111 haplotype reconstruction.

112 *Carcharhinus melanopterus* inhabits reef flats and sheltered lagoons^{21,22} in the Indian
113 and Pacific Oceans (the “Indo-Pacific”). While their biology has been widely studied their
114 population dynamics and dispersal patterns are largely unknown. Nevertheless, insights
115 into their movements and population health status can be gained through assessments of
116 their evolutionary history. The species is considered globally near threatened (NT) by the
117 International Union for Conservation of Nature (IUCN) Red list, with a decreasing
118 population²⁶. *Carcharhinus melanopterus* appears to be locally abundant in many areas but
119 recent fishing and anthropogenic pressures on reef environments may have affected their
120 distribution and population dynamics²⁶⁻²⁸. A recent bottleneck was detected in Moorea²⁹
121 suggesting that decreases in effective population size may also have occurred in some parts
122 of its range. In light of this observation, we used an ABC framework incorporating extensive
123 simulation to explore the population history of *C. melanopterus*, focusing particularly on
124 signals in the data that might reflect recent changes in population size. In doing so, we also
125 use this study system as a test case to explore the extent to which metapopulation structure
126 can hinder the detection of a recent bottleneck and propose an empirical approach and
127 sampling scheme that addresses this issue.

128

129 **Results**

130 *Genetic diversity and data summary*

131 The demographic history of the metapopulation of *C. melanopterus* was examined using two
132 datasets: the first based on estimates from a “single deme” (SID) in Northern Australia, the
133 second based on a collection of samples taken from various locations in the Indo-Pacific,
134 which we refer to as the “scatter” sample (SCD) (see Fig. 1, Table 1 and Supplementary

135 Table 1). We sequenced 18 individuals in total with an average coverage of 75×. After
136 applying strict filters and removing duplicate reads (see SI), we obtained sequence data for
137 995 loci for SCD and 998 loci for SID, spanning 606,647 bp and 632,160 bp respectively (Table
138 1 and Supplementary Table S2). The average length of the regions is ~600bp and details of
139 the length distribution for SCD and SID are shown in Supplementary Fig. 1. Overall, 2605
140 and 1946 high quality single nucleotide polymorphisms (SNPs) were called for the scatter
141 and single deme dataset, respectively. As expected, we detected a higher number of SNPs
142 from SCD (considering samples from different locations) than SID, and SNP density is
143 higher in introns than in exons for both datasets (no significant difference was found
144 between SNP density in intron 5' and intron 3'; Supplementary Table S2). We also performed
145 a Principal Component Analysis (PCA) including all 18 samples to assess the level of
146 population structure within our dataset (Supplementary Fig. S2). At least three distinctive
147 clusters are suggested by the first two principal components, explaining ~40% of the
148 variance and separating Australia from the Indo-Pacific and Oman samples. Additional
149 population sub-structure within the Australian and the Indonesian clusters is highlighted by
150 the third principal component (Supplementary Fig. S2, panel b and d). Overall, this confirms
151 that the signal in our dataset is incompatible with panmixia and thus a simple demographic
152 model of a single isolated population is not appropriate to fully describe the demographic
153 history of these samples.

154

155 *Demographic inferences: indirect approach*

156 First, simple population genetics models were applied to indirectly estimate the
157 demographic parameters of the metapopulation by contrasting the characteristics of the

158 genealogies under the two sampling schemes, SID and SCD. A constant-size model (model
159 COS) and a demographic-change model (model CHG1) were initially compared (Fig. 2,
160 panel a and b). Both models were found to have equal posterior probability (0.5) in SID,
161 suggesting a pattern compatible with either a weak signal of expansion (N_{mod} was slightly
162 higher than N_{anc} , Supplementary Table S3) or a constant population size. This is consistent
163 with the ancestral versus modern effective population size ratio (resize) computed under the
164 CHG1 model (median = 0.48, 95% credible interval (CI) ranging from 0.03 to 2.56). In
165 contrast to these results, model CHG1 showed a posterior probability of 0.89 when tested in
166 SCD and a 95% CI of the resize between 0.03 and 0.54, suggesting a strong signal of
167 expansion. This was confirmed when the effective population size was estimated at different
168 time points in the past with an ABC-skyline reconstruction: a sharp signal of expansion was
169 detected around ~90,000 generations before present, while a constant-size population was
170 clearly observed in SID (Fig. 3). Under CHG1 we were able to detect the time of the
171 expansion of the metapopulation but more recent changes in population size cannot be
172 excluded *a priori*. Therefore, a model with two demographic changes (model CHG2; Fig. 2,
173 panel c) was also tested in order to account for possible more recent demographic events
174 (e.g. post-glacial expansion), as described for other marine species^{30,31}. Posterior probabilities
175 of model CHG1 and model CHG2 are similar for both datasets (Supplementary Table S4).
176 When the effective population size was estimated at different time points in the past, CHG1
177 and CHG2 showed the same pattern (Fig. 3 and Supplementary Fig. S3). The estimation of
178 T_{cl} in CHG2 for both datasets reflects the estimate of T_{c} in CHG1 suggesting that CHG2
179 identifies only the change in effective population size corresponding to the demographic
180 expansion that was also identified by CHG1.

181

182 *Demographic inferences: direct approach*

183 Considering the results from the indirect approach, we therefore explored the demographic
184 signals in the data by fitting a non-equilibrium finite island model²⁰ (model FIM; Fig. 2,
185 panel d) with 100 demes to both datasets. The choice of an island model is justified by the
186 absence of panmixia in the PCA and by the near-uniform G_{st} matrix (Supplementary Table
187 S5), suggesting similar genetic distances across all samples except Oman. Model FIM is
188 defined by three parameters, namely Nm , the time of the onset of the island T_i and the
189 effective size of the ancestral deme, N_{anc} . Demographic parameters of the metapopulation
190 were estimated directly for the two different sampling schemes, SCD and SID. The FIM
191 model showed the highest posterior probability in the model choice for both SID and SCD
192 (Table 2). Cross-validation of the model choice confirmed that FIM is distinguished from
193 COS and CHG1 with high probability for both datasets (see SI and Supplementary Table S6).
194 Importantly, the two datasets produced similar estimates of both Nm and T_i (Table 3). We
195 further confirmed this result by applying FIM to the complete dataset of 18 samples
196 (SCD+SID without shared samples) (Table 3). With this value of $Nm \approx 40$ in the
197 metapopulation, the signature of the ancestral expansion at T_i is lost in SID but not in SCD,
198 consistent with expectation based on results in simulation studies of range expansion (see
199 ^{19,32,33}). Indeed, for similar Nm values the gene genealogy of lineages sampled from a single
200 deme is a mixture of recent coalescent events occurring during the *scattering* phase and more
201 ancient coalescent events occurring during the *collecting* phase. This generates a gene
202 genealogy similar to that typical of a constant size population. Conversely, a scatter sample
203 will show a signature of expansion as most of the coalescent events will occur at the time of

204 the foundation of the metapopulation. Cross-validation of Nm suggests that the estimation is
205 reliable and robust, with bias of the mode around -0.02 and -0.01 for SCD and SID
206 respectively (Supplementary Table 7). We estimated T_i to have occurred ~59,000 generations
207 ago (the average between the point estimates for the onset of the island from SCD and SID;
208 Table 3). We observed an overestimation of the onset of the expansion in CHG1 (as T_c)
209 compared to FIM (as T_i) in SCD. Cross-validation suggests a more robust and accurate
210 estimate of T_i than T_c , with an associated bias of 0.48 and 1.154 respectively (Supplementary
211 Table S7).

