1    Title: **Post-glacial expansion dynamics, not refugial isolation, shaped the genetic structure of a**

2    **migratory bird, the Yellow Warbler (*Setophaga petechia*).**

3    **Authors:** Eleanor F. Miller* [a], Michela Leonardi [a], Robert Beyer [a], Mario Krapp [a], Marius Somveille [a, b,]

4    [c], Gian Luigi Somma [a], Pierpaolo Maisano Delser [a]**, Andrea Manica [a]**

5    [a] Department of Zoology, University of Cambridge, Downing Street, Cambridge, CB2 3EJ, UK.

6    [b] Department of Biology, Colorado State University, Fort Collins, CO, 80521 USA

7    [c] Centre for Biodiversity and Environment Research, Department of Genetics, Evolution and

8    Environment, University College London, London WC1E 6BT, UK

9    ** Equal contribution

10    * Corresponding Author: Eleanor F. Miller, Department of Zoology, University of Cambridge,

11    Downing Street, Cambridge, CB2 3EJ, UK. Email: em618@cam.ac.uk

12    ORCIDs:

13    E.F.M.  0000-0002-3213-5714

14    M.L.    0000-0001-8933-9374

15    M.S.    0000-0002-6868-5080

16    G.L.S.   0000-0003-0383-5719

17    M.K.    0000-0002-2599-0683

18    A.M.    0000-0003-1895-450X

19    P.M.D.  0000-0002-1844-1715

20      **Abstract**

21      During the glacial periods of the Pleistocene, swathes of the Northern Hemisphere were covered by

22      ice sheets, tundra and permafrost leaving large areas uninhabitable for temperate and boreal

23      species.  The glacial refugia paradigm proposes that, during glaciations, species living in the Northern

24      Hemisphere were forced southwards, forming isolated, insular populations that persisted in disjunct

25      regions known as refugia.  According to this hypothesis, as ice sheets retreated, species recolonised

26      the continent from these glacial refugia, and the mixing of these lineages is responsible for modern

27      patterns of genetic diversity.  However, an alternative hypothesis is that complex genetic patterns

28      could also arise simply from heterogenous post-glacial expansion dynamics, without separate

29      refugia. Both mitochondrial and genomic data from the North American Yellow warbler (*Setophaga*

30      *petechia)* shows the presence of an eastern and western clade, a pattern often ascribed to the

31      presence of two refugia. Using a climate-informed spatial genetic modelling (CISGeM) framework,

32      we were able to reconstruct past population sizes, range expansions, and likely recolonisation

33      dynamics of this species, generating spatially and temporally explicit demographic reconstructions.

34      The model captures the empirical genetic structure despite including only a single, large glacial

35      refugium. The contemporary population structure observed in the data was generated during the

36      expansion dynamics after the glaciation and is due to unbalanced rates of northward advance to the

37      east and west linked to the melting of the icesheets. Thus, modern population structure in this

38      species is consistent with expansion dynamics, and refugial isolation is not required to explain it,

39      highlighting the importance of explicitly testing drivers of geographic structure.

## Introduction

40

41     It has frequently been shown that seemingly continuously distributed populations in the Northern

42     Hemisphere harbour geographic structure in their genetic diversity. Indeed, within North America,

43     many widespread and migratory passerines exhibit clear differences in both migration patterns and

44     genomic diversity between eastern and western populations e.g. (1–3). This pattern has been

45     interpreted as the consequence of glaciations, during which species were forced southwards,

46     forming isolated, insular populations that persisted in disjunct regions known as refugia (4,5).

47     According to this narrative, as ice-sheets retreated, species recolonised the continent from these

48     glacial refugia, and the subsequent mixing of these lineages is responsible for modern patterns of

49     genetic diversity.

50     However, even though the cycles of expansion and contraction could have fragmented ranges,

51     leading to multiple glacial refugia in some species, multiple glacial refugia have not been

52     demonstrated for all species e.g. (6,7). Indeed, it is becoming clear that glaciations in North America

53     might not have driven range fragmentation as ubiquitously as it has previously been assumed, e.g.

54     (8,9).

55     What other processes might then have shaped the genetics of modern populations? Range

56     expansions have been shown to have the potential to leave profound signatures in the genetic

57     structure of metapopulations through repeated founder events (10). An extreme consequence of

58     this process is gene surfing, when rare variants can become common through stochastic sampling

59     during a founder event, and then be spread widely at high frequency during the subsequent

60     expansion. An important role for the recolonization dynamics in shaping modern-day population

61     structuring has been recently put forward for a trans-continentally distributed species, the painted-

62     turtle, *Chrysemys picta* (11). Reid et al. (11) demonstrated that, for this species, genetic

63     differentiation during range expansion and isolation-by-distance are more likely to have driven

64     modern-day population diversity than isolation in allopatric refugia.

65    The difficulty in quantifying the role of the range change dynamics on the genetic structure of

66    species is that the results are highly dependent on the detail of the dynamics. Whilst it is

67    straightforward to build simple spatial models that represent a range expansion, capturing the

68    spatial and temporal heterogeneities of the real process is challenging. A possible solution is to use

69    climate informed spatial genetic models (CISGeMs), which use climate reconstructions to condition

70    the local demography of individual populations within a map, and quantify the demographic

71    parameters by Approximate Bayesian Computation comparing predicted and observed genetic

72    quantities (see Fig. 1). This approach has been successfully used to reconstruct the dynamics of the

73    out of Africa expansion of humans (12,13).

74    In this paper, we use the CISGeM framework to explore the past range dynamics of the Yellow

75    Warbler. This species is an abundant passerine species with a large continuous contemporary range

76    and clear geographic population structuring, for which range-wide genomic data are available (14).

77    Here we test whether today's patterns of genetic structure in the North American Yellow Warbler

78    (*Setophaga petechia*) can be best explained by recolonization from isolated glacial refugia, or if,

79    more simply, heterogenous post-glacial expansion dynamics, without separate refugia, may have

80    been enough to result in observed patterns today. Firstly, we describe genetic patterns that are

81    found in the Yellow Warbler today from an empirical RAD-seq dataset.  Then, we fit a spatially

82    explicit model of population growth and expansion that accounts for past climatic variation to the

83    dataset.  By simulating the genetics and fitting to the observations with an Approximate Bayesian

