# MOLECULAR ECOLOGY

## Comparing RADseq and microsatellites to infer complex phylogeographic patterns, an empirical perspective in the Crucian carp, *Carassius carassius*, L.

SCHOLARONE™
Manuscripts

1 # Comparing RADseq and microsatellites to infer
2 # complex phylogeographic patterns, an empirical
3 # perspective in the Crucian carp, *Carassius carassius,*
4 # L.

5

6 **Authors:** [1]Daniel L Jeffries, [2]Gordon H Copp, [1]Lori Lawson Handley, [3]K. Håkan Olsén, [4]Carl D
7 Sayer, [1]Bernd Hänfling

8

9 [1] *Evolutionary Biology Group, School of Biological, Biomedical and Environmental Sciences,*
10 *Hardy Building, University of Hull, Hull, HU6 7RX, UK*
11 [2]*Salmon & Freshwater Team, Cefas, Pakefield Road, Lowestoft, Suffolk NR33 0HT, UK, and*
12 *Department of Life and Environmental Sciences, Faculty of Science and Technology, Bournemouth*
13 *University, Poole, UK*
14 [3]*School of Natural Science, Technology and Environmental Studies, Södertörn University, Alfred*
15 *Nobels allé 7, Flemingsberg, 141 89 Huddinge, Sweden*
16 [4] *Environmental Change Research Centre, Department of Geography, University College London,*
17 *Pearson Building, Gower Street, London, WC1E 6BT, UK*

18

22

23 **Corresponding author:** Daniel Jeffries (dljeffries86@gmail.com) *Evolutionary Biology Group,*
24 *School of Biological, Biomedical and Environmental Sciences, Hardy Building, University of Hull,*
25 *Hull, HU6 7RX, UK*

26

27 **Running title:** The complex phylogeography of the crucian carp

28

29 # Abstract

30   The conservation of threatened species must be underpinned by phylogeographic knowledge. This

31   need is epitomised by the freshwater fish *Carassius carassius*, which is in decline across much of

32   its European range. Restriction site associated DNA sequencing (RADseq) is increasingly used for

33   such applications, however RADseq is expensive, and limitations on sample number must be

34   weighed against the benefit of large numbers of markers. This trade-off has previously been

35   examined using simulation studies, however, empirical comparisons between these markers,

36   especially in a phylogeographic context, are lacking. Here, we compare the results from

37   microsatellites and RADseq for the phylogeography of *C. carassius* to test whether it is more

38   advantageous to genotype fewer markers (microsatellites) in many samples, or many markers

39   (SNPs) in fewer samples. These datasets, along with data from the mitochondrial cytochrome b

40   gene, agree on broad phylogeographic patterns; showing the existence of two previously

41   unidentified *C. carassius* lineages in Europe; one found throughout northern and central-eastern

42   European drainages, and a second almost exclusively confined to the Danubian catchment. These

43   lineages have been isolated for approximately 2.15 M years, and should be considered separate

44   conservation units. RADseq recovered finer population structure and stronger patterns of IBD than

45   microsatellites, despite including only 17.6% of samples (38% of populations and 52% of samples

46   per population). RADseq was also used along with Approximate Bayesian Computation to show

47   that the postglacial colonisation routes of *C. carassius* differ from the general patterns of freshwater

48   fish in Europe, likely as a result of their distinctive ecology.

49 # Introduction

50　Phylogeographic studies have revealed that the contemporary distributions of European taxa and

51　their genetic diversity have been largely shaped by the glacial cycles of the Pleistocene epoch, and

52　in particular by range shifts during recolonisation from glacial refugia (Hewitt 1999). In freshwater

53　fishes, the dynamics of recolonisation are tightly linked to the history of river drainage systems

54　(Bianco 1990; Bănărescu 1990, 1992; Bernatchez & Wilson 1998; Reyjol *et al.* 2006). For example,

55　watersheds pose a significant barrier to fish dispersal, often resulting in strong genetic structuring

56　across separate drainage systems (Durand *et al.* 1999; Hänfling *et al.* 2002). However, during

57　glacial melt periods, ephemeral rivers and periglacial lakes can arise, providing opportunities for

58　colonisation (Gibbard *et al.* 1988) of otherwise isolated drain basins (Grosswald 1980; Arkhipov *et*

59　*al.* 1995). These processes have resulted in complicated recolonisation scenarios in Europe, which,

60　in contrast to North America (Bernatchez & Wilson 1998), appear to possess few general patterns

61　of population structure. Furthermore, previous phylogeographic studies have predominantly focused

62　on highly mobile, obligatory or facultatively lotic species, with more sedentary, lentic species being

63　largely overlooked.

64

65　The crucian carp, *Carassius carassius* (Linnaeus 1758), is native to parts of central, eastern and

66　northern Europe and almost exclusively restricted to lentic ecosystems, including lakes, ponds and

67　river floodplains (Copp 1991; Copp *et al.* 2008). *C. carassius*, has recently experienced sharp

68　declines in the number and sizes of populations throughout its native range, leading to some local

69　population extinctions. The reasons for these declines include habitat loss through drought and

70　terrestrialisation in England (Copp 1991; Wheeler 2000; Sayer *et al.* 2011), acidification

71　(Holopainen & Oikari 1992), poor water quality in the Danube river catchment (Navodaru *et al.*

72　2002), and hybridisation with several non-native species (Copp *et al.* 2010; Savini *et al.* 2010;

73　Mezhzherin *et al.* 2012; Wouters *et al.* 2012; Rylková *et al.* 2013). The susceptibility of *C.*

74    *carrassius* to genetic isolation and bottlenecks is compounded by small population sizes (Hänfling

75    *et al.* 2005) and low dispersal (Holopainen *et al.* 1997). Strong geographic structure is therefore

76    likely in this species. Although the threats to *C. carassius* populations are recognised on a regional

77    level (Lusk *et al.* 2004; Mrakovčić *et al.* 2007; Wolfram & Mikschi 2007; Simic, V *et al.* 2009;

78    Copp & Sayer 2010), a global conservation strategy is missing. Broad scale phylogeographic data

79    and definition of evolutionary significant units are essential for informing unified conservation

80    efforts for this species (Frankham *et al*. 2002).