212

213 *Detection of a recent bottleneck*

214 We used ABC to test whether the absence of a bottleneck in our dataset reflects the real
215 demographic history of *C. melanopterus*, or a lack of power to detect it. To this end, we
216 simulated pseudo observed data sets (*pods*) under a modified FIM model (FIM-BOTT)
217 characterised by two additional parameters (Supplementary Fig. S4, panel b): T_{bott} (the time
218 when an instantaneous decrease in Nm occurred) and I_{bott} (the ratio of the ancestral Nm to
219 the new one). When $I_{\text{bott}}=1$ FIM-BOTT reduces to FIM (Fig. 2 panel d). FIM-BOTT is defined
220 by two major events: the onset of the island at T_i and the reduction in connectivity (i.e., a
221 reduction in the effective population size of a deme, a reduction in the migration rate, or
222 both) at T_{bott} . We simulated 1,000 *pods* under model FIM-BOTT for each parameter
223 combination and we reanalysed them using the CHG2 model. CHG2 allows for two changes
224 in effective population size through time and it could therefore recover both events of FIM-
225 BOTT. We mimicked both the SID and the SCD sampling schemes and we tested twelve
226 combinations of I_{bott} and T_{bott} , namely two intensities of Nm reduction (100× and 1000×) at six

227 time points (10, 50, 100, 200 500, 1000 generations ago). In this way, we were able to examine
228 their respective behaviours when a recent bottleneck (here represented as a decrease in Nm)
229 was simulated in order to obtain insight regarding the optimal sampling strategy to detect it.
230 We plotted for comparison the results obtained by analysing *pods* simulated under FIM. For
231 $T_{\text{bott}} \leq 50$ generations no signal of the bottleneck was detected for either sampling scheme at
232 any I_{bott} (Fig. 4, Supplementary Fig. S5 and S6). With progressively older T_{bott} , we observed a
233 general reduction of the estimated N_e when compared to the one estimated from *pods*
234 simulated under FIM (Fig. 4, Supplementary Fig. S5, S6). This reduction was detected first in
235 SID and only later (~500 generations ago) in SCD (Supplementary Fig. S5 and S6). However,
236 a change in N_e through time (expected under a bottleneck) was never observed for any
237 combination of I_{bott} and T_{bott} in SCD (Fig. 4 and Supplementary Fig. 5 and 6). SID was more
238 affected by the change in connectivity, but a signature of a bottleneck appeared only for
239 $I_{\text{bott}}=1000\times$ and $T_{\text{bott}} \geq 200$ (Supplementary Fig.S6). The capacity to detect a bottleneck can
240 depend both on the real properties of the gene genealogy and on the inferential method
241 used to analyse the data (the ABC skyline in this case). Therefore, we tested the power of our
242 ABC-skyline approach with *pods* simulated under an unstructured model with bottleneck
243 intensities that were comparable to those simulated in the metapopulation scenario.
244 Specifically, we performed simulations under model CHG1-BOTT (Supplementary Fig. S4,
245 panel a), an unstructured population with demography estimated under CHG1 (in our SCD
246 sample) and the same combinations of I_{bott} and T_{bott} as above. Here, the parameter I_{bott} refers
247 to the decrease in N_e , while for FIM-BOTT it referred to the decrease in Nm . N_e estimated
248 from *pods* simulated under CHG1-BOTT drops compared to the values estimated from *pods*

249 simulated under CHG1 as recently as 10 generations ago (Fig. 4). Moreover, a clear signature
250 of a bottleneck appears as soon as $T_{\text{bott}}=50$ (Fig. 4).

251 Similar to the reasoning followed when studying the real data, we further analysed
252 the *pods* simulated under FIM-BOTT for all combinations of I_{bott} , T_{bott} and sampling scheme
253 with the FIM model. We plotted the median of the posterior distributions of Nm estimated
254 from each *pods* (hereafter, Nm_{est}). To quantify the reduction of Nm_{est} due to the simulated
255 decrease in Nm , we show as comparison Nm_{est} obtained from *pods* simulated under FIM (Fig.
256 5, Supplementary Fig. S7). First of all, we note that both sampling schemes can correctly
257 recover Nm when *pods* are simulated with the FIM model, indicating that this is an
258 appropriate inferential procedure (Fig. 5). We then observed that the two sampling scheme
259 respond differently when *pods* are simulated under FIM-BOTT. Generally, T_{bott} seems to
260 have a stronger effect on Nm_{est} than does I_{bott} (Fig. 5, Supplementary Fig. S7). SCD does not
261 show any significant reduction in Nm_{est} until 500 generations ago, while SID starts showing a
262 decrease almost immediately, after just 10 generations. The decline of Nm_{est} in SID is faster
263 when $I_{\text{bott}}=1,000$ than for $I_{\text{bott}}=100$. In both cases, the distributions of Nm_{est} for SID and SCD up
264 to 500 generations ago show a significant difference. At 1000 generations, both datasets
265 show the same pattern with severe decreases in Nm_{est} (Fig. 5, Supplementary Fig. S7).