84    Computation framework, we investigate to what extent these recolonization dynamics could explain
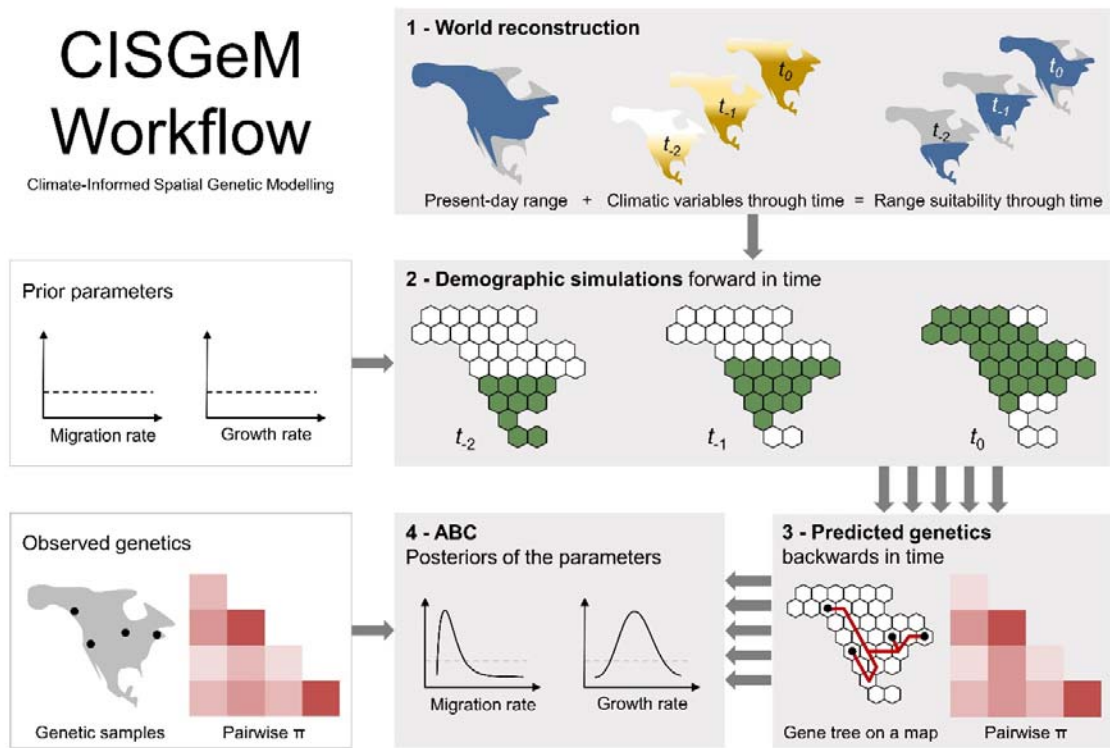
85    modern genomic patterns.

**Figure 1.** A schematic representation of the Climate-Informed Spatial Genetic Modelling framework implemented in this paper.

## Results

**Observed genetics:**

The 200 samples included in our study came from 21 sites across the modern breeding range of the North American Yellow Warbler.  Sample sizes per site ranged from 6 to 20 individuals (see Materials and Methods). Analysis of the genetic structure (15) of the Yellow Warbler population revealed a clear longitudinal divide, with distinct East and West clusters that converge in the centre of the continent.  Populations were with a mixture proportion of less than 70% for either of the two clusters ('East' and 'West') were grouped in the 'Central' category. This pattern is congruent with both the distribution of mitochondrial haplotypes (16) and patterns of migratory connectivity (17) in this species.
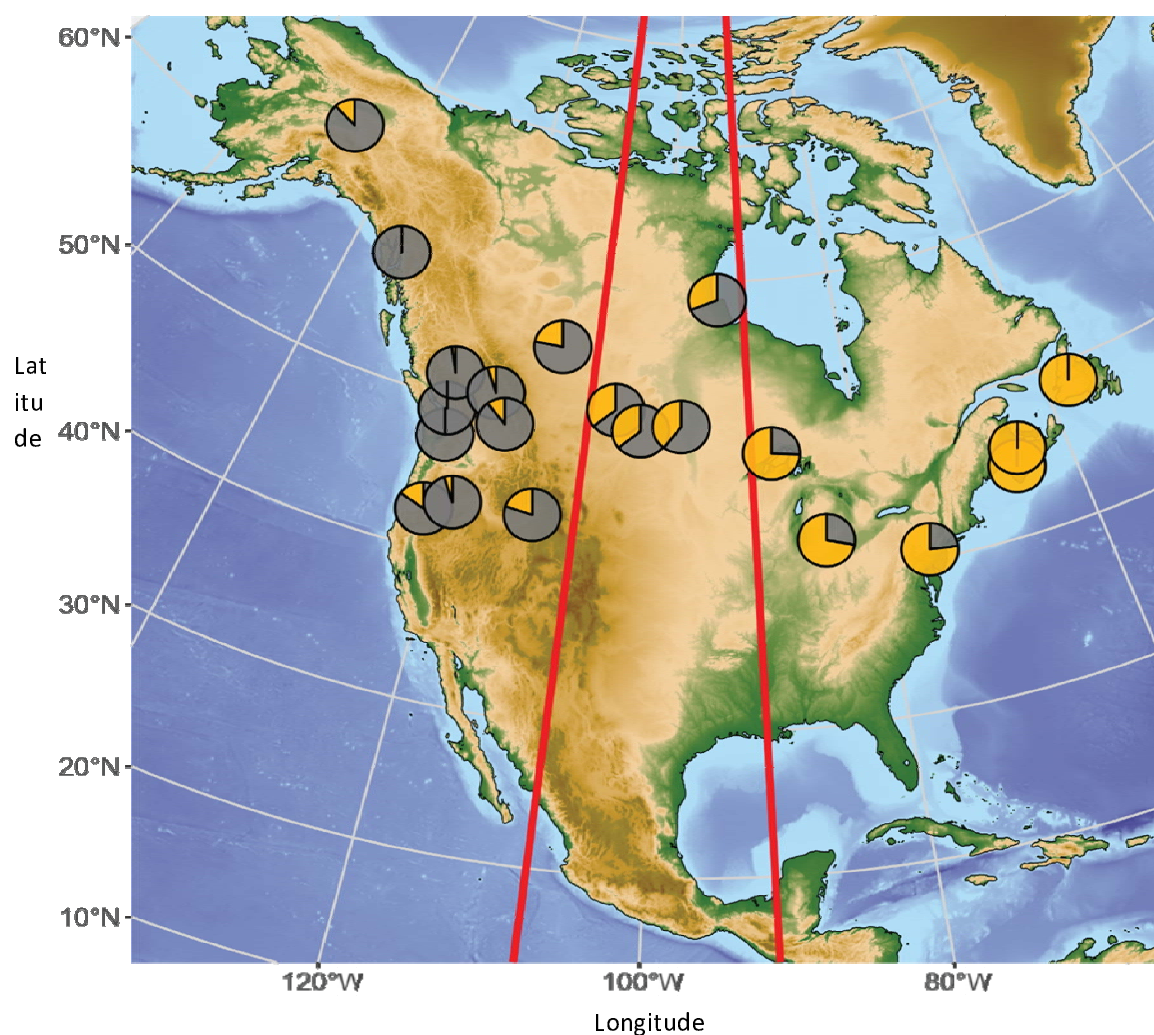
**Figure 2.** Genetic clustering results ($K$ = 2) results for all 21 populations in our study. Red lines separate the populations longitudinally into West, Central, and East.

98 **Species Distribution Modelling for world reconstruction:**

99 CISGeM requires the reconstruction of range suitability maps through time (step 1 in Fig. 1). We built

100 a Species Distribution Model (18) for Yellow Warblers based on modern data, and projected back in

101 time using paleoclimate reconstructions. The raw species occurrences data to define the present-day

102 range, downloaded from the Global Biodiversity Information facility (GBIF, our data can be found at

103 https://doi.org/10.15468/dl.jfkwcg), totalled 1,573,147 data points. After filtering for coordinate

104 accuracy, allowing an attributed error of 1km maximum, and filtering to only include points found

105 within the BirdLife breeding and resident geographical ranges (BirdLife International and Handbook

106 of the Birds of the World 2018), we were left with 177,202 data points. As SDM works on

107 presence/absence data and not frequencies, we retained only one presence per 0.5° grid cell,

108 further refining this dataset down to 3,364 observations. With these observations we selected the

109 four most informative, uncorrelated (threshold=0.7), bioclimatic variables to base our model on:

110 Leaf Area Index (LAI), BIO7 (Temperature Annual Range), BIO8 (Mean Temperature of Wettest

111 Quarter), BIO14 (Precipitation of Driest Month). Observations were further thinned based on a

112 maximum distance between points of 70 km, leaving 1,188 presences; this procedure is used to

113 correct for uneven sampling biases (Steen et al. 2020). We fitted SDMs to predict the probability of

114 occurrence in each grid cell. We used an ensemble of four different algorithms: generalised linear

115 models (GLM, (19)), generalized boosting method (GBM,(20)), generalised additive models (GAM,

116 (21)), and random forest (22). Models were run performing spatial cross validation with 80% of the

117 data used to train the algorithm and the remaining 20% to test it.