81

82    Phylogeographic data have traditionally been collected using mitochondrial gene regions and/or

83    nuclear markers such as AFLPs and microsatellites. However, cost and time often limits the number

84    of these nuclear markers used, which can result in low power for addressing phylogeographic

85    questions (Cornuet & Luikart 1996; Luikart & Cornuet 2008; Landguth *et al.* 2012; Peery *et al.*

86    2012; Hoban *et al.* 2013). Single nucleotide polymorphisms (SNPs) are increasingly used in

87    phylogeography for assessments of population structure (for example see Morin *et al.* 2010;

88    Emerson *et al.* 2010; Hess *et al.* 2011; Hauser *et al.* 2011). However, being bi-allelic, SNP loci

89    contain less information than highly polymorphic microsatellites (Coates *et al.* 2009) and therefore

90    large numbers of SNPs are needed to provide adequate statistical power. SNP discovery and assay

91    development, which has been costly and slow in the past, has recently been greatly facilitated by the

92    invention of restriction site associated DNA sequencing (RADseq, (Miller *et al.* 2006)), which

93    enables the fast identification of thousands of orthologous SNP markers in non-model organisms.

94    Nevertheless, although next generation sequencing costs are falling, RADseq remains a relatively

95    expensive approach, which often constrains the number of biological samples that can be included

96    in a given study. Researchers are, therefore, faced with a trade-off between the number of samples

97    and loci during study design. The optimal balance between the two is likely to be based on several

98    important but often unknown properties of the study system in question, for example the strength of

99    population structure (i.e. $F_{ST}$). Identifying these properties and comparing the relative strengths and

100    weaknesses of different molecular markers have recently been highlighted as priority topics in

101    landscape genetics and phylogeography (Epperson *et al.* 2010; Balkenhol & Landguth 2011).

102    Recent simulation studies have provided some important insights into this trade-off, for example,

103    Schwartz & McKelvey (2009) find that patchy geographic sampling along an IBD gradient could

104    result in falsely identified distinct lineages, whereas Landguth *et al.* (2012) find that increasing the

105    number of loci can strengthen the correlation between genetic and geographic distance for a given

106    sample set. To date, comprehensive empirical comparisons between microsatellite and SNP markers

107    in a phylogeographic context are lacking (but see Bradbury *et al.* 2015).

108

109    In the present study, we use a combination of mitochondrial DNA (mtDNA), microsatellites and

110    genome-wide SNPs obtained from RADseq in order to: 1) produce a comprehensive

111    phylogeography for *C. carassius* as a basis for Europe-wide conservation strategies, 2) test

112    competing scenarios of postglacial recolonisation that have potentially contributed to the

113    contemporary distribution of the species, and 3) compare the power of microsatellites and RADseq

114    based population structure analyses, in the context of the first two objectives. In this third aim, we

115    specifically ask, whether the benefits gained by the high numbers of markers obtained from

116    RADseq outweigh the potential loss of power associated by the reduction in the number of samples

117    in our system.

## Materials and Methods

118

### *Sample collection and DNA extraction*

119

120    *C. carassius* is a Cyprinid native to much of continental Europe; latitudinally from the North Sea

121    and Baltic Sea basins, through central Europe north of the Alps down to the Ponto-Caspian region

122    and longitudinally from Belgium and perhaps northern France into Siberia (Lelek 1980). However,

123    the true extent of this native range is unknown, largely due to difficulties in morphologically

124    distinguishing it from three closely related, introduced and widespread species: *Carassius auratus*,

125    *Carassius gibelio*, and *Cyprinus carpio* (Wheeler 2000; Hickley & Chare 2004). We initially

126    collected 1354 samples from 72 populations across 13 European countries, but due to frequent

127    hybridisation between the *C. carassius* and the three species mentioned above, it was necessary to

128    identify and remove hybrids from this sample set. To this end, all samples were first genotyped at 6

129    species diagnostic microsatellite loci. We removed all samples identified as hybrids from the dataset

130    and, to safeguard against cryptic hybridisation, we also removed all *C. carassius* that were

131    sympatric with hybrids (see SI text for full details of species identification and hybrid detection).

132    This left 867 *C. carassius* samples from 57 populations across the species' distribution in central

133    and northern Europe (Table 1, Fig. 1). Sample sizes ranged from n=4 to n=37, with a mean of n=17

134    (Table 1). Fish were anaesthetised by a UK Home Office (UKHO) personal license holder (GHC) in

135    a 1 mL L$^{-1}$ bath of 2-phenoxyethanol prior to collection of a 1 cm$^2$ tissue sample from the lower-

136    caudal fin, and wounds were treated with a mixture of adhesive powder (Orahesive) and antibiotic

137    (Cicatrin) (Moore *et al.* 1990). Tissue samples were immediately placed in ≥95% ethanol, and

138    stored at -20$^{o}$C. DNA was extracted from 2–4 mm$^2$ of each tissue sample using either the Gentra

139    Puregene DNA isolation kit or the DNeasy DNA purification kit (both Qiagen, Hilden, Germany).

140    For the RADseq library, DNA was quantified using the Quant-iT™ PicoGreen® dsDNA Assay kit

141    (Invitrogen) and normalised to concentrations ≥50 ng ml$^{-1}$. Gel electrophoresis was then used to

142    check that DNA extractions contained high molecular weight DNA.

143

*Molecular markers and methods*

145    Three types of molecular markers were used in this study. Mitochondrial DNA sequencing was

146    used to identify highly distinct lineages and to date the divergence between them through

147    phylogenetic analysis. Two sets of nuclear markers; microsatellites and RADseq-derived SNPs,

148    were used to investigate more recent and complex structure in a population genetics framework and

149    to compare the relative power of each marker to do so.

150

*Mitochondrial DNA amplification*

A total of 83 *C. carassius* individuals, randomly chosen from a subset of 30 populations, which

were chosen to represent all major catchment areas and the widest possible geographic range (min.

n = 1, max. n = 4, mean n = 2.7), were sequenced at the cytochrome b (*cyt*b) gene (Table 1). PCR

reactions were carried out following the protocol in Takada *et al.* (2010) using the forward and

reverse primers L14736-Glu and H15923-Thru on an Applied Biosciences® Veriti Thermal Cycler.

PCR products were sequenced in both directions on an ABI3700 by Macrogen Europe. The forward

and reverse *cyt*b sequence reads were aligned using a GenBank sequence from the UK (accession

no. JN412539, Table 1) as a reference and ambiguous nucleotides were manually edited using

CodonCode aligner v.2.0.6 (CodonCode Corporation).

161

*Microsatellite amplification*

Of the 867 samples identified as pure *C. carassius,* 19 samples were in populations with sample

numbers which were too low to be useful for population genetics analyses (< 4). The remaining

848 samples, from 49 populations, were genotyped at 13 microsatellite loci, including the six

species diagnostic loci used for hybrid identification (Supporting Information (SI) Table 1).