266

267 **Discussion**

268

269 In this study we present a novel approach to population genomics that is suitable for
270 application to non-model organisms, particularly those for which it is difficult to obtain
271 large sample sizes. The novelty stems from both the technology applied to produce NGS

272 sequence data and the sampling scheme adopted for the population genomics inferences.
273 First, we used a recently developed target gene capture approach paired with NGS²³ that, so
274 far, has only been used for phylogenomics inference³⁴. Second, we exploit intrinsic
275 differences in the shape of gene genealogies that result from the *collecting* and *scattering*
276 phases of metapopulations^{16,18,32} to demonstrate that it is possible to make demographic
277 inferences indirectly based on a few carefully chosen individuals when using a sufficiently
278 large number of loci. We chose to test the utility of our framework using the black tip reef
279 shark, *C. melanopterus*, because we felt that this species was a good example of the practical
280 challenges that are faced by population and conservation geneticists that are working with
281 non-model systems. Firstly, individuals of this species have a small home range and exhibit
282 strong site fidelity, making it a good example of a non-model organism that exhibits
283 metapopulation structure. Indeed, this species has previously been shown to exhibit
284 extensive population structure and restricted movements in French Polynesia³⁵⁻³⁷, the Great
285 Barrier Reef³⁸ and Western Australia³⁹. Secondly, they mostly occur in remote regions of the
286 world and, as is the case for the majority of free ranging marine animals, it is logistically
287 difficult to sample them in large numbers. Finally, the International Union for the
288 Conservation of Nature has listed *C. melanopterus* as “near threatened”, meaning that any
289 information that can be gleaned regarding the demography and structure of this species can
290 have implications for management. These attributes make *C. melanopterus* a well-suited
291 system for to test the efficacy of our approach for reconstructing the evolutionary dynamics,
292 historical demography and associated conservation implications in a non-model species that
293 has a high degree of metapopulation sub-structure, but from which widespread sampling is
294 difficult.

295 Our approach derives its effectiveness from the separation of time scales (i.e.,
296 *scattering vs collecting* phases), first described elegantly by Wakeley (1999), for the gene
297 genealogies in an island model and later shown to be valid for many metapopulation
298 models, i.e., stepping stone¹⁸, spatially continuous⁴⁰ and range expansion¹⁹. Our indirect
299 approach based on the comparison of the genealogical history at the two sampling levels
300 appears robust to mis-specification of the metapopulation model, but it cannot provide
301 precise estimates of the demographic parameters. To further characterise the underlying
302 demography, we analysed our data directly using a non-equilibrium finite island model.
303 While other metapopulation models might better fit the data, they come at the risk of
304 overparameterisation (i.e., non-symmetric island model, spatial models as a stepping stone).
305 Herein, we demonstrate statistically that even a relatively simple model such as the FIM
306 provides a better proxy to the data and more realistic description of a species history than
307 models that ignore population structure (Table 2).

308 The relative length of the *scattering vs collecting* phase determines the shape of the
309 gene genealogy and is mostly influenced by Nm ^{16,19}. This holds true for lineages sampled
310 within a deme (our SID) and lineages scattered throughout the range of the species (our
311 SCD). If an expansion occurred in the species through the colonisation of new habitat, high
312 Nm values will determine a signature of population growth at all sampling levels¹⁸. For
313 decreasing Nm , the signature of expansion is first lost when sampling lineages from a single
314 deme and then in a scattered sample. We found a significant signature of expansion in SCD
315 and a constant population size in SID. Ray et al. (2003) suggested that at $Nm < 50$ the
316 genealogy in a single deme would mostly resemble that of a constant size population (or
317 even of a population reducing in size) in a range expansion model. When we fit the FIM

318 model independently to both SCD and SID, we found a metapopulation characterised by a
319 value of $Nm \sim 40$ (95% CI: 20-110), with an onset of formation around $\sim 59,000$ (95% CI: 6,000-
320 165,000) generations ago. This level of connectivity is consistent with philopatry in this
321 species, as previously described in Moorea and Tetiaroa (French Polynesia)³⁵. Unfortunately
322 our sampling scheme does not allow us to test males and females separately, but it would be
323 interesting to investigate possible differences in the level of connectivity related to sex-
324 specific migration. We further apply the same FIM model to the whole dataset and we
325 obtained similar values of both Nm and T_i (Table 3). In summary, the indirect approach of
326 comparing gene genealogies under different sampling schemes and directly applying the
327 FIM model (applied to three largely independent datasets) provided consistent results,
328 suggesting that this approach is effective for reconstructing the demography of this species.
329 We stress that even though there is no spatial component in the FIM model, the time of
330 formation of the metapopulation is analogous to the expansion time in a range expansion
331 model. Therefore, we speculate that this time estimate probably represents the colonisation
332 of the Indian Ocean by *C. melanopterus*, through a range expansion. Vignaud et al.²⁹ found
333 evidence for genetic structure and isolation by distance at a larger geographic scale (i.e.,
334 Pacific and Indian Ocean) for this species. Clearly, a FIM model would not be a good
335 approximation in this case and a non-equilibrium stepping-stone or a range expansion
336 model might be more appropriate. However, our individuals (with the exception of Oman)
337 come from a restricted geographic area so that the spatial component of the demographic
338 model can be neglected. Indeed, we simulated from the posterior distribution of the FIM the
339 expected F_{st} between two demes at 14 independent microsatellites loci with the same
340 mutation rate as in Vignaud et al.²⁹. We found the F_{st} distribution to be compatible with the

341 value they found between East and West Australia, which are more distant than the samples
342 from our study area. This further suggests that our model is appropriate for our specific
343 dataset, but we are aware that it might not provide an accurate description of the *C.*
344 *melanopterus* demography over its whole range.

345 We made some simplifying assumptions when analysing the *C. melanopterus* data.
346 We used a single step change of effective population size in our indirect approach and a
347 constant Nm when fitting the FIM model. Obviously, the real evolutionary histories of
348 species are likely to be more complex and it may be challenging to distinguish between local
349 events (such as a bottleneck or expansion in a deme) and changes in connectivity/number of
350 demes in metapopulations^{16,41}. Climate change and anthropogenic activity are reported to
351 have impacted several marine species, including some sharks^{30,42-44}. A recent bottleneck in *C.*
352 *melanopterus* was identified in Moorea and attributed to human activity²⁹. In an effort to
353 distinguish between patterns that might be caused by metapopulation structure from those
354 caused by localized bottlenecks we subjected the data to analysis using both SID and SCD.
355 We did not find any evidence for a more recent event than the metapopulation expansion of
356 *C. melanopterus* (Supplementary Fig. S3), ruling out the possibility of a demographic change
357 related to the post-glacial maximum (~20-15 KYA) as has been proposed for other sharks³⁰.
358 The distribution of *C. melanopterus* is confined to reef flats and sheltered lagoons and the
359 effects of the last glacial maximum appear to have had little impact on this species.
360 However, recent and mild bottlenecks are more challenging to detect using coalescent
361 analyses^{45,46}. Indeed, many marine species do not show clear signatures of a bottleneck
362 despite being threatened by overfishing for many generations^{47,48}. Nevertheless, meta-
363 analysis studies have shown general reduction in diversity in overfished species⁴⁹. These

364 issues motivated us to investigate how best to detect a recent reduction in effective size in a
365 metapopulation. To this end, we performed a simulation study focusing on a
366 metapopulation with the same demography as that estimated in our *C. melanopterus* data but
367 experiencing a recent bottleneck. Obviously, a more extensive investigation of parameter
368 space would be interesting but is beyond the scope of this paper.