118 At present, the predicted potential distribution matches well the best range estimates for the

119 species (Supplementary Fig. 1.). Based on paleoclimate and vegetation reconstruction (see Materials

120    and Methods), range projections from the present day back to 50 thousand years ago suggest that

121    the distribution of habitat suitable for the Yellow Warbler expanded and contracted, to various

122    degrees, multiple times.  The potential range underwent a substantial contraction into the south of

123    the continent at the peak of the Last Glacial Maximum (LGM), ~21kya, before beginning to re-

124    expand (Supplementary Fig. 2. A-B), but the range never separated into distinct eastern and western

125    refugia. Following the LGM, the retreat of the Cordillerian Ice Sheet was asymmetric: in the west, the

126    ice started to retreat at about 18ka (23) with the opening of a corridor that progressively expanded

127    to the higher latitudes, whereas the eastern and central part of the Laurentide Ice Sheet began

128    retreating much later (24). By 13kya this deglaciated terrain became habitable for Yellow Warblers

129    according to our SDM (Supplementary Fig. 2. C).  From then on, as the ice sheets retreated further,

130    habitat to the east of the continent and in the central area deglaciated, becoming increasingly viable

131    (Supplementary Fig. 2. D-F).

132    **Climate Informed Spatial Genetic Model**

133    The reconstructed range suitability maps over time were used as an input for CISGeM.  In this

134    framework, the genetics of multiple populations can be modelled within a spatially explicit

135    reconstruction of the world where the suitability of each deme changes through time according to

136    the SDM back-cast suitability scores (Fig. 1). Using an Approximate Bayesian Computation

137    framework, we fitted basic demographic parameters such as population growth rate and migration,

138    as well as the link between SDM suitability scores and local population sizes.  The mean pairwise

139    genetic differentiation ($\pi$) between populations in each of the three clades (East, Central and West)

140    were used as summary statistics that had to be matched by the model. We performed a Monte-

141    Carlo sweep of the input parameters (Table 1.), generating a total of 61,504 simulations.

142    Visual inspection of the values of pairwise differentiation among the clades revealed that the model

143    was able to recapitulate the observations in a realistic fashion (see Supplementary Fig. 3. A for a PCA

144    plot of the values predicted by the model vs the observations, and Supplementary Fig. 4. for all

145    individual summary statistics). We then formally tested model fit with '*gfit*' from the '*abc*' package

146    (25) in R (see Materials and Methods for details). This test verifies that the distance between the

147    observed and the simulated data is not significantly larger than the distance of a random simulation

148    to other simulations (and thus that the model is able to capture the patterns seen in the data): our

149    model recovered a p value of 0.379 which implies a good fit (Supplementary Fig. 3. B).

150    We used a random forest algorithm (ABC-RF) (26) to generate posterior probabilities of the input

151    demographic parameters given the observed levels of pairwise population differentiation

152    (Supplementary Fig. 5.). The metapopulation dynamics was characterised by an expansion dynamics

153    with moderate to strong bottlenecks (as determined by a relatively low directed expansion

154    coefficient, $m_d$, which defines the proportion of individuals that move into an unoccupied area,

155    Supplementary Fig. 5. C), followed by limited subsequent migration (low values in the undirected

156    expansion coefficient rates $m_r$, Supplementary Fig. 5. E, in accordance with observations that this

157    species tends to be philopatric in its breeding range). Such signals suggest an expansion

158    characterised by sequential founder events that would have set up a pattern of isolation by distance

159    along the colonisation routes that was preserved by the limited migration afterwards.

160    From the top 2.5% best fitting simulations (n=1055 runs), we reconstructed the demography of the

161    species through space and time. The average demographic profile, calculated as a weighted mean of

162    population size across these simulations, shows that the Yellow Warbler was forced to contract its

163    range at the peak of the Last Glacial Maximum (~21kya) as the ice sheets grew across the north of

164    the continent (Fig. 3. A&B). At this point, the population existed in a restricted but broadly

165    continuous range in the south of the continent. As the climate ameliorated, northward range

166    expansion became possible. However, the pattern of recolonization was uneven. By 13kya, our

167    model reconstructs an expansion mostly following the corridor that opened between the Laurentide

168    and Cordilleran icesheets on the west of the continent, whilst expansion on the eastern side was

169    limited (Fig. 3. C). The western spread continued at a pace with the melting of the Cordilleran ice

170    sheet (Fig. 3. D), but the eastern expansion lagged behind due to the slower melting of the

171    Laurentide icesheet (Fig. 3. E). The central and eastern part of the continent were fully colonised by

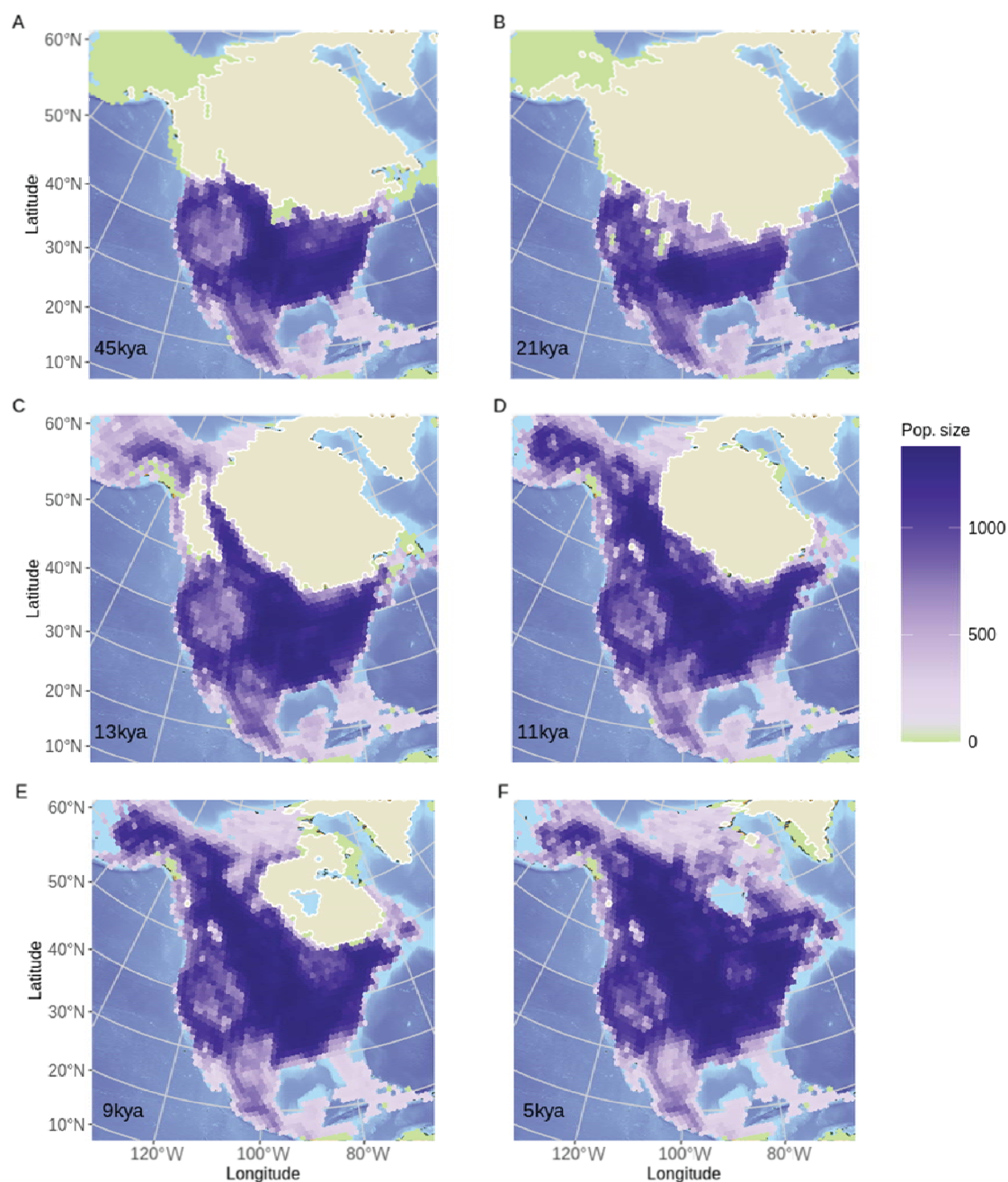172    5kya, when the ice sheets were fully melted (Fig. 3. F).