Microsatellites were amplified in three multiplex PCR reactions, using the Qiagen multiplex PCR

mix with manufacturer's recommended reagent concentrations, including Q solution and 1 μl of

template DNA. Primer concentrations for each locus are provided in SI Table 1 and PCRs were

performed on an Applied Biosciences® Veriti Thermal Cycler. The annealing temperature used was

54°C for all reactions, and all other PCR cycling parameters were set to Qiagen multiplex kit

recommended values. PCR products were run on a Beckman Coulter CEQ 8000 genome analyser

173 using a 400 bp size standard and microsatellite alleles scored using the Beckman Coulter CEQ8000

174 software.

175

176 *RADseq*

177 A total of 160 individuals (18 populations, min. $n = 8$, max. $n = 10$, mean $n = 8.9$), identified as

178 pure *C. carassius* with the diagnostic microsatellites, were used in the RADseq (Table 1). These

179 samples were chosen to represent a wide geographic range and all major phylogeographic clusters

180 identified using the microsatellite data. These samples were split across 13 libraries prepared at

181 Edinburgh Genomics (University of Edinburgh, UK) according to the protocol in Davey *et al.*

182 (2012) using the enzyme *Sbf*1. Libraries were then sequenced using paired-end sequencing across

183 five lanes of two Illumina HiSeq 2000 flowcells (Edinburgh Genomics).

184

185 *Data analyses*

186 *Phylogenetic analysis of mtDNA*

187 In addition to the 83 sequenced samples (SI Table 2), we retrieved 19 published *C. carassius* and

188 three *C. carpio cyt*b sequences from GenBank to be used as an outgroup. The *C. carpio* samples

189 were chosen to include samples from multiple, distant lineages of *C. carpio* located in Japan,

190 Greece and India. All sequences used were validated through cross checking with their original

191 publications (Table 1). Sequence alignment was performed in MEGA6 (Tamura *et al.* 2013) using

192 default settings, and DNAsp v.5.0 (Librado & Rozas 2009) was used to calculate sequence

193 divergence and to identify haplotypes.

194

195 Haplotypes of all *C. carassius* samples and the three *C. carpio* outgroup individuals were exported

196 to BEAST v.1.7.5 (Drummond *et al.* 2012) for phylogenetic analyses in order to identify the major

197    phylogenetic lineages within European *C. carassius*. Phylogenetic model testing with jModeltest2

198    v.2.1.7 (Guindon *et al.* 2003; Darriba *et al.* 2012) using Akaike information criterion (AIC),

199    Bayesian Information Criteria (BIC) and the decision-theoretic performance-based (DT) approach

200    showed that HKY (Hasegawa *et al.* 1985) was the most appropriate substitution model for our

201    dataset. Using this model, the splits between the major phylogenetic clades were then dated using a

202    relaxed molecular clock method in BEAST. The widely-used Dowling *et al.* (2002) cyprinid *cyt*b

203    divergence rate of 1.05% pairwise sequence divergence / MY was used after converting to a per

204    lineage value of 0.0053 mutations/site/MY for use in BEAST. We used a 'coalescent: constant size'

205    tree prior, which assumes an unknown but constant population size backwards in time, as

206    recommended for intraspecific phylogenies (Drummond *et al.* 2012) . MCMC chain lengths were 1

207    x $10^7$ with samples taken every 1000 iterations. A gamma site heterogeneity model was used, with

208    the default of four categories. Substitution rates, rate heterogeneity and base frequencies were

209    unlinked between each codon position to allow substitution rate to vary between them. Default

210    values were used for all other parameters and priors.

211

212    *Population structure and diversity analyses using microsatellites*

213    Allele dropout and null alleles in the microsatellite data were tested using Microchecker (Van

214    Oosterhout *et al.* 2004). FSTAT v. 2.9.3.2 (Goudet 1995) was then used to check for linkage

215    disequilibrium (LD) between loci (using 10,000 permutations), deviations from Hardy-Weinberg

216    equilibrium (HWE) within populations (126500 permutations) and for all population genetic

217    summary statistics. Genetic diversity within populations was estimated using Nei's estimator of

218    gene diversity ($H_o$) (Nei 1987) and Allelic richness ($A_r$), which was standardised to the smallest

219    sample size (n = 4) using rarefaction (Petit *et al.* 1998). Pairwise $F_{ST}$ values were calculated

220    according to (Weir & Cockerham 1984) and 23520 permutations and sequential Bonferroni

221    correction were used to test for significance of $F_{ST}$.

222

223 IBD was investigated using a Mantel test in the adegenet v1.6 (Jombart & Ahmed 2011) package in

224 R v3.0.1 (R Core Team 2013). We then tested for an association between $A_r$ and longitude and

225 latitude, which is predicted under a stepping-stone colonisation model (Ramachandran *et al.* 2005;

226 Simon *et al.* 2014), using linear regression analysis in R.

227

228 Population structure was then further examined using Discriminant Analyses of Principal

229 Components (DAPC) also in adegenet (DAPC, see SI text and Jombart *et al.* 2010 for more details).

230 DAPC has been shown to perform as well or better than the commonly used program,

231 STRUCTURE (Pritchard *et al.* 2000) for both simple and complex models of population structure

232 (Jombart *et al.* 2010). Furthermore, unlike STRUCTURE, DAPC is free of underlying assumptions

233 of Hardy-Weinberg equilibrium, which are likely to be violated when effective population sizes are

234 small, as is often the case in *C. carassius* (Hänfling *et al.* 2005).

235

236 In preliminary DAPC analysis using all 49 *C. carassius* populations, Sweden (SWE9) was found to

237 be so genetically distinct from the rest of the data set that it masked the variation between the other

238 populations. This population was therefore omitted from further DAPC analyses. To infer the

239 appropriate number of genetic clusters in the data, we used BIC scores (SI Fig. 5a), in all cases

240 choosing lowest number of genetic clusters from the range suggested. Spline interpolation

241 (Hazewinkel 1994) was then used to identify the appropriate number of principal components to use

242 in the subsequent discriminant analysis (SI Fig. 5a).

243

244 *RADseq data filtering and population structure analysis*

245 The quality of the RADseq raw read data was examined using FastQC (Andrews 2010), the dataset

246 was then cleaned, processed and SNPs were called using the Stacks pipeline v 1.19 (Catchen *et al.*

247 2011). Preliminary tests were carried out in order to identify optimal Stacks parameters (See SI

248 text). Final parameter values for the respective Stacks module were as follows; ustacks: M=2, m=8,

249   removal (-r) and deleveraging (-d) algorithms were also used; cstacks: N=2 (n populations = 18, n

250   individuals = 160); populations module: one SNP per RAD locus was used (--write_single_snp) and

251   SNPs were only retained if they were present in 70% of individuals (r=0.7) in at least 17 out of the

252   18 populations in the study (p=17), which allows for mutations in restriction sites that may cause

253   loci to dropout in certain lineages. All other parameters were kept at default values. Finally, we

254   filtered out loci which had a heterozygosity of > 0.5 and $F_{IS} < 0.0$ in one or more populations in

255   order to control for the possibility of erroneously merging ohnologs resulting from the multiple

256   genome duplications that have occurred the in *Cyprinus* and *Carassius* genera (Henkel *et al.* 2012;

257   Xu *et al.* 2014). The resulting refined SNP set was then used in subsequent phylogeographic

258   analyses. The R package Adegenet v. 1.42 was used to calculate $H_o$ and pairwise $F_{ST}$, test for IBD

259   and genetic clusters were inferred using DAPC.