369 This simulation study had three goals: i) to better understand the recent
370 demographic history of *C. melanopterus*; ii) to characterise the behaviour of metapopulations
371 experiencing recent bottlenecks; iii) to evaluate the relative merits of our different sampling
372 schemes. Analyses of the simulated data paralleled those adopted for the real data, i.e.,
373 using the direct and indirect approaches. The indirect approach revealed two important
374 features: i) SCD fails to detect evidence of a bottleneck, even at high intensities, showing a
375 pattern that is similar to the no-bottleneck case (Fig. 4); ii) SID similarly does not show
376 evidence of bottleneck, but it shows a reduction in genetic diversity compared to the no-
377 bottleneck case at both intensities (100× and 1000×), within just 10 generations. The genetic
378 consequence of the bottleneck on a SID sample is the reduction in the estimated effective
379 size compared to the no-bottleneck scenario (Fig. 4, Supplementary Fig. S5 and S6).
380 However, we could not have detected such a decrease of the effective size without some
381 baseline reference values, suggesting that it is difficult to assess whether a species is
382 critically endangered when it is part of a structured metapopulation. This might explain
383 why many overfished species do not show a genetic signature of bottleneck despite being
384 threatened.

385 The direct approach confirmed that the two sampling schemes behave differently
386 after a bottleneck. We computed Nm using the FIM model and we found, as expected, that a

387 stable metapopulation shows a similar estimate of Nm under both sampling schemes (Fig. 5
388 and Supplementary Fig. S7). Conversely, a decrease in the Nm estimated from SID, but not
389 from SCD, is observed under a bottleneck scenario (Fig. 5 and Supplementary Fig. S7),
390 suggesting that the comparison of the two datasets may help detecting recent changes in
391 effective population size. These results (both the ABC-skyline models and the Nm estimates)
392 can be interpreted from a coalescent point of view: a recent bottleneck will produce a burst
393 of coalescent events only at the deme level (i.e., during the *scattering* phase), having no
394 impact at all on the *collecting* phase even for high intensity of reduction. A scattered sample,
395 which does not have the *scattering* phase, will not be affected by recent bottleneck events
396 while a local single deme sample will. These findings suggest that metapopulations are
397 more resilient than unstructured populations to recent changes in effective size (or
398 connectivity), implying that taking population structure into account is crucial when
399 inferring the demography of a species. Most species are structured, including those that are
400 endangered. In conservation studies we need to detect changes in effective population size
401 on the order of tens to, at most, hundreds of generations. Our method can be therefore be
402 usefully applied in a conservation setting as it is able to detect such recent changes. We
403 expect that it can be easily applied to both critically endangered (e.g. the daggernose shark)
404 and commercially exploited (e.g.. the Atlantic cod) species to better design appropriate
405 conservation management plans. We now envisage exploring the power of our approach
406 under a wider spectrum of parameters in order to provide guidelines for conservation
407 studies, but this is beyond the scope of this work. In the case of the *C. melanopterus* data
408 examined in the current study, the single deme sampled from Northern Australia has not
409 experienced a recent bottleneck. However, we cannot exclude the possibility that this may

410 have occurred for the other demes of the metapopulation, where conditions may have been
411 more extreme than those simulated here.

412

413 In summary, we present a recently developed target gene capture approach used
414 here for the first time to perform inferences at the intra-specific level. By exploiting
415 theoretical and simulation results on the coalescent history of metapopulations, we show
416 how few carefully collected specimens can provide information about the demography of a
417 species when many loci are available. This strategy can minimise sampling effort and
418 associated cost of data production, making population genomic analysis of non-model
419 organisms more feasible. We showed that this approach allows refined estimates of
420 connectivity among demes relative to traditional population genetics approaches and that
421 the sampling scheme adopted is particularly effective when investigating recent changes in
422 effective population size and, potentially, migration rates. This opens new avenues for
423 theoretical developments on the conservation management and recovery of endangered
424 species.

425

426 **Materials and Methods**

427 *Sampling, library preparation and sequencing*

428 A total of 18 samples of *C. melanopterus* were included from 8 different locations (Fig. 1):
429 Northern Territory and Queensland, Australia (N=11), Oman (N=1), Philippines (N=1),
430 Thailand (N=1), Malaysia (N=1), South Kalimantan, Indonesia (N=1), West Kalimantan,
431 Indonesian Borneo (N=1) and East Kalimantan, Indonesia (N=1). Samples were divided in
432 two datasets called “scatter” (SCD) and “single deme” (SID) with 9 and 11 samples

433 respectively (Table 1, Supplementary Table S1). The SCD dataset aims to mimic a “scatter
434 sample”, sensu Wakeley, where all samples come from independent geographical locations
435 distributed in close proximity to each other. We also included one sample from Oman in
436 order to infer the demographic history of the metapopulation at a larger scale and to cover a
437 more extensive area of the Indian Ocean. Muscle tissue was extracted from dead specimens
438 collected from local fish markets and then stored in 95% ethanol. Genomic DNA was
439 extracted using the E.Z.N.A Tissue DNA Kit (Omega Bio-Tek, Inc Norcross, GA) according
440 to the manufacturer’s instructions. Illumina sequencing libraries (500bp) were prepared and
441 amplified by PCR prior to two rounds of target gene capture (targeting 1077 autosomal
442 regions) using a ‘touchdown’ DNA hybridization approach²³. Resulting enriched libraries
443 were re-amplified to incorporate a sample specific index, pooled and sequenced paired-end
444 (250 bp reads) on an Illumina MiSeq benchtop sequencer (Illumina, Inc, San Diego, CA).
445 Sequence reads associated with each sample were identified and sorted by their respective
446 indices²³.

447

448 *De novo reference sequence assembly*

449 Briefly, sequence read data from several individuals was used to build a haploid reference
450 sequence for the 1077 target exons and associated introns. Adapters and low quality reads
451 were removed from the raw sequence read data and contigs were assembled *de novo*. On
452 target material was identified by determining the orthology of assembled contigs to a set of
453 core orthologs identified across model vertebrates (see SI for details). For each exon and
454 intron (5’ and 3’) the longest contig across all sequenced individuals was retained and used
455 to assemble a reference sequence that is 1,293,710 bp including exons and introns (5’ and 3’).

456

457 *Read mapping, variant calling and filtering*

458 A Burrows-Wheeler Aligner (BWA) ⁵⁰ was used to map reads to the reference sequence
459 generated by *de novo* assembly ⁵¹. Duplicate read marking was performed with Picard
460 v.1.129 (<http://broadinstitute.github.io/picard/>), local realignment with the Genome Analysis
461 ToolKit (GATK) v3.3-0 ⁵² and genotype calling with samtools v1.1 ⁵³. Raw variants were
462 filtered using custom Perl (v5.18.2) and R (3.0.2) scripts (Supplementary Table S8). Filtering
463 was based on strand bias, minimum depth 6×, removal of triallelic sites, removal of sites
464 with heterozygous calls in more than 80% of samples and removal of any missing calls in
465 order to have a dataset without any missing genotypes, reducing background noise and
466 uncertainty (see SI for details).