**Figure 3.** Weighted mean population size (per deme) of Yellow Warbler at A) 45kya, B) 21kya, C) 13kya, D) 11,

E) 9kya, F) 5kya, from 1055 simulations retained during the parameter estimation. Dark grey regions are areas

uninhabitable for yellow warblers at the given point in time.

173    The importance of this asymmetric expansion in setting up the patterns of genetic diversity across

174    the range of Yellow warblers can be seen by mapping geolocation of common ancestor (CA) events

175    that occurred between populations. These events allow us to reconstruct gene flow through time, as

176    shaped by colonisations and subsequent connectivity, revealing how the patterns of diversity have

177    emerged. Using two populations from two regions each time, we plotted locations of the CA events

178    between the East and West, Central and East, and Central and West regions (Fig. 4. A). When we

179    considered populations from the East and West (Fig. 4. B), we can see that common ancestor events

180    between these two clusters show a "v" shape that matches closely the shape of the ice sheets at

181    13kya, when the postglacial expansion occurred. Importantly, the same pattern was also found

182    when we considered only West and Central populations (Fig. 4. C), albeit with a greater intensity of

183    events in the corridor between the two populations.  Even though we did not have any population

184    from the East, common ancestors event reveal that central populations are linked to that area, thus

185    representing a mix of the western and eastern arms of the expansion. The same is true when we

186    considered only Central and East populations Fig. 4. D).  Together with the reconstructed

187    demography in Fig. 3., this pattern shows the importance of the early expansion up the west coast,

188    followed by subsequent expansion up the east coast, in setting up an initial divergence of the clades,

189    which then mixed in the central region comparatively recently.  The signature left by the relatively

190    strong founder events that occurred during the expansion have not yet been eroded by the

191    relatively low levels of migration, explaining the current patterns of genetic diversity and structure in

192    this species.

**Figure 4.** The location of common ancestor (CA) events are plotted across a map of North America. A) is an elevation map of the region with all six sampling locations labelled. In B-D colour density represents proportion of total CA events on the map that occur in each deme. B) is an analysis based on two populations each from West (blue) and Central (purple) regions, C) is based on Central and East (red) regions, and finally D) the West and East regions.

193    **Discussion**

194    In this study, we examined the relative roles of different forces that may have driven modern day

195    genetic structuring in a widespread species. We used a set of complementary datasets to explore

196    structure in the North American Yellow Warbler (*Setophaga petechia*), a common passerine species.

197    By integrating genetic data and climatic and environmental variables through time into a spatially-

198    explicit modelling framework (CISGeM), we were able to build a detailed reconstruction of the

199    population dynamics for this species, stretching back through the last fifty thousand years. Our

200    model was able to reconstruct population size changes, track potential range expansions, and

201    simulate recolonisation dynamics, whilst capturing the genetic structure found in the modern

202    population. With this information, we were able to explore the extent to which expansion dynamics

203    could explain modern genomic patterns of the Yellow Warbler.

204    East-west population structure, as found in the Yellow Warbler, is not an uncommon pattern in

205    North America. These genetic differences, as well as variation in other traits such as migratory

206    behaviour, are often considered to support the existence of isolated refugia during glaciations (e.g.

207    (16,27)). However, recent work on refugia has shown that the patterns of diversity found in the

208    Northern Hemisphere only fit the expectations from cyclic expansion-contraction fragmenting

209    ranges and driving genetic variation at a coarse level (8,9,28).

210    By explicitly modelling the recolonization dynamics, we have demonstrated a plausible explanation

211    for the formation of genetic structure over time, without the need of multiple glacial refugia. The

212    dynamics of the modelled Yellow Warbler recolonization show that, to a large extent, this passerine

213    species tracked the uneven (asynchronous) retreat of the Laurentide Ice Sheet, with a longitudinally

214    unequal progression northward (Fig. 3.). Despite the species exhibiting a single large glacial

215    refugium, the asymmetrical pattern of re-expansion generates the genetic structure of east, west,

216    and central population clusters found in the empirical genetic data. This implies a key role for post-

217    glacial re-expansion in shaping modern-day populations.

218    The important role of re-expansion dynamics has recently been highlighted in a range of different

219    species e.g. (8,11,28), though it would be naïve to assume that the complex patterns of diversity

220    found in real populations could be easily explained by a single mechanistic process (29).  Our work

221    highlights that, at the very least, modern population diversity and structure may have originated

222    from a combination of different processes, each of which needs to be carefully considered.

223    We acknowledge that, within this framework, we were unable to consider the possible influence of

224    biotic interactions which may have impacted the pattern of recolonization (30).  Our model also

225    works with demes that are discrete spatial units of a fixed size, allowing for a step change in the

226    likelihood of common ancestor events occurring within the deme and outside it.  Moving away from

227    the discretisation of space could help further 'naturalise' our model, and indeed models that

228    incorporate continuous space are rapidly advancing (31).  However, there are still major

229    computational challenges to overcome before these tools would be suitable for an area on the scale

230    of this study.

231    Whilst theories that describe broad patterns have been crucial to increasing our understanding of

232    the likely impacts environmental changes have had on populations, we now realise that North

233    American avifauna is probably a composite of species with different histories (32).  Species have

234    responded individually to the rapid climate changes faced in the Pleistocene and therefore we would

235    not wish to claim our findings refute the existence and effect of North American glacial refugia for

236    birds.  However, now the resources and techniques exist to study the idiosyncratic responses of

237    different species, and it will be possible to assess the importance of isolated refugia in shaping the

238    genetic structure of species. Furthermore, an increased understanding on the different population

239    dynamics that underlined species responses to the large climatic changes that occurred over the last

240    glacial cycle might provide an important tool to refine our ability to predict the responses of species

241    to anthropogenic change in the future.

242    **Materials and Methods**

243    **Study species**

244    The North American Yellow Warbler (*Setophaga petechia*) is a small, riparian, migratory passerine.

245    Today, this common species is widely distributed across the continent.  However, despite its large

246    and well-connected contemporary range, the Yellow Warbler exhibits spatial structure across its

247    range, including multiple mitochondrial clades (16) and clear isolation by distance (33). Although not

248    a species of concern, the Yellow Warbler has recorded a declining population trend in the North

249    American Breeding Bird Survey between 1966-2015, triggering several studies looking into the

250    species ability to cope in the face of a rapidly changing climate (33,34). One such study by Bay et al.