260

261   *Reconstructing postglacial colonisation routes in Europe*

262   DIYABC v. 2.0 (Windows, Cornuet *et al.* 2014) was used to reconstruct the most likely *C.*

263   *carassius* recolonisation routes through Europe after the last glacial maximum. We used the

264   RADseq data set for this analysis as it showed a much clearer pattern of population structure than

265   the microsatellite data in DAPC analyses (see Results). Furthermore, preliminary DIYABC

266   analyses using microsatellites failed to identify a scenario which was significantly more likely than

267   its counterparts, suggesting low power in this dataset for reconstructing complex phylogeographic

268   patterns over long timescales.

269

270   As DIYABC is a computationally intensive method, it was necessary to perform analyses on a

271   subset of 1000 randomly-selected SNP loci from the full RADseq dataset to reduce computation

272   time. This SNP subset was first analysed with DAPC to confirm that it produced the same

273   population structure as the full dataset and was then used to compare the likelihood of a number of

274   user defined colonisation scenarios (i.e. a specific population tree topology, together with the

275    parameter prior distributions that are associated with it). First, 1 million datasets were simulated for

276    each scenario. These simulated summary statistic datasets represented the theoretical expectation

277    under each scenario, and were compared to the same summary statistics calculated from the

278    observed data, in order to identify the most likely of the tested scenarios. In DIYABC, two methods

279    of comparison between simulated and observed datasets are used; logistic regression and "direct

280    approach", the latter method identifies the scenario that produces the largest proportion of the *n*

281    number of closest scenarios to the observed, where *n* is specified by the user. The goodness-of-fit of

282    scenarios was also assessed using the model checking function implemented in DIYABC (Cornuet

283    *et al.* 2014). In all analyses, the single-sample summary statistics used were the mean and variance

284    of gene diversity across all polymorphic loci and the mean gene diversity across all loci. The two-

285    sample summary statistics used were the mean and variance of $F_{ST}$ and Nei's distance for loci with

286    $F_{ST}$ greater than zero between two samples and the mean $F_{ST}$ and Nei's distance for all loci. Finally,

287    for scenarios including admixture events, the maximum likelihood estimates of admixture

288    proportions were also used. See Cornuet *et al.* (2014) for the exact equations used and their

289    implementation in DIYABC.

290

291    To reduce the number and complexity of possible scenarios, we split DIYABC analysis into three

292    stages (Table 2). In stage 1, we tested 11 broad scale scenarios (Scenarios 1 -11, SI Fig. 1).

293    Populations were grouped into three pools in order to reduce the number and complexity of possible

294    scenarios (Table 2); Pool 1 – all northern European populations (npops = 17, *n* = 155), Pool 2 – Don

295    population (npops = 1, *n* = 9), Pool 3 – Danubian population (npops = 1, *n* = 6). In six scenarios (1,

296    2, 8-11), northern European and the Don population diverged from each other more recently than

297    from Danubian populations. These scenarios differ in the patterns of effective population size

298    change and the presence or absence of a bottleneck. In scenarios 3 and 4, northern European and

299    Danubian populations are more closely related to each other than to the Don population.  And in the

300    remaining three scenarios, one pool of populations is the product of an admixture event between the

301    other two. Population poolings and scenarios were both chosen on the basis of the broad

302    phylogeographic structure identified in the mtDNA and RADseq population structure analysis (see

303    Results).

304

305    In the second and third stages, we performed a finer scale analysis, focussing on the 17 northern

306    European populations alone. Populations were again pooled, this time into six groups, on the basis

307    of both population structure and geography (Table 2). In stage 2 we tested five scenarios (Scenarios

308    12-16, see SI Fig. 2a for graphical description of each scenario), with no bottlenecks included,

309    which represented the major topological variants that were most likely, given population structure

310    results from DAPC. We then identified the most likely of these scenarios in DIYABC and took this

311    forward into the final stage of the analysis where we tested 6 multiple bottleneck combinations (SI

312    Fig. 2b) around this scenario. This three stage approach allowed us to systematically build a

313    complex scenario for the European colonisation of *C. carassius*. Finally, we used the posterior

314    distributions of the time parameters, simulated using the scenario identified as most likely in stages

315    one and three, to estimate the times of the major lineage splits in European *C. carassius*. These

316    parameters, calculated by DIYABC in generations, were converted to years using an average

317    generation time of 2 years (Tarkan *et al.* 2010).

318

319    *Comparison of microsatellite and RADseq data*

320    Finally, we compared the results derived from population structure analyses on microsatellite and

321    RADseq data to assess their suitability for addressing our phylogeographic question. It is important

322    to note that differences between the full microsatellite and RADseq datasets could be attributable to

323    one or a combination of the following; the number of populations, the geographic distribution of

324    populations, the number of samples per population, the number of markers, or the information

325    content of the marker type. To disentangle these sources of variation, we created two microsatellite

326    data subsets; M2, which included only individuals used in RADseq, (excluding three individuals for

327    which microsatellite data was incomplete, $n$ = 146, npops = 19), and M3, which contained all

328    individuals for which microsatellite data was available in populations that were used in RADseq ($n$

329    = 313, npops = 19;

330    Table 3). This gave us three pairs of datasets for comparison: 1) RADseq Vs. M2: same individuals

331    but different marker types, 2) M1 vs M2: full microsatellite dataset versus a subset of the

332    populations, and 3) M2 vs M3: same populations but different number of individuals per

333    population. This strategy enabled us to test for the influence of marker, sampling of populations and

334    individuals per population respectively. Comparisons were performed between datasets on

335    heterozygosities and pairwise $F_{STS}$ using both Pearson's product-moment correlation coefficient and

336    paired Student's t-tests in R. IBD results were compared using Mantel tests (Jombart & Ahmed

337    2011), and DAPC results were compared on the basis of similarity of number of inferred clusters

338    and cluster sharing between populations.

339

340    # Results

341    *Phylogenetic analyses of mitochondrial data*

342    The combined 1090 bp alignment of 100 *cyt*b *C. carassius* mtDNA sequences yielded 22

343    haplotypes, which were split across two well supported and highly differentiated phylogenetic

344    lineages (Fig. 2, SI Table 3). Lineage 1 was found in all northern European river catchments

345    sampled, as well as eastern European (Dnieper) and southeastern European (Don and Volga)

346    catchments, whereas Lineage 2 was almost exclusively confined to the River Danube catchment.