467

468 *Descriptive analyses*

469 Principal component analysis (PCA) was performed with the function “prcomp” in the R
470 environment ⁵⁴. A Gst distance matrix⁵⁵ was calculated with Arlequin 3.5⁵⁶ (see SI for details).

471

472 *Demographic Inference*

473 We developed an ABC ²⁴ framework to estimate parameters and compare demographic
474 models. The folded site frequency spectrum (SFS) and total number of SNPs were used as
475 summary statistics, to avoid phasing issues. Mutation and recombination rates were allowed
476 to vary across loci using hyperprior distributions (see SI for details). Generation time was set
477 to seven years ³⁶. Four demographic models were tested for both datasets (Fig. 2). Model
478 COS represents a constant-size population model described only by the effective population

479 size (N) (Fig. 2, panel a). Model CHG1 represents a single instantaneous demographic-
480 change model occurring at the time T_c (Fig. 2, panel b). N_{anc} and N_{mod} are the ancestral and
481 modern effective population sizes respectively. Values of the ancestral versus modern
482 effective population size ratio ($resize, N_{anc}/N_{mod}$) greater than 1 are compatible with a
483 reduction in effective population size while values lower than 1 suggest a population
484 expansion. If this ratio is equal to 1, the scenario is compatible with a constant-size
485 population. Model CHG2 represents a demographic-change model in which two events,
486 occurring at times T_{c1} and T_{c2} , cause changes in population size described by three effective
487 population sizes (N_{anc} , N_{int} and N_{mod} for ancestral, intermediate and modern respectively)
488 (Fig. 2, panel c). Model FIM represents a non-equilibrium finite island model with 100 demes
489 ($N_1 \dots N_{100}$) described by Nm as the product between the effective population size (N) and the
490 migration rate (m) (Fig. 2, panel d). N_{anc} is the effective size of the ancestral deme from which
491 the island originated at the time T_i . One hundred demes was chosen in order to approximate
492 the large number of demes that are necessary to describe the coalescent history in
493 metapopulation models in terms of the *scattering* and *collecting* phases²⁰.

494 We generated 100,000 simulations for each demographic model using fastsimcoal2
495 v2.5.1⁵⁷. Each simulation includes 995 and 998 gene genealogies for SCD and SID
496 respectively, to be consistent with our real data. Prior distributions are displayed in Table 3,
497 Supplementary Table S3 and S9. The posterior probabilities of each model were calculated
498 by a weighted multinomial logistic regression⁵⁸ for which we retained the best 25,000
499 simulations. The parameters of the best-fitting model were estimated from the 5,000
500 simulations closest to the observed dataset using a local linear regression according to²⁴.
501 Posterior distributions for model CHG1 and FIM are shown in Supplementary Fig. 8 and 9.

502 Both CHG1 and CHG2 models were used to graphically reconstruct the variation of
503 effective population size through time. To this end, for each combination of parameters
504 retained by the ABC algorithm (5,000 in our case), we recorded the effective size at specific
505 time points. The median value of the posterior distribution of the effective size at each time
506 point was calculated together with the 95% credible interval and plotted against time to
507 obtain an ABC-skyline reconstruction (Fig. 3, 4, Supplementary Fig. S3, S5 and S6). Time
508 points were defined by randomly extracting values from an exponential distribution with a
509 rate calibrated with the upper bound (97.5%) of the estimated mean TMRCA across loci. In
510 this way all generated time points were not greater than 97.5 % of the estimated mean
511 TMRCA across loci. Moreover, recent time points (0, 25, 50, 100, 200, 300, 400 and 500
512 generations ago) were manually added to increase the resolution towards recent events.
513 Analyses were performed in the R environment⁵⁴ with the library *abc*⁵⁹.

514 We performed cross-validation for both model selection and parameter estimation by
515 randomly generating *pods* from the prior distributions under each model. For each cross-
516 validation experiment we generated 1000 *pods* and we applied the same inferential
517 procedure as for the observed data (see SI, Supplementary Fig. S10 and S11). In addition, a
518 posterior predictive test⁶⁰ was carried out to test whether the data can be reproduced under
519 a specific demographic model. Bayesian p-values computed from the posterior distribution
520 of the number of polymorphic sites showed that none of the four models in both datasets
521 could be rejected (SI, Supplementary Fig. S12).

522

523 *Recent Bottleneck*

524 We tested for signatures of a recent population bottleneck in the metapopulation data set for
525 both sampling schemes. We restricted our ABC analyses to a metapopulation with
526 demographic parameters consistent with those estimated from the empirical data examined
527 in this study. We simulated *pods* under a modified FIM model (FIM-BOTT) characterised by
528 two additional parameters (Supplementary Fig. S4, panel b): T_{bott} (the time when an
529 instantaneous decrease in Nm occurred) and I_{bott} (the ratio of the ancestral Nm to the new
530 one). When $I_{\text{bott}}=1$ FIM-BOTT reduces to FIM (Fig. 2, panel d). We tested twelve
531 combinations of parameters, namely two I_{bott} ($100\times$ and $1000\times$) and six T_{bott} (10, 50, 100, 200,
532 500 and 1000 generations ago). For each combination of parameters, we simulated 1000 *pods*
533 for both SCD and SID (Supplementary Fig. S4, panel b). To further understand the effect of
534 metapopulation structure on bottlenecks, we performed additional simulations of an
535 unstructured population under a modified CHG1 model (CHG1-BOTT, Supplementary Fig.
536 S4, panel a). In this case, *pods* were simulated with demographic parameters estimated in our
537 real SCD data under model CHG1, to which we added a recent bottleneck characterised by
538 I_{bott} and T_{bott} . Note that here I_{bott} represents the ratio between the ancestral effective
539 population size N_e and the new one (while in FIM-BOTT it represents the ratio of the
540 ancestral Nm to the new one). We tested the same twelve parameters combinations of I_{bott}
541 and T_{bott} that we previously used when simulating FIM-BOTT. CHG1-BOTT reduces to
542 CHG1 when $I_{\text{bott}}=1$. All *pods* simulated under both FIM-BOTT (for both sampling schemes)
543 and CHG1-BOTT were analysed using the ABC-skyline produced by the CHG2 model
544 (Fig.2, panel c). We applied the same settings used for the real data and for each
545 combination of parameters we plotted the average of the N_e through time across the 1,000
546 *pods* (see SI). All bottleneck scenarios were compared to the ABC-skyline reconstructed from

547 *pods* simulated under FIM or CHG1. Finally, *pods* simulated under FIM-BOTT were also
548 analysed using FIM to estimate Nm . We compared these values to those estimated from *pods*
549 generated under FIM.