251    (33) built RAD-seq data from individuals sampled across the species' range in order to explore

252    potential population trends in response to future climate scenarios.  Such data was made available

253    on GenBank and forms the basis of our empirical dataset here.

254    **Raw genetic data**

255    RAD sequence data for North American Yellow Warblers (*Setophaga petechia*) from 21 populations

256    (33) were downloaded from the NCBI Sequence Read Archive (SRA). From the 269 accessions

257    associated with the Bay et al. paper we chose to focus on only the individuals included in the original

258    analysis (n = 223), individuals for which full information about their breeding population was

259    available. A further 22 samples were dropped as the file sizes were under 75MB and, therefore,

260    were likely to have low coverage. One final exclusion was made, GenBank accession number

261    SRR6366039, as the sample was found to be an outlier with a measure of diversity higher than the

262    range of all other samples, despite comparable levels of coverage and number of sites. This left 200

263    samples for further analysis. These individuals were sampled from across the modern population

264    range, providing a good overview of the population genetics of this species, see Fig 2. for sampling

265    locations.

266    **Clustering analysis**

267    RAD-seq methods are known to create specific biases in estimated allele frequencies, potentially

268    affecting downstream analysis of the data (35). Using allele frequencies derived directly from the

269    sequence data in a genotype-free method has been shown to account for RAD-seq specific issues,

270    improving population genetic inferences (35). Therefore, we used Analyses of Next-Generation

271    Sequencing Data (ANGSD) (36,37) to infer genotype likelihoods directly from aligned BAM files.

272    Filters were set to only include SNPs with a p value of $< 2 \times 10^{-6}$ and only keep sites with at least 100

273    informative individuals.  These ANGSD genotype likelihood values were then used as input for

274    NGSadmix to calculate population admixture, setting a $K$ (presumed cluster number) value of 2 and

275    keeping minimum informative individuals at 100.

276    **Observed genetics for CISGeM**

277    In order to calculate pairwise π (the average number of differences between two sequences,

278    normalised by the number of available positions), we first calculated genotype likelihoods in ANGSD.

279    Input files were aligned BAM files, we used the samtools genotype likelihood method and inferred

280    the major and minor allele from these likelihoods, with the command below:

281    angsd -GL 1 -out genolike -doGlf 1 -doMajorMinor 1 -bam bam.filelist

282    We then computed pairwise π from the ANGSD output; since our population genetic simulations

283    (see below) modelled haploid samples (as it is the case of most genetic simulators, e.g. msprime

284    (38)), we used the below formula:

$$\frac{pAB_{AB} \cdot 0.5 + pAA_{BB} + pAA_{AB} \cdot 0.5 + pAB_{AA} \cdot 0.5}{nSites}$$

285    In order to make the modelling computationally feasible, we then investigated how many samples

286    were needed to get a reliable estimate of π for each population (Supplementary Fig. 6.). This analysis

287    showed that five diploid individuals, or ten chromosomes, provided a reasonable compromise for

288    noise.  All estimates of pairwise π were therefore re-computed with only five individuals per

289    population. As estimates were consistent with the values from the full dataset (Supplementary Fig.

290 7.), for computation efficiency of the model, all future analyses were based on this subset of the

291 data.

292 **Species Distribution Modelling for world reconstruction**

293 The range and population size of a species changes in time and space according to fluctuations in

294 resources and environmental conditions. In order to build a spatially explicit model it is first

295 necessary to use Species Distribution Modelling (SDM) to reconstruct how population ranges and

296 demographics may have changed through time. For this study an SDM analysis was undertaken using

297 an R (39) pipeline.

298 **Climate reconstructions:**

299 Climate data for North America were drawn from a 0.5° resolution dataset for 19 bioclimatic

300 variables; Net Primary productivity (NPP), Leaf Area Index (LAI) and all the BioClim variables (40)

301 with the exclusion of BIO2 and BIO3; covering the last 50,000 years in 1,000 year time steps from the

302 present to 22kya and in 2,000years time steps before that date (41). This dataset was originally

303 constructed from a combination of HadCM3 climate simulations of the last 120,000 years (42), high-

304 resolution HadAM3H simulations of the last 21,000 years (43), and empirical present-day data. The

305 data had been downscaled and bias-corrected using the Delta Method (44). Bioclimatic variables

306 through time were then used as input data to inform the SDM.

307 **SDM data preparation:**

308 Species occurrences data for the present day were initially downloaded from the GBIF database

309 (https://www.gbif.org), the original downloads are available at the following DOI: S. petechia

310 10.15468/dl.jfkwcg (GBIF.org). These data were then filtered based on the attributed accuracy of the

311 coordinates (maximum error: 1 km) and additionally, only points that were within Birdlife breeding

312 and resident geographical ranges (45) were retained. Remaining occurrences were then matched to

313    the 0.5° resolution grid used for the palaeoclimatic reconstructions and, as the method works on

314    presence/absence data and not frequency, only one presence per grid cell was kept.

315    This cleaned observation dataset was then used to define a set of informative bioclimatic variables

316    with the most influence on the species distribution for use in the Species Distribution Model (SDM),

317    through visual check of how much the distribution of the variable values differed between the

318    observation points and the whole area.  We selected the variables with highest differences between

319    the two curves, which are most likely to be relevant for the species, and then, in order to avoid using

320    highly correlated variables, which may increase noise in the data, we constructed a correlation

321    matrix between the variables associated with each of the retained observations.  Where two values

322    were highly correlated, the variable with the lowest overall correlation across the matrix was kept,

323    allowing us to select a set of uncorrelated variables (threshold = 0.7) leaving us with the following

324    ones to be used for SDM modelling: LAI (leaf area index), BIO7 (Temperature Annual Range), BIO8

325    (Mean Temperature of Wettest Quarter), BIO14 (Precipitation of Driest Month).

326    Geographic biases in sampling effort are common when observation data are collected

327    opportunistically, such as the data in the GBIF database.  In order to reduce this bias, we thinned our

328    dataset using the R package *spThin* (46) enforcing a minimum distance of 70 km between

329    observations.  Given the random nature of removing nearest-neighbour data points, we repeated

330    this step 100 times ('rep' = 100) retaining for further analysis the result with the maximum number

331    of observations after thinning.

332    **SDM modelling:**

333    The SDM was built with the R package biomod2 (47) following the same procedure used in Miller et

334    al. (48).  The thinned observation dataset was used as presences whilst the landmass of North

335    America was considered as background.  The same number of pseudo-absences as presences were

336    then drawn five separate times, at random, from outside the BirdLife resident and breeding masks:

337    creating five independent datasets for analysis. For each data set, following Bagchi et al. (49),

338    models were then run independently using four different algorithms: generalised linear models

339    (GLM), generalized boosting method (GBM), generalised additive models (GAM), and random forest.

340

341    Spatial cross-validation was used to evaluate the model; 80% of the data were used to train the

342    algorithm and the remaining 20% to test it. Initially, both the presences and the five pseudoabsences

343    datasets were subdivided in 14 latitudinal bands using the R package BlockCV (50). Each band was

344    given a 'band ID number', looping sequentially through numbers 1-5 until all bands were labelled.

345    Then the bands were assembled into five working data splits grouped by their band ID (numbers 1-

346    5). This was performed to maximise the probability of having at least some presences in all five data

347    splits as a data split cannot be used for evaluation if it contains only absences. Each of the four

348    models (GLM, GBM, GAM, and random forest) were then run five times (once for each

349    pseudoabsence run), using in turn four of the five defined data splits to calibrate and one to evaluate

350    based on TSS (threshold = 0.7).