347    There were, however, a few exceptions to this clear geographical split; two individuals, one from

348    the Elbe and one from the Rhine in northern Germany, belonged to mtDNA Lineage 2, as did one

349    individual from the River Lahn river catchment in western Germany. Also one population in the

350    Czech Republic, located on the border between the Danube and Rhine river catchments, was found

351    to contain individuals belonging to lineages 1 and 2.

352

353    The mean number of nucleotide differences within lineages 1 and 2 was 2.25 and 2.00, respectively,

354    which equated to a sequence divergence 0.2% and 0.18%, respectively. Between the two lineages

355    there was an average of 22.5 nucleotide differences (2.06% mean sequence divergence), with 19 of

356    these being fixed. BEAST molecular clock analysis dated the split between lineages 1 and 2 to be

357    1.30–3.22 million years ago (MYA), with a median estimate of 2.15 MYA (Fig. 2).

358

359    *Nuclear marker datasets and quality checking*

360    Microchecker showed no consistent signs of null alleles or allele dropout in microsatellite loci and

361    no significant LD was found between any pairs of loci. No populations showed significant deviation

362    from Hardy-Weinberg proportions (adjusted nominal level 0.0009).

363

364    After filtering raw RADseq data, *de novo* construction of loci across the 19 populations produced

365    35 709 RADseq loci that were present in at least 70% of individuals in at least 17 populations.

366    These loci contained a total of 29 927 polymorphic SNPs (approx. 0.84 SNPs per locus). Only the

367    first SNP in each RADseq locus was retained, to avoid confounding signals of LD. This yielded a

368    total of 18 908 loci with a mean coverage of 29.07 reads (SI Fig. 3b). Finally 5719 of these SNP

369    loci were filtered out due to high (> 0.5) heterozygosity and/or $F_{IS}$ of < 0.0 in at least one

370    population. In doing so, we removed many high coverage tags (SI Fig. 3a), which was consistent

371    with over-merged ohnologs having higher coverage (*i.e.* reads from more than two alleles) than

372    correctly assembled loci. The final dataset therefore contained 13189 SNP loci, with a mean

373    coverage of 27.72 reads.

374

375    *Within population diversity at nuclear loci*

376    Observed heterozygosity ($H_o$), averaged across all microsatellite loci within a population, ranged

377    from 0.06 (SWE9) to 0.44 (BLS), with a mean of 0.25 across all populations (SD = 0.105), and was

378    highly correlated with $A_r$ ($t$ = 19.67, $P$ < 0.001, df = 40), which ranged from 1.26 (FIN1) to 2.96

379   (POL3) with a mean of 1.92 (SD = 0.51). Mean $H_o$ averaged across all RADseq loci for all

380   populations was 0.013 (SD = 0.013), ranged from 0.001 to 0.057 and was significantly correlated

381   with $H_o$ from microsatellite loci at populations shared between both datasets (r = 0.69, $t$ = 3.74, $P$ =

382   0.002, df = 15). Microsatellite $A_r$ significantly decreased along an east to west longitudinal gradient

383   (adj. $R^2$ = 0.289, $P$ < 0.001, SI Fig. 4b) consistent with decreasing diversity along colonisation

384   routes. However, $A_r$ did not decrease with increasing latitude (adj $R^2$ =-0.007, P = 0.414, SI Fig. 4a).

385   We also repeated this analysis after removing samples from mtDNA Lineage 2 in the Danube

386   catchment. Again there was no relationship between $A_r$ and latitude ($R^2$ =-0.023, $P$ = 0.254, SI Fig.

387   4c), but the relationship between $A_r$ and longitude was strengthened (adj. $R^2$ = 0.316, $P$ < 0.001, SI

388   Fig. 4d).

389

390   *Population Structure in Europe based on nuclear markers*

391   Population structure was strong, as predicted. Using the full (M1) microsatellite dataset, mean

392   pairwise $F_{ST}$ was 0.413 (min = 0.0; BEL2 and BEL3), max = 0.864 (NOR2 vs GBR2), with 861 of

393   the 1128 pairwise population comparisons being significant $F_{ST}$ ($P$ < 0.05, SI Table 4). Pairwise $F_{ST}$

394   calculated from the RADseq dataset also showed strong structure (SI Table 5), ranging from 0.067

395   (DEN1, DEN2) to 0.699 (NOR2, GBR4), and these values were highly correlated with the same

396   population comparisons in the M3 microsatellite dataset (r = 0.66, $t$ = 9.01, $P$ < 0.01, df = 104).

397

398   BIC scores obtained from initial DAPC analyses of the microsatellite dataset, using all 49

399   populations, indicated that between 11 and 19 genetic clusters (SI Fig. 5a) would be an appropriate

400   model of the variation in the data. As a conservative estimate of population structure, we chose 11

401   clusters for use in the discriminant analysis, retaining eight principal components as recommended

402   by the spline interpolation a-scores (SI Fig. 5a). This initial analysis showed that populations

403   belonging to Cluster 10 (RUS1, Don river catchment) and Cluster 11 (GER3, GER4, CZE1,

404    Danubian catchment) were highly distinct from clusters found in northern Europe (Fig. 1b). Since

405    the marked genetic differentiation between these three main clusters masked the more subtle

406    population structure among northern European populations (see Fig. 1b), we repeated the DAPC

407    analysis without the populations from the Danube and Don (RUS1, GER3, GER4, CZE1, Fig. 1b).

408    The results of this second DAPC analysis revealed an IBD pattern of population structure, across

409    Europe (Fig. 1). Mantel tests excluding the Danubian and Don populations corroborated these

410    results; showing significant correlation with geographic distance in northern Europe (adjusted $R^2 =$

411    0.287, $P < 0.001$, SI Fig. 6a), with Danubian populations shown to be more diverged than their

412    geography would predict (data not shown).

413

414    In the RADseq DAPC analysis, BIC scores suggested between four and ten genetic clusters, a lower

415    number than that inferred from the microsatellite data set. Again we chose the lowest number of

416    suggested clusters (four) clusters to take forward in the analysis (SI Fig. 5b). Following spline

417    interpolation, we retained six principal components and kept two of the linear discriminants from

418    the subsequent discriminant analysis (SI Fig. 5b). The inferred population structure showed that the

419    Danubian population (HUN2) and the Don population (RUS1) were highly diverged from the

420    northern European clusters. Unfortunately, HUN2 is not present in the microsatellite dataset for

421    direct comparison, however both datasets, and the mtDNA data show the same pattern of high

422    divergence between northern Europe and Danubian populations. DAPC analyses of RADseq data

423    again showed an IBD pattern in northern European populations, which was confirmed with Mantel

424    tests when the Danubian population HUN2 was excluded (adjusted $R^2 = 0.722$, $P < 0.001$; SI Fig.