550

551 References

- 552 1 Nadachowska-Brzyska, K. *et al.* Demographic divergence history of pied flycatcher and
553 collared flycatcher inferred from whole-genome re-sequencing data. *PLoS Genet* **9**,
554 e1003942, (2013).
- 555 2 Romiguier, J. *et al.* Comparative population genomics in animals uncovers the determinants
556 of genetic diversity. *Nature* **515**, 261-263, (2014).
- 557 3 Emerson, K. J. *et al.* Resolving postglacial phylogeography using high-throughput sequencing.
558 *Proc Natl Acad Sci U S A* **107**, 16196-16200, (2010).
- 559 4 Li, S. & Jakobsson, M. Estimating demographic parameters from large-scale population
560 genomic data using Approximate Bayesian Computation. *BMC Genet* **13**, 22, (2012).
- 561 5 Roux, J., Rosikiewicz, M. & Robinson-Rechavi, M. What to compare and how: Comparative
562 transcriptomics for Evo-Devo. *Journal of Experimental Zoology Part B: Molecular and*
563 *Developmental Evolution*, (2015).
- 564 6 Arnold, B., Corbett-Detig, R. B., Hartl, D. & Bomblies, K. RADseq underestimates diversity and
565 introduces genealogical biases due to nonrandom haplotype sampling. *Mol Ecol* **22**, 3179-
566 3190, (2013).
- 567 7 Huang, H. & Knowles, L. L. Unforeseen Consequences of Excluding Missing Data from Next-
568 Generation Sequences: Simulation Study of RAD Sequences. *Syst Biol*, (2014).
- 569 8 Leache, A. D. *et al.* Phylogenomics of Phrynosomatid Lizards: Conflicting Signals from
570 Sequence Capture versus Restriction Site Associated DNA Sequencing. *Genome Biol Evol* **7**,
571 706-719, (2015).
- 572 9 Jones, M. R. & Good, J. M. Targeted capture in evolutionary and ecological genomics. *Mol*
573 *Ecol*, (2015).
- 574 10 Eriksson, A. & Manica, A. The doubly conditioned frequency spectrum does not distinguish
575 between ancient population structure and hybridization. *Mol Biol Evol* **31**, 1618-1621,
576 (2014).
- 577 11 Peter, B. M., Wegmann, D. & Excoffier, L. Distinguishing between population bottleneck and
578 population subdivision by a Bayesian model choice procedure. *Mol Ecol* **19**, 4648-4660,
579 (2010).
- 580 12 Heller, R., Chikhi, L. & Siegmund, H. R. The confounding effect of population structure on
581 Bayesian skyline plot inferences of demographic history. *PLoS One* **8**, e62992, (2013).
- 582 13 Chikhi, L., Sousa, V. C., Luisi, P., Goossens, B. & Beaumont, M. A. The confounding effects of
583 population structure, genetic diversity and the sampling scheme on the detection and
584 quantification of population size changes. *Genetics* **186**, 983-995, (2010).
- 585 14 Currat, M. & Excoffier, L. Strong reproductive isolation between humans and Neanderthals
586 inferred from observed patterns of introgression. *Proc Natl Acad Sci U S A* **108**, 15129-
587 15134, (2011).
- 588 15 Francois, O., Blum, M. G., Jakobsson, M. & Rosenberg, N. A. Demographic history of
589 european populations of *Arabidopsis thaliana*. *PLoS Genet* **4**, e1000075, (2008).

- 590 16 Wakeley, J. Nonequilibrium migration in human history. *Genetics* **153**, 1863-1871, (1999).
- 591 17 Wakeley, J. The coalescent in an island model of population subdivision with variation
592 among demes. *Theor Popul Biol* **59**, 133-144, (2001).
- 593 18 Stadler, T., Haubold, B., Merino, C., Stephan, W. & Pfaffelhuber, P. The impact of sampling
594 schemes on the site frequency spectrum in nonequilibrium subdivided populations. *Genetics*
595 **182**, 205-216, (2009).
- 596 19 Ray, N., Currat, M. & Excoffier, L. Intra-deme molecular diversity in spatially expanding
597 populations. *Mol Biol Evol* **20**, 76-86, (2003).
- 598 20 Wakeley, J. Segregating sites in Wright's island model. *Theor Popul Biol* **53**, 166-174, (1998).
- 599 21 Papastamatiou, Y. P., Lowe, C. G., Caselle, J. E. & Friedlander, A. M. Scale-dependent effects
600 of habitat on movements and path structure of reef sharks at a predator-dominated atoll.
601 *Ecology* **90**, 996-1008, (2009).
- 602 22 Papastamatiou, Y. P., Friedlander, A. M., Caselle, J. E. & Lowe, C. G. Long-term movement
603 patterns and trophic ecology of blacktip reef sharks (*Carcharhinus melanopterus*) at Palmyra
604 Atoll. *J Exp Mar Biol Ecol* **386**, 94-102, (2010).
- 605 23 Li, C., Hofreiter, M., Straube, N., Corrigan, S. & Naylor, G. J. Capturing protein-coding genes
606 across highly divergent species. *Biotechniques* **54**, 321-326, (2013).
- 607 24 Beaumont, M. A., Zhang, W. & Balding, D. J. Approximate Bayesian computation in
608 population genetics. *Genetics* **162**, 2025-2035, (2002).
- 609 25 Bertorelle, G., Benazzo, A. & Mona, S. ABC as a flexible framework to estimate demography
610 over space and time: some cons, many pros. *Mol Ecol* **19**, 2609-2625, (2010).
- 611 26 Heupel, M. *Carcharhinus melanopterus*, <www.iucnredlist.org> (2009).
- 612 27 Field, I. C. *et al.* Changes in size distributions of commercially exploited sharks over 25 years
613 in northern Australia using a Bayesian approach. *Fish Res* **125**, 262-271, (2012).
- 614 28 Henderson, A., Al-Oufi, H. & McIlwain, J. Survey, status and utilization of the elasmobranch
615 fishery resources of the Sultanate of Oman. (Sultan Qaboos University, Muscat, 2007).
- 616 29 Vignaud, T. M. *et al.* Blacktip reef sharks, *Carcharhinus melanopterus*, have high genetic
617 structure and varying demographic histories in their Indo-Pacific range. *Mol Ecol* **23**, 5193-
618 5207, (2014).
- 619 30 Portnoy, D. S. *et al.* Contemporary population structure and post-glacial genetic demography
620 in a migratory marine species, the blacknose shark, *Carcharhinus acronotus*. *Mol Ecol* **23**,
621 5480-5495, (2014).
- 622 31 Marko, P. B. *et al.* The 'Expansion-Contraction' model of Pleistocene biogeography: rocky
623 shores suffer a sea change? *Mol Ecol* **19**, 146-169, (2010).
- 624 32 Mona, S., Ray, N., Arenas, M. & Excoffier, L. Genetic consequences of habitat fragmentation
625 during a range expansion. *Heredity (Edinb)* **112**, 291-299, (2014).
- 626 33 Wegmann, D., Currat, M. & Excoffier, L. Molecular diversity after a range expansion in
627 heterogeneous environments. *Genetics* **174**, 2009-2020, (2006).
- 628 34 Li, C. *et al.* DNA capture reveals transoceanic gene flow in endangered river sharks. *Proc Natl*
629 *Acad Sci U S A* **112**, 13302-13307, (2015).
- 630 35 Mourier, J. & Planes, S. Direct genetic evidence for reproductive philopatry and associated
631 fine-scale migrations in female blacktip reef sharks (*Carcharhinus melanopterus*) in French
632 Polynesia. *Mol Ecol* **22**, 201-214, (2013).
- 633 36 Mourier, J., Mills, S. C. & Planes, S. Population structure, spatial distribution and life-history
634 traits of blacktip reef sharks *Carcharhinus melanopterus*. *J Fish Biol* **82**, 979-993, (2013).
- 635 37 Vignaud, T., Clua, E., Mourier, J., Maynard, J. & Planes, S. Microsatellite analyses of blacktip
636 reef sharks (*Carcharhinus melanopterus*) in a fragmented environment show structured
637 clusters. *PLoS One* **8**, e61067, (2013).