351    Finally, a full ensemble combining all algorithms and pseudoabsences runs (51) was created, using

352    only models with TSS > 0.7, averaged using four different statistics: mean, median, committee

353    average and weighted mean. The statistic showing the highest TSS, the mean, was then used to

354    predict the probability of occurrence in each grid cell.  This was then projected for all available time

355    slices from the present to 50 thousand years ago.

356            **CISGeM Demography:**

357    CISGeM's demographic module consists of a spatial model that simulates long-term and global

358    growth and migration dynamics of Yellow Warblers. These processes depend on a number of

359    parameters (see Table 1.), which we later estimate statistically based on empirical genetic data.

| $\alpha$ | Allometric scaling factor for population size |
|---|---|
| $\beta$ | Allometric scaling exponent for population size |
| $t_0$ | Temporal origin of the population |
| $r$ | Intrinsic growth rate |
| $m_r$ | Non-directed (random) mobility parameter |
| $m_d$ | Directed mobility parameter |

360 **Table 1.** Details of parameters used in CISGeM.

361 The model operates on a global hexagonal grid of 40962 cells that represent the whole world (the

362 distance between the centre of two hexagonal cells is 241 ±15 km); 2422 grid cells make up North

363 America. Each time step represents 1 year, the generation time of Yellow Warblers. Each time step

364 of a simulation begins with the computation of the carrying capacity of each grid cell, i.e. the

365 maximum number of YWs theoretically able to live in the cell for the environmental resources at the

366 given point in time. Here, we estimate the carrying capacity in a grid cell $x$ at a time $t$ as

$$K(x, t) \stackrel{\text{def}}{=} \begin{cases} \alpha \cdot p(x, t)^{\beta}, \text{ if } x \text{ is on land at time } t \\ 0, \text{ else} \end{cases}$$

367 where $p(x, t)$ denotes the probability of a species inhabiting cell $x$ at time $t$ (see section 'Species

368 Distribution Modelling'). The particular function used here was chosen based on analysis of SDM

369 projections and census data of Holarctic birds (R. Green, pers. comm.).

370 The estimated carrying capacities are used to simulate spatial population dynamics as follows. We

371 begin a simulation by initialising a population of yellow warblers in a grid cell $x_0$ (represents the

372 spatial origin of yellow warbler in our model) at a point in time $t_0$ with $K(x_0, t_0)$ individuals.

373 At each subsequent time step between $t_0$ and the present, CISGeM simulates two processes: the

374 local growth of populations within grid cells, and the spatial migration of individuals across cells. We

375 used the logistic function to model local population growth, estimating the net number of individuals

376 by which the population of size $N(x, t)$ in the a $x$ at time $t$ increases within the time step as

377
$$r \cdot N(x,t) \cdot \left(1 - \frac{N(x,t)}{K(x,t)}\right),$$

378    where $r$ denotes the intrinsic growth rate. Thus, growth is approximately exponential at low

379    population sizes, before decelerating, and eventually levelling off at the local carrying capacity.

380    Across-cell migration is modelled as two separate processes, representing a non-directed, spatially

381    uniform movement into all neighbouring grid cells on the one hand, and a directed movement along

382    a resource availability gradient on the other hand. Under the first movement type, the number of

383    individuals migrating from a cell $x$ into any one of the up to six neighbouring cells is estimated as

384
$$m_r \cdot N(x,t),$$

385    where $m_r$ is a mobility parameter. This mechanism is equivalent to a spatially uniform diffusion

386    process, which has previously been used to model random movement in other species (52). Under

387    the second movement type, an additional number of individuals moving from a grid cell $x_1$ to a

388    neighbouring cell $x_2$ is estimated as

$$m_d \cdot N(x_1,t) \cdot max\left(0, \frac{K(x_2,t) - N(x_2,t)}{K(x_2,t)} - \frac{K(x_1,t) - N(x_1,t)}{K(x_1,t)}\right)$$

389    The number $\frac{K(x,t)-N(x,t)}{K(x,t)}$ represents the relative availability of unused resources in the cell $x$ at time

390    $t$, equalling 1 if all natural resources in $x$ are potentially available for yellow warblers ($N(x,t) = 0$),

391    and 0 if all resources are used ($N(x,t) = K(x,t)$). Thus, individuals migrate in the direction of

392    increasing relative resource availability, and the number of migrants is proportional to the steepness

393    of the gradient. The distinction between directed and non-directed movement allows us to examine

394    to which extent migration patterns can be explained by random motion alone or requires us to

395    account for more complex responses to available resources.

396    For some values of the mobility parameters $m_r$ and $m_d$, it is possible for the calculated number of

397    migrants from a cell to exceed the number of individuals in that cell. In this scenario, the number of

398    migrants into neighbouring cells are rescaled proportionally such that the total number of migrants

399    from the cell is equal to the number of individuals present.

400    Similarly, it is in principle possible that the number of individuals present in a cell after all migrations

401    are accounted for (i.e., the sum of local non-migrating individuals, minus outgoing migrants, plus

402    incoming migrants from neighbouring cells) exceeds the local carrying capacity. In this case,

403    incoming migrants are rescaled proportionally so that the final number of individuals in the cell is

404    equal to the local carrying capacity. In other words, some incoming migrants perish before

405    establishing themselves in the destination cell, and these unsuccessful migrants are not included in

406    the model's output of migration fluxes between grid cells. In contrast, non-migrating local residents

407    remain unaffected in this step. They are assumed to benefit from a residential advantage (53), and

408    capable of outcompeting incoming migrants.

409    CISGeM's demographic module outputs the number of individuals in each grid cell, and the number

410    of migrants between neighbouring grid cells, across all time steps of a simulation. These quantities

411    are the used to reconstruct genetic lineages.

412    **CISGeM predicted genetics**

413    Once a global population demography has been constructed, gene trees are simulated. This process

414    is dependent on the population dynamics recorded in the demography stage and assumes local

415    random mating according to the Wright-Fisher dynamic. From the present, ancestral lines of

416    sampled individuals are tracked back through the generations, recording which cell each line belongs

417    to. Every generation, the lines are randomly assigned to a gamete from the individuals within its

418    present cell. If the assigned individual is a migrant or coloniser, the line moves to the cell of origin for

419    that individual before 'reproduction'. Whenever two lines are assigned to the same parental

420    gamete, this is recorded as a coalescent event, and the two lines merge into a single line

421    representing their common ancestor. This process is repeated until all the lineages have met,

422    reaching the common ancestor of the whole sample. If multiple lineages are still present when the

423    model reaches the generation and deme from which the demography was initialised, the lines enter

424    a single ancestral population ($K_0$) until sufficient additional coalescent events have occurred for the

425    gene tree to close.

426        **ABC parameter estimation**

427    Parameter space was explored with a Monte Carlo sweep in which demographic parameters were

428    randomly sampled from flat prior ranges: directed expansion coefficient [0.0,0.14], undirected

429    expansion coefficient [0.0,0.04], intrinsic growth rate [0.02,0.15], allometric scaling exponent [0.1,1],

430    and allometric scaling factor [20,5000] on a $\log_{10}$ scale.  A fixed mutation rate of $2.3 \times 10^{-9}$ μ/Site/Year

431    was used (54).

432    Model fit was initially calculated within an Approximate Bayesian Computation (abc) framework

433    using the results of the Monte Carlo sweep. To compute summary statistics, populations were

434    clustered into three groups representing the West, Central, and East regions of the North American

435    continent, based on the NGSadmix outputs.  The mean pairwise π for populations was then

436    computed within each group and between each pair of groups, giving us a total of 6 summary

437    statistics.