425    6b).

426

427    *Postglacial recolonisation of C. carassius in Europe*

428    DAPC results of the 1000 SNP RADseq dataset used in DIYABC showed that it produced the same

429    population structure as the full RADseq dataset (SI Fig. 7). For the broad-scale scenario tests in

430    stage one of the DIYABC analysis, both logistic regression and direct approach identified Scenario

431    9 as being most likely to describe the true broad-scale demographic history (SI Fig. 8). Model

432    checking showed that the observed summary statistics for our data fell well within those of the

433    posterior parameter distributions for scenario 9 (SI Fig. 8c). Scenario 9 agrees with the mtDNA

434    results, suggesting that the Danubian populations have made no major contribution to the

435    colonisation of northern Europe. The median posterior distribution estimate of the divergence time

436    between Danubian and northern European populations is 2.18 MYA (95% CI = 1.03 – 5.12 MYA),

437    assuming a two-year generation time (Tarkan *et al.* 2010)), which is strikingly similar to that of

438    mtDNA dating analysis. Scenario 9 also suggests that the northern European populations

439    experienced a population size decline after the split of Pool 1 from the population in the Don river

440    catchment, which lasted approximately 8920 years (95% CI = 616 – 13700 years) and reduced $N_e$

441    by 32%.

442

443    In stage two of the DIYABC analysis, we tested the major variant scenarios for the colonisation of

444    northern Europe. In assessing the relative probabilities of scenarios, there was some discrepancy

445    between the direct approach, which revealed Scenario 14 to be most likely, and the logistic

446    regression, which favoured Scenario 13 (with Scenario 14 being the second most likely). However,

447    the goodness-of-fit model checking showed that the observed dataset fell well within the posterior

448    parameter distributions for Scenario 14 (SI Fig. 9a), but not for Scenario 13 (not shown). Therefore,

449    Scenario 14 was carried forward into stage three in which we tested six more scenarios (SI Fig. 2b)

450    to compare combinations of bottlenecks using the same population tree topology as in Scenario 14.

451    Direct approach, logistic regression and model checking all found scenario 14d to be the most likely

452    (SI Fig. 9b), we therefore accepted this as the scenario for the colonisation of *C. carassius* in

453     northern Europe (SI Fig. 9b). This scenario infers an initial split between two sub-lineages in

454     northern Europe approximately 33 600 YBP (Fig. 4), one of which re-colonised northwest Europe

455     and one that re-colonised Finland through the Ukraine and Belarus. Scenario 14d also inferred a

456     secondary contact between these sub-lineages approximately 15 940 YBP, resulting in the

457     populations currently present in Poland; these admixed populations provided the source of one

458     colonisation across the Baltic into Sweden, and a second route was inferred into southern Sweden

459     from Denmark (Table 3, SI Fig. 9b).

460

461     *Comparing microsatellite datasets and RAD sequencing data*

462     The results from the RADseq ($n$ = 149, npops = 16) dataset and the full microsatellite dataset (M1,

463     $n$ = 848, npops = 49) largely agreed on the inferred structure and cluster identity of populations.

464     However, there were some important differences between them. Firstly, the IBD pattern of

465     population structure in northern Europe was much stronger in the RADseq data ($R^2$ = 0.722, $P$ <

466     0.001, SI Fig. 6) compared to the M1 dataset ($R^2$ = 0.287, $P$ < 0.001, excluding Danubian

467     populations and SWE9 from both datasets,  SI Fig. 6). Secondly, clusters inferred by the RADseq

468     DAPC analysis are much more distinct, *i.e.* there is much lower within-cluster, and higher between-

469     cluster variation in the RADseq results than in the M1 dataset results (Fig. 3).

470

471     As the properties of the RADseq and M1 datasets differ in four respects, namely marker type,

472     number of populations, number of samples per population (Table 3) and uniformity of sampling

473     locations, (SI Fig. 10) it was not possible to identify the cause of discrepancies in their results.

474     Therefore, below we report the results from the pair-wise dataset comparisons, which isolate the

475     effects of these parameter differences.

476

477　1) *M1 Vs. M3:* the effect that the number of populations and the uniformity of sampling locations

478　might have on inferred population structure. The geographic distribution of sampling locations was

479　more clustered in M1 (full microsatellite dataset) than in M3 (containing microsatellite for samples

480　in populations used in RADseq (SI Fig. 10), and IBD patterns were considerably stronger in the M3

481　subset (adj. $R^2$ = 0.447, $P$ < 0.001) than in the full M1 dataset (adj. $R^2$ = 0.287, $P$ < 0.001). In

482　contrast DAPC results were very similar between datasets, with inferred cluster number, structure

483　and population identity of clusters generally agreeing well (Fig. 1, Fig. 3c).

484

485　*2) M2 Vs. M3*: the effect of reducing the number of samples per population on the inferred

486　population structure. The number of samples per population in the M2 subset (microsatellite data

487　only for the samples used in RADseq, mean = 9.125 ± 0.8) was significantly lower than that of the

488　M3 subset (mean, 19.6 ± 9.0, $t$ = -4.66, df = 15, $P$ < 0.001), as was the number of alleles per

489　population (M2 mean = 24.4 ± 7.3, M3 mean = 27.4 ± 8.1, $t$ = -5.72, df = 15, $P$ < 0.001). Population

490　heterozygosities were significantly different between M2 and M3 (M2 mean = 0.21, M3 mean =

491　0.23, $t$ = -2.4, df = 15, $P$ = 0.012), but highly correlated (r = 0.94, $t$ = -11.13, $P$ < 0.001, df = 15).

492　Pairwise $F_{ST}$s were very strongly correlated (r = 0.97, $t$ = 46.26, $P$ < 0.001, df = 105), but again, still

493　significantly different between the two datasets (M2 mean = 0.46, M3 mean = 0.49 , $t$ = -6.21, $P$ <

494　0.001, df = 15, Table 4). The patterns of IBD were almost identical for M2 ($R^2$ = 0.455, $P$ < 0.001)

495　and M3 ($R^2$ = 0.447, $P$ < 0.001,  SI Fig. 6) and population structure inferred by DAPC was again

496　similar. BIC scores suggested a similar range of cluster number for M2 and M3, the smallest of

497　which was nine in both cases.