- 638 38 Chin, A., Tobin, A. J., Heupel, M. R. & Simpfendorfer, C. A. Population structure and
639 residency patterns of the blacktip reef shark *Carcharhinus melanopterus* in turbid coastal
640 environments. *J Fish Biol* **82**, 1192-1210, (2013).
- 641 39 Speed, C. W. *et al.* Spatial and temporal movement patterns of a multi-species coastal reef
642 shark aggregation. *Mar Ecol Prog Ser* **429**, 261-U618, (2011).
- 643 40 Wilkins, J. F. A separation-of-timescales approach to the coalescent in a continuous
644 population. *Genetics* **168**, 2227-2244, (2004).
- 645 41 Mazet, O., Rodriguez, W. & Chikhi, L. Demographic inference using genetic data from a single
646 individual: Separating population size variation from population structure. *Theor Popul Biol*
647 **104**, 46-58, (2015).
- 648 42 Hernández, S. *et al.* Demographic history and the South Pacific dispersal barrier for school
649 shark (*Galeorhinus galeus*) inferred by mitochondrial DNA and microsatellite DNA mark. *Fish*
650 *Res* **167**, 132-142, (2015).
- 651 43 Vignaud, T. M. *et al.* Genetic structure of populations of whale sharks among ocean basins
652 and evidence for their historic rise and recent decline. *Mol Ecol* **23**, 2590-2601, (2014).
- 653 44 O'Leary, S. J. *et al.* Genetic Diversity of White Sharks, *Carcharodon carcharias*, in the
654 Northwest Atlantic and Southern Africa. *Journal of Heredity*, (2015).
- 655 45 Girod, C., Vitalis, R., Leblois, R. & Freville, H. Inferring population decline and expansion from
656 microsatellite data: a simulation-based evaluation of the Msvar method. *Genetics* **188**, 165-
657 179, (2011).
- 658 46 Roman, J. & Palumbi, S. R. Whales before whaling in the North Atlantic. *Science* **301**, 508-
659 510, (2003).
- 660 47 Riccioni, G. *et al.* Spatio-temporal population structuring and genetic diversity retention in
661 depleted Atlantic Bluefin tuna of the Mediterranean Sea. *P Natl Acad Sci USA* **107**, 2102-
662 2107, (2010).
- 663 48 Marra, A., Mona, S., Sa, R. M., D'Onghia, G. & Maiorano, P. Population Genetic History of
664 *Aristeus antennatus* (Crustacea: Decapoda) in the Western and Central Mediterranean Sea.
665 *PLoS One* **10**, e0117272, (2015).
- 666 49 Pinsky, M. L. & Palumbi, S. R. Meta-analysis reveals lower genetic diversity in overfished
667 populations. *Mol Ecol* **23**, 29-39, (2014).
- 668 50 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
669 *Bioinformatics* **25**, 1754-1760, (2009).
- 670 51 Birol, I. *et al.* De novo transcriptome assembly with ABySS. *Bioinformatics* **25**, 2872-2877,
671 (2009).
- 672 52 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
673 next-generation DNA sequencing data. *Genome Res* **20**, 1297-1303, (2010).
- 674 53 Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and
675 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987-
676 2993, (2011).
- 677 54 R: A Language and Environment for Statistical Computing (R Foundation for Statistical
678 Computing, 2014).
- 679 55 Nei, M. Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci U S A* **70**,
680 3321-3323, (1973).
- 681 56 Excoffier, L. & Lischer, H. E. Arlequin suite ver 3.5: a new series of programs to perform
682 population genetics analyses under Linux and Windows. *Mol Ecol Resour* **10**, 564-567,
683 (2010).
- 684 57 Excoffier, L., Dupanloup, I., Huerta-Sanchez, E., Sousa, V. C. & Foll, M. Robust demographic
685 inference from genomic and SNP data. *PLoS Genet* **9**, e1003905, (2013).
- 686 58 Beaumont, M. A. in *Simulation, Genetics, and Human Prehistory* 135-154 (McDonald
687 Institute for Archaeological Research, 2008).

688 59 Csillery, K., Francois, O. & Blum, M. G. B. abc: an R package for approximate Bayesian
689 computation (ABC). *Methods Ecol Evol* **3**, 475-479, (2012).
690 60 Gelman, A., Carlin, J., Stern, H. & Rubin, D. *Bayesian Data Analysis*. (CRC Press, 2004).

691

692 **Acknowledgments**

693 The authors thank Giorgio Bertorelle and Oscar Lao for helpful comments and suggestions
694 and Laurent Excoffier for sharing fastsimcoal v2.5.2 before publication. Data analysis was
695 performed on the computer cluster Genotoul bioinformatics platform
696 (www.bioinfo.genotoul.fr). This work was supported by the Agence Nationale de la
697 Recherche (ANR-12-BSV7-0012 to P.M.D.) and by the US National Science Foundation (DEB-
698 1132229).

699

700 **Author Contributions**

701 S.M., S.P. and G.N. designed research; M.H, C.L. and S.C. contributed with new reagents
702 and analytical tools; P.M.D., S.C., S.M. and G.N. performed research; P.M.D. and S.M.
703 analysed data; M.V. gave comments on the manuscript and project; P.M.D., S.C., S.M and
704 G.N. wrote the paper.