438    We performed parameter estimation with the R package '*abc*' (25) using a local linear abc algorithm,

439    setting the tolerance to 0.025.  For each of the simulations retained by the abc analysis,

440    demographic simulations were then recorded and combined to create an average, representative,

441    profile of the population's demographic history.

442        **ABC model fitting: gfit & gfitpca**

443    We also confirmed the quality of the model fit using formal hypothesis testing approaches from the

444    R package '*abc*' (25).  Firstly we used the '*gfit*' function (55) to confirm that our model outperformed

445    a series on null models.  In this function the goodness of fit test statistic, or D-statistic, is the median

446    Euclidean distance between the observed summary statistics and the nearest (accepted) summary

447    statistics.  For comparison, a null distribution of D is then generated from summary statistics of 1000

448    pseudo-observed datasets.  A goodness of fit p-value can then be calculated as the proportion of D

449    based on pseudo-observed data sets that are larger than the empirical value of D.  Consequently, a

450    non-significant p-value signifies that the distance between the observed and accepted summary

451    statistics is not larger than the expectation, confirming that the model fits the observed data well.

452    We then further performed an a priori goodness of fit test using the *'gfitpca'* function which

453    captures and plots the two first components obtained with a principle component analysis.  We used

454    a *'cprob'* value of 0.1, 0.15, and 0.2, leaving a different proportion of points from the model outside

455    the displayed envelope. The observed summary statistics is then marked to check that it is contained

456    within these envelopes, indicating a good fit.

457        **ABC model fitting: abcrf**

458    We further evaluated model fit and posterior distributions with an abc random forest (RF) approach

459    implemented via the R package *'abcrf'* (26,56). Forests of 1,000 trees were used.

460

461 **Author contribution:**

462 E.F.M. and A.M. devised the project. E.F.M. ran the genetic analysis under the supervision of P.M.D.

463 and A.M.. M.L. ran the Species Distribution Modelling. M.K. and R.B. and A.M. developed the

464 modelling software with help from E.F.M, P.M.D. and G.L.S.. M.S. helped with the interpretation of

465 results. E.F.M and A.M. wrote the manuscript with feedback from all the co-authors.

466 **Competing interests:**

467 The authors declare no competing interests.

468   1.   Kelly JF, Hutto RL. An East-West Comparison of Migration in North American Wood Warblers.

469        Condor. 2005;107(2):197–211.

470   2.   Lovette IJ, Clegg SM, Smith TB. Limited Utility of mtDNA Markers for Determining

471        Connectivity among Breeding and Overwintering Locations in Three Neotropical Migrant

472        Birds. Conserv Biol. 2004;18(1):156–66.

473   3.   Peters JL, Gretes W, Omland KE. Late Pleistocene divergence between eastern and western

474        populations of wood ducks (Aix sponsa) inferred by the "isolation with migration" coalescent

475        method. Mol Ecol. 2005;14(11):3407–18.

476   4.   Hewitt GM. Postglacial re-colonisation of European biota. Biol J Linn Soc. 1999;68(May):87–

477        112.

478   5.   Haffer J. Speciation in amazonian forest birds. Vol. 165, Science. 1969. p. 131–7.

479   6.   Davis LA, Roalson EH, Cornell KL, Mcclanahan KD, Webster MS. Genetic divergence and

480        migration patterns in a North American passerine bird: Implications for evolution and

481        conservation. Mol Ecol. 2006;15(8):2141–52.

482   7.   Colbeck GJ, Gibbs HL, Marra PP, Hobson K, Webster MS. Phylogeography of a widespread

483        North American migratory songbird (Setophaga ruticilla). J Hered. 2008;99(5):453–63.

484   8.   Bemmels JB, Dick CW. Genomic evidence of a widespread southern distribution during the

485        Last Glacial Maximum for two eastern North American hickory species. J Biogeogr.

486        2018;45(8):1739–50.

487   9.   Lumibao CY, Hoban SM, McLachlan J. Ice ages leave genetic diversity 'hotspots' in Europe but

488        not in Eastern North America. Ecol Lett. 2017;20(11):1459–68.

489   10.  Excoffier L, Foll M, Petit RJ. Genetic consequences of range expansions. Annu Rev Ecol Evol

490        Syst. 2009;40:481–501.

491     11.     Reid BN, Kass JM, Wollney S, Jensen EL, Russello MA, Viola EM, et al. Disentangling the

492             genetic effects of refugial isolation and range expansion in a trans-continentally distributed

493             species. Heredity (Edinb) [Internet]. 2019;122(4):441–57. Available from:

494             http://dx.doi.org/10.1038/s41437-018-0135-5

495     12.     Raghavan M, Steinrücken M, Harris K, Schiffels S, Rasmussen S, DeGiorgio M, et al. Genomic

496             evidence for the Pleistocene and recent population history of Native Americans. Science (80-

497             ). 2015;349(6250).

498     13.     Eriksson A, Betti L, Friend AD, Lycett SJ, Singarayer JS, von Cramon-Taubadel N, et al. Late

499             Pleistocene climate change and the global expansion of anatomically modern humans. Proc

500             Natl Acad Sci [Internet]. 2012 Oct 2;109(40):16089–94. Available from:

501             http://www.pnas.org/cgi/doi/10.1073/pnas.1209494109

502     14.     Bay R, Harrigan R, Underwood V Le, Gibbs HL, Smith TB, Ruegg K. Response to Comment on

503             "Genomic signals of selection predict climate-driven population declines in a migratory bird"

504             Science. 2018;361(August):2–4.

505     15.     Skotte L, Korneliussen TS, Albrechtsen A. Estimating individual admixture proportions from

506             next generation sequencing data. Genetics. 2013;195(3):693–702.

507     16.     Boulèt M, Gibbs HL, Hobson KA. Integrated analysis of genetic, stable isotope, and banding

508             data reveal migratory connectivity and flyways in the northern yellow warbler (Dendroica

509             petechia; aestiva group). Ornithol Monogr. 2006;61(July 2015):29–78.

510     17.     Bay RA, Karp DS, Saracco JF, Anderegg WRL, Frishkoff L, Wiedenfeld D, et al. Genetic variation

511             reveals individual-level climate tracking across the full annual cycle of a migratory bird.

512             bioRxiv [Internet]. 2020;(preprint). Available from:

513             https://www.biorxiv.org/content/10.1101/2020.04.15.043331v1.abstract

514     18.     Elith J, Leathwick JR. Species Distribution Models: Ecological Explanation and Prediction

515   Across Space and Time. Annu Rev Ecol Evol Syst. 2009;40(1):677–97.

516 19. McCullagh P, Nelder JA. Generalized Linear Models, 2nd Edn. Vol. 39. Chapman and Hall;

517   1990. 385 p.

518 20. Elith J, Leathwick JR, Hastie T. A working guide to boosted regression trees. J Anim Ecol.

519   2008;77(4):802–13.

520 21. Hastie TJ, Tibshirani. RJ. Generalized additive models. Chapman and Hall; 1990.

521 22. Breiman L. Random Forest. Mach Learn [Internet]. 2001 Oct 31;5–32. Available from:

522   https://www.taylorfrancis.com/books/9780429890277/chapters/10.1201/9780429469275-8

523 23. Darvill CM, Menounos B, Goehring BM, Lian OB, Caffee MW. Retreat of the Western

524   Cordilleran Ice Sheet Margin During the Last Deglaciation. Geophys Res Lett.

525   2018;45(18):9710–20.

526 24. Margold M, Stokes CR, Clark CD. Reconciling records of ice streaming and ice margin retreat

527   to produce a palaeogeographic reconstruction of the deglaciation of the Laurentide Ice Sheet.