498

499　*3) RADseq Vs. M3:* The effect of the number and the type of markers used on the phylogeographic

500　results. We compared the results from the RADseq and M2 datasets, which contain exactly the

501　same samples (with the exception of three individuals missing in M2). Significant correlations were

502　again found between heterozygosities estimated for the two datasets (r = 0.69, $t$ = 3.73, $P$ = 0.002,

503   df = 15) and pair-wise $F_{ST}$s (r = 0.70, $t$ = 10.09, $P$ < 0.001, df = 105), but RADseq data yielded

504   much lower pairwise $F_{ST}$s (mean RAD = 0.29, mean M2 = 0.46, $t$ = 13.74, $P$ < 0.001, df = 15).

505   DAPC analysis of RADseq data resolved populations into much more distinct clusters (Figs. 3a,

506   3b), and the IBD pattern found was considerably stronger in the RADseq ($R^2$ = 0.722, $P$ < 0.001)

507   dataset compared to M2 ($R^2$ = 0.455, $P$ < 0.001,  SI Fig. 6).

508

# Discussion

510   In this study, we aimed to simultaneously produce a phylogeographic framework on which to base

511   conservation strategies for *C. carassius* in Europe, and compare the relative suitability of genome-

512   wide SNP markers and microsatellite markers for such an undertaking. Through comparison of the

513   inferred population structure from microsatellite and genome-wide SNP data, we show that there

514   are important differences in the results from each data type, attributable predominantly to marker

515   type, rather than within population sampling or spatial distribution of samples. However, despite

516   these differences, all three data types used (mitochondrial, microsatellite and SNP data) agree that,

517   unlike many other European freshwater fish for which phylogeographic data is available, *C.*

518   *carassius* has not been able to cross the Danubian catchment boundary into northern Europe. This

519   has resulted in two, previously unknown, major lineages of *C. carassius* in Europe, which we argue

520   should be considered as separate conservation units.

521

*Phylogeography and postglacial recolonisation of C. carassius in Europe*

523   The most consistent result across all three marker types (mtDNA sequences, microsatellites and

524   RADseq) was the identification of two highly-divergent lineages of *C. carassius* in Europe. The

525   distinct geographic distribution of these lineages; Lineage 1 being widely distributed across north

526   and eastern Europe and Lineage 2 generally only in the River Danube catchment, indicates a long-

527    standing barrier to gene flow between these geographic regions. Bayesian inference based on

528    mtDNA phylogeny and ABC analysis of RADseq data showed remarkable agreement, estimating

529    that these lineages have been isolated for 2.15 MYA (95% CI = 1.30–3.22) and 2.18 (95% CI = 2 –

530    6.12) MYA respectively, which firmly places the event at the beginning of the Pleistocene (2.6

531    MYA; (Gibbard & Head 2009). This pattern differs substantially from the general phylogeographic

532    patterns observed in other European freshwater fish. Indeed, previous studies have shown that the

533    Danube catchment has been an important source for the postglacial recolonisation of freshwater fish

534    into northern Europe or during earlier interglacials in the last 0.5 MYA. For example, bullhead

535    *Cottus gobio* (Hänfling & Brandl 1998; Hänfling *et al.* 2002), chub *Leuciscus cephalus* (Durand *et*

536    *al.* 1999), Eurasian perch *Perca fluviatilis* (Nesbø *et al.* 1999), riffle minnow *Leuciscus souffia*

537    (Salzburger *et al.* 2003), grayling *Thymallus thymallus* (Gum *et al.* 2009), European barbel *Barbus*

538    *barbus* (Kotlík & Berrebi 2001), and roach *Rutilus rutilus* (Larmuseau *et al.* 2009) all crossed the

539    Danube catchment boundary into northern drainages such as those of the rivers Rhine, Rhône and

540    Elbe during the mid-to-late Pleistocene. The above species occur in lotic habitats, and most are

541    capable of relatively high dispersal. In contrast *C. carassius* has a very low propensity for dispersal,

542    and a strict preference for the lentic backwaters, isolated ponds and small lakes (Holopainen *et al.*

543    1997; Culling *et al.* 2006; Copp 1991). We therefore hypothesise that these ecological

544    characteristics of *C. carassius* have reduced its ability to traverse the upper Danubian watershed,

545    which lies in a region characterised by the Carpathian Mountains and the Central European

546    Highlands. This region may have acted as a barrier to the colonisation of *C. carassius* into northern

547    European drainages during the Pleistocene. It should be noted, however, that the phylogeography of

548    two species, the spined loach *Cobitis taenia* and European weatherfish *Misgurnus fossilus*, does not

549    support this hypothesis as a general pattern for floodplain species (Janko *et al.* 2005; Culling *et al.*

550    2006). The former is the only species that we know of other than *C. carassius* showing long-term

551    isolation between the Danube and northern European catchments, but has lotic habitat preferences

552    and good dispersal abilities (Janko *et al.* 2005; Culling *et al.* 2006), whereas the latter inhabits

553   similar ecosystems as *C. carassius*, with low dispersal potential, but has colonised northern Europe

554   from the Danube catchment (Bohlen *et al.* 2006, 2007).

555

556   There is one notable exception to the strict separation between Danubian and northern European *C.*

557   *carassius* populations. The population CZE1, located in the River Lužnice catchment (Czech

558   Republic), which drains into the River Elbe, clusters with Danubian populations in both the

559   microsatellite and mtDNA data. This sample site, from the River Lužnice, is very close to the

560   Danubian catchment boundary and is situated in a relatively low lying area. Therefore, some recent

561   natural movements across the watershed between these river catchments, either through river

562   capture events or ephemeral connections, could have been possible. A similar pattern has been

563   shown in some European bullhead *Cottus gobio* populations along the catchment Danube/Rhine

564   catchment border (Riffel & Schreiber 1995). We also observed the presence of two mtDNA

565   haplotypes from Lineage 2 in some individuals from northern German populations (GER1, GER2,

566   GER8), however, one of these haplotypes was shared with Danubian individuals and the results

567   were not confirmed by nuclear markers. Overall this is most likely to be the result of occasional

568   human mediated long-distance dispersal for the purposes of intentional stocking.