705

706 **Competing interests**

707 The authors declare no competing financial interests.

708

709 **Data accessibility**

710 A vcf file containing all variants used in this study is being submitted to the Dryad database

711 (<http://datadryad.org/>).

712

713 **Fig. 1: Sampling locations with sample size.** Samples in blue were included in the scatter
714 dataset (SCD, 9 samples). Samples in orange were considered for the single deme dataset
715 (SID, 11 samples). Map of sampling locations was generated with the library “rworldmap”
716 with R software (v3.0.2, <https://cran.r-project.org/>)⁵⁴.

717 **Fig. 2: Demographic models tested for both sampling schemes.** a) COS, constant-size
718 model; b) CHG1, demographic-change model (one demographic change); c) CHG2,
719 demographic-change model (two demographic changes); d) FIM, non-equilibrium finite
720 island model. N_{mod} : modern effective population size; N_{anc} : ancestral effective population
721 size; T_c , T_{c1} and T_{c2} : time of the demographic change (in generations); T_i : time of the onset of
722 the island (in generations).

723 **Fig. 3: Skyline reconstruction of the effective population size through time.** a) single deme
724 (SID, in orange); b) scatter (SCD, in blue). Both skylines were reconstructed up to the mode
725 of the mean TMRCA across loci. Median values are shown (darker lines) with the 95% high
726 posterior density interval (lighter lines).

727 **Fig. 4: Average skyline reconstruction under model CHG2 of pseudo-observed data sets**
728 **(pods) simulated with model FIM and FIM-BOTT in a metapopulation (for SCD and SID)**
729 **and in an unstructured population (model CHG1 and CHG1-BOTT).** Data were simulated
730 with $I_{\text{bott}}=1000$ at: a) $T_{\text{bott}}=10$ generations; b) $T_{\text{bott}}=50$ generations; c) $T_{\text{bott}}=100$ generations
731 before present. An average skyline reconstruction is shown across 1000 simulations. Solid
732 lines: scenario with bottleneck (model FIM-BOTT and CHG1-BOTT); dashed line: scenario
733 without bottleneck (model FIM and CHG1). Median values are shown (darker lines) with
734 the 95% high posterior density interval (lighter lines).

735 **Fig. 5: Distribution of the median of Nm estimated under model FIM from pseudo-**
736 **observed data sets (pods) generated with FIM-BOTT and FIM.** Data were simulated for
737 both sampling schemes with $I_{\text{bott}}=100$ and various T_{bott} (10, 50, 100, 200, 500 and 1000
738 generations ago). Dotted lines represent the value of Nm used to simulate pods under FIM
739 model.

740

741

742

743

744

745 **Table 1: summary of the “scatter” (SCD) and “single deme” (SID) dataset.**

Dataset	No. of samples	No. of regions	Bp sequenced	No. of SNPs	SNP density region (kb)
SCD	9	995	606,647	2605	4.3
SID	11	998	632,160	1946	3.1

746

747 **Table 2: models posterior probability calculated as in ⁵⁸ using the closest 25,000**
748 **simulations. SCD: scatter dataset; SID: single deme dataset.**

749

Dataset	COS	CHG1	FIM
SCD	0.00	0.40	0.60
SID	0.24	0.31	0.45

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773 **Table 3: parameter estimation under the finite island model (FIM).** N_{anc} : effective
 774 population size of the ancestral deme; T_i : time of the onset of the island (in generations); Nm :
 775 product of the effective population size N and the migration rate m for each deme; SCD:
 776 scatter dataset (9 individuals); SID: single deme dataset (11 individuals); pooled dataset:
 777 SCD+SID without shared samples (18 individuals).
 778

Model FIM	Median	Mode	0.025 ^a	0.975 ^a	Prior ^b
SCD					
N_{anc}	33,567	42,704	1604	62,242	U: 100-100,000
T_i	72,760	56,354	24,124	141,628	U: 100-200,000
Nm	52.85	47.17	38.80	109.77	U: 0.05-250
SID					
N_{anc}	31,016	30,737	1571	86,026	U: 100-100,000
T_i	63,431	61,671	6095	163,634	U: 100-200,000
Nm	42.04	38.61	28.51	117.80	U: 0.05-250
Pooled dataset					
N_{anc}	52,541	58,148	8176	84,492	U: 100-100,000
T_i	65,004	18,329	4255	176,692	U: 100-200,000
Nm	58.7	51.4	38.3	150.9	U: 0.05-250

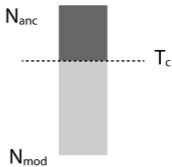
779 ^aUpper and lower limits of the 95% credible interval about the estimated mode.

780 ^bU, uniform probability, in the range of the two values.

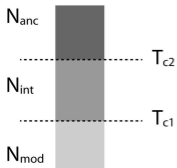
a) Model COS



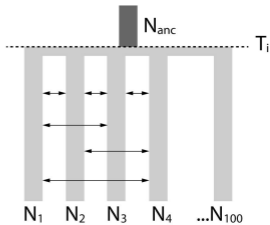
b) Model CHG1



c) Model CHG2

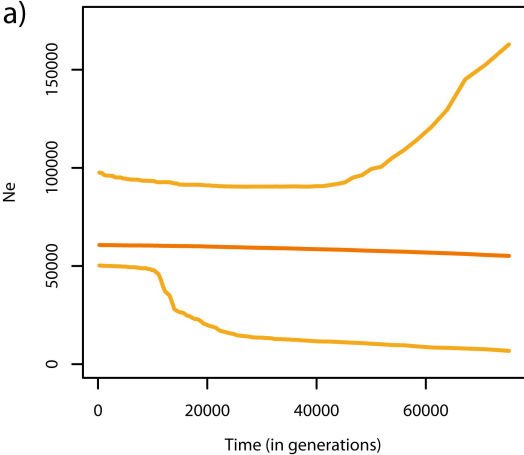


d) Model FIM



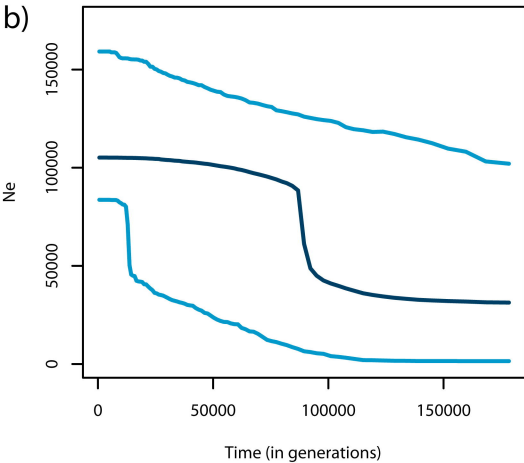
Deme - SID

a)



Scatter - SCD

b)



Metapopulation

Unstructured population

Scatter - SCD

Deme - SID