528   Quat Sci Rev [Internet]. 2018;189:1–30. Available from:

529   https://doi.org/10.1016/j.quascirev.2018.03.013

530 25. Csilléry K, François O, Blum MGB. Abc: An R package for approximate Bayesian computation

531   (ABC). Methods Ecol Evol. 2012;3(3):475–9.

532 26. Pudlo P, Marin JM, Estoup A, Cornuet JM, Gautier M, Robert CP. Reliable ABC model choice

533   via random forests. Bioinformatics. 2016;32(6):859–66.

534 27. Ruegg KC, Smith TB. Not as the crow flies: A historical explanation for circuitous migration in

535   Swainson's thrush (Catharus ustulatus). Proc R Soc B Biol Sci. 2002;269(1498):1375–81.

536 28. Markova S, Hornikova M, Lanier HC, Henttonen H, Searle JB, Weider LJ, et al. High genomic

537   diversity in the bank vole at the northern apex of a range expansion: the role of multiple

538     colonizations and end-glacial refugia. Mol Ecol [Internet]. 2020;0–3. Available from:

539     https://onlinelibrary.wiley.com/doi/abs/10.1002/jmv.25688

540  29.  Cheviron ZA, Hackett SJ, Capparella AP. Complex evolutionary history of a Neotropical

541     lowland forest bird (Lepidothrix coronata) and its implications for historical hypotheses of the

542     origin of Neotropical avian diversity. Mol Phylogenet Evol. 2005;36(2):338–57.

543  30.  Pearson RG, Dawson TP. Predicting the impacts of climate change on the distribution of

544     species: are bioclimate envelope models useful? Glob Ecol Biogeogr [Internet]. 2003

545     Sep;12(5):361–71. Available from: http://doi.wiley.com/10.1046/j.1466-822X.2003.00042.x

546  31.  Haller BC, Messer PW. SLiM 3: Forward Genetic Simulations Beyond the Wright-Fisher Model.

547     Mol Biol Evol. 2019;36(3):632–7.

548  32.  Zink RM. Comparative phylogeography in North American birds. Evolution (N Y).

549     1996;50(1):308–317.

550  33.  Bay RA, Harrigan RJ, Underwood V Le, Gibbs HL, Smith TB, Ruegg K. Genomic signals of

551     selection predict climate-driven population declines in a migratory bird. Science (80- ). 2018

552     Jan 5;359(6371):83–6.

553  34.  Mazerolle DF, Dufour KW, Hobson KA, Haan HE Den. Effects of large-scale climatic

554     fluctuations on survival and production of young in a Neotropical migrant songbird, the

555     yellow warbler Dendroica petechia. J Avian Biol [Internet]. 2005 Feb;36(2):155–63. Available

556     from: https://doi.org/10.1111/j.0908-8857.2005.03289.x

557  35.  Warmuth VM, Ellegren H. Genotype-free estimation of allele frequencies reduces bias and

558     improves demographic inference from RADSeq data. Mol Ecol Resour. 2019 May

559     17;19(3):586–96.

560  36.  Nielsen R, Korneliussen T, Albrechtsen A, Li Y, Wang J. SNP calling, genotype calling, and

561     sample allele frequency estimation from new-generation sequencing data. PLoS One.

562        2012;7(7).

563    37.    Korneliussen TS, Albrechtsen A, Nielsen R. ANGSD: Analysis of Next Generation Sequencing

564        Data. BMC Bioinformatics. 2014 Dec 25;15(1):356.

565    38.    Kelleher J, Etheridge AM, McVean G. Efficient Coalescent Simulation and Genealogical

566        Analysis for Large Sample Sizes. PLoS Comput Biol. 2016;12(5):1–22.

567    39.    R Core Team. R: A Language and Environment for Statistical Computing [Internet]. 2019.

568        Available from: http://www.r-project.org/

569    40.    Hijmans RJ, Cameron SE, Parra JL, Jones PG, Jarvis A. Very high resolution interpolated

570        climate surfaces for global land areas. Int J Climatol. 2005;25(15):1965–78.

571    41.    Beyer RM, Krapp M, Manica A. High-resolution terrestrial climate, bioclimate and vegetation

572        for the last 120,000 years. Sci Data. 2020;7(1):1–9.

573    42.    Singarayer JS, Valdes PJ. High-latitude climate sensitivity to ice-sheet forcing over the last 120

574        kyr. Quat Sci Rev [Internet]. 2010;29(1–2):43–55. Available from:

575        http://dx.doi.org/10.1016/j.quascirev.2009.10.011

576    43.    Armstrong E, Hopcroft PO, Valdes PJ. Reassessing the Value of Regional Climate Modeling

577        Using Paleoclimate Simulations. Geophys Res Lett. 2019;46(21):12464–75.

578    44.    Beyer R, Krapp M, Manica A. An empirical evaluation of bias correction methods for

579        palaeoclimate simulations. Clim Past. 2020;16(4):1493–508.

580    45.    BirdLife International and Handbook of the Birds of the World. Bird species distribution maps

581        of the world. Version 2018.1. 2018.

582    46.    Aiello-Lammens ME, Boria RA, Radosavljevic A, Vilela B, Anderson RP. spThin: An R package

583        for spatial thinning of species occurrence records for use in ecological niche models.

584        Ecography (Cop). 2015;38(5):541–5.

585    47.    Thuiller W, Georges D, Engler R, Breiner F. biomod2: Ensemble Platform for Species

586            Distribution Modeling [Internet]. 2019. Available from: https://cran.r-

587            project.org/package=biomod2

588    48.    Miller EF, Green RE, Balmford A, Beyer R, Somveille M, Leonard M, et al. mtDNA-based

589            reconstructions of change in effective population sizes of Holarctic birds do not agree with

590            their reconstructed range sizes based on paleoclimates. bioRxiv. 2019;

591    49.    Bagchi R, Crosby M, Huntley B, Hole DG, Butchart SHM, Collingham Y, et al. Evaluating the

592            effectiveness of conservation site networks under climate change: Accounting for

593            uncertainty. Glob Chang Biol. 2013;19(4):1236–48.

594    50.    Valavi R, Elith J, Lahoz-Monfort JJ, Guillera-Arroita G. blockCV: An r package for generating

595            spatially or environmentally separated folds for k-fold cross-validation of species distribution

596            models. Methods Ecol Evol. 2019;10(2):225–32.

597    51.    Araújo MB, New M. Ensemble forecasting of species distributions. Trends Ecol Evol.

598            2007;22(1):42–7.

599    52.    Kot M. Elements of Mathematical Ecology [Internet]. Cambridge University Press; 2001. 453

600            p. Available from:

601            https://books.google.co.uk/books?id=Zh3GNd9M1oUC&printsec=frontcover&source=gbs_ge

602            _summary_r&cad=0

603    53.    Pérez-Tris J, Tellería JL. Migratory and sedentary blackcaps in sympatric non-breeding

604            grounds: Implications for the evolution of avian migration. J Anim Ecol. 2002;71(2):211–24.

605    54.    Smeds L, Qvarnström A, Ellegren H. Direct estimate of the rate of germline mutation in a bird.

606            Genome Res. 2016;26(9):1211–8.

607    55.    Lemaire L, Jay F, Lee I-H, Csilléry K, Blum MGB. Goodness-of-fit statistics for approximate

608            Bayesian computation. 2016;1–30. Available from: http://arxiv.org/abs/1601.04096

609     56.     Raynal L, Marin JM, Pudlo P, Ribatet M, Robert CP, Estoup A. ABC random forests for

610             Bayesian parameter inference. Bioinformatics. 2019;35(10):1720–8.

611

612