569

570   Population structure within Lineage 1 is characterised by a pattern of IBD and a loss of allelic

571   richness from eastern to western Europe. This is consistent with the most likely colonisation

572   scenario identified by the DIYABC analysis, indicating a general southeast to northwest expansion

573   from the Ponto-Caspian region towards central and northern Europe (Fig. 4). The Ponto-Caspian

574   region, and in particular the Black Sea basin, was an important refugium for freshwater fishes

575   during the Pleistocene glacial cycles, and a similar colonisation route has been inferred for many

576   other freshwater species in northern Europe (Nesbø *et al.* 1999; Durand *et al.* 1999; Culling *et al.*

577   2006; Costedoat & Gilles 2009). The DIYABC analysis also suggests that there was an interval of >

578   200 000 years between the split of the Don population (≈ 270 000 years ago) and the next split in

579   the scenario (approx. 33 600 years ago), which marks the main expansion across central and

580   northern Europe. It appears that no further population divergence can be dated back to the time

581   interval between the Riss/Saalian and the Würm/Weichelian glacial periods. This may be because

582   the range of *C. carassius* has not undergone a major change during that time interval, but it is more

583   likely that the signal of expansion during the Riss-Würm interglacial has been eradicated through a

584   subsequent range contraction during the Würm/Weichelian glacial period. The model also suggests

585   that the Würm/Weichelian period was accompanied by a sustained but moderate reduction in

586   population size over almost 9000 years (Bottleneck A, Fig. 4), which may reflect general population

587   size reductions during the Riss glaciations or a series of shorter bottlenecks during subsequent range

588   expansion (Ramachandran et al. 2005, Simon et al 2015, Hewitt 2000).

589

590   DIYABC analyses inferred the colonisation of northern Europe by two sub-lineages within the

591   mtDNA Lineage 1, which were isolated from each other approximately 33 600 years ago. These

592   sub-lineages may reflect two glacial refugia resulting from the expansion of the Weichselian ice cap

593   to its maximum extent roughly 22 000 years ago (see hypothetical refugia II and III in Fig. 4). The

594   western sub-lineage underwent a second long period of population decline (Bottleneck B, Fig. 4),

595   which may again represent successive founder effects during range expansion. There is then

596   evidence of secondary contact between these sub-lineages (node b, approximately ≈ 15 940 years

597   ago), contributing to the genetic variation now found in Poland. This inferred admixture event may

598   represent one of the numerous inundation and drainage capture events, which resulted from the

599   melting of the Weichselian ice cap, that are known to have occurred around this time (Grosswald

600   1980; Gibbard *et al.* 1988; Arkhipov *et al.* 1995). However, as the colonisation of Europe was

601   likely to have occurred via the expansion of colonisation fronts (*i.e.* dashed contour lines in Fig. 4),

602   rather than along linear paths, it could also be indicative of the known IBD gradient between the

603   inferred western and eastern sub-lineages. Such a gradient (eg. between northwestern and

604 northeastern Europe) may give false signals of admixture between intermediate populations, such as

605 those in Poland.

606

607 The colonisation of the Baltic sea basin also seems to have been complex, with three independent

608 routes inferred by DIYABC scenario 14d; one recent route through Denmark into southern Sweden,

609 one to the east of the Baltic Sea, through Finland, and one across the Baltic Sea, from populations

610 related to those in Poland (Pool 4). The first of these agrees well with the findings of Janson *et al*.

611 (2014), whereby populations, including SWE8 from our study (SK3P in Janson *et al.* 2014), in this

612 region were found to be distinct from those in central Sweden. The eastern route shows similarities

613 to the colonisation patterns of *P. fluvilatilis*, which is hypothesised to have had a refugium east of

614 Finland (Nesbø *et al.* 1999) during the most recent glacial period. This is certainly also plausible in

615 *C. carassius* and may account for the distinctiveness of Finnish populations seen in microsatellites

616 and RADseq DAPC analysis. The last colonisation route, across the Baltic Sea from mainland

617 Europe, may have coincided with the freshwater Lake Ancylus stage of the Baltic Sea's evolution,

618 which existed from ≈ 10 600 to 7 500 years ago (Björck 1995; Kostecki 2014). The Lake Ancylus

619 stage likely provided a window for the colonisation of many of the species now resident in the

620 Baltic, and has been proposed as a possible window for the colonisation of *T. thymallus* (Koskinen

621 *et al.* 2000), *C. taenia*, (Culling *et al.* 2006), *C. gobio* (Kontula & Väinölä 2001) and four

622 *Coregonus* species (Svärdson 1998). Consistent with this, we found strong similarity between

623 populations from Fasta Åland, southern Finland and central Sweden, suggesting that shallow

624 regions in the central part of Lake Ancylus (what is now the Åland Archipelago), may have

625 provided one route across Lake Ancylus.

626 It is also likely that the contemporary distribution of *C. carassius* in the Baltic has been influenced

627 by human translocations. *C. carassius* were often used as a food source in monasteries in many

628 parts of Sweden (Janson *et al.* 2014), and the Baltic island of Gotland (Rasmussen 1959; Svanberg

629 *et al.* 2013) was an important trading port of the Hanseatic League – a commercial confederation

630    that dominated trade in northern Europe from the 13<sup>th</sup> to 17<sup>th</sup> centuries. Previous data suggest that

631    *C. carassius* was transported from the Scania Province, southern Sweden, where *C. carassius*

632    aquaculture was common at least during the 17<sup>th</sup> century, to parts further north  (Svanberg *et al.*

633    2013; Janson *et al.* 2014).

634

635    *Implications for the conservation of* C. carassius *in Europe*

636    The two *C. carassius* lineages exhibit highly-restricted gene flow between them and are the highest

637    known organisational level within the species. They therefore meet the genetic criteria for

638    Evolutionarily Significant Units (ESUs) as described in (Fraser & Bernatchez 2001). This is

639    especially important in light of the current *C. carassius* decline in the Danubian catchment

640    (Bănărescu 1990; Navodaru *et al.* 2002; Lusk *et al.* 2010; Savini *et al.* 2010). The conservation of

641    *C. carassius* in central Europe must therefore take these catchment boundaries into consideration, as

642    opposed to political boundaries. A first step would be to include *C. carassius* in Red Lists, not only

643    for individual countries, but at the regional (e.g. European Red List of Freshwater Fishes; (Freyhof

644    & Brooks 2011) and global (IUCN 2015) scales, and we hope that the evidence presented here will

645    facilitate this process. Within the northern European lineage, the Baltic Sea basin shows high levels

646    of population diversity, likely owing to its complex colonisation history. As such, the Baltic

647    represents an important part of the *C. carassius* native range. Although *C. carassius* is not currently

648    thought to be threatened in the Baltic region, *C. gibelio* is invading this region and is considered a

649    threat (Urho & Lehtonen; Deinhardt 2013).

650

*Microsatellites vs RADseq for phylogeography*

652  Broad conclusions drawn from each of our RADseq-derived SNPs, full or partial microsatellite

653  datasets are consistent, demonstrating deep divergence between northern and southern European

654  populations and an IBD pattern of population structure in northern Europe. This similarity in spatial

655  signal between marker types was also observed by (Bradbury *et al.* 2015). However, two striking

656  differences exist in the phylogeographic results produced by RADseq compared to those of the

657  microsatellite datasets. Firstly, the IBD pattern inferred from RADseq data was considerably

658  stronger than for any of the microsatellite datasets. This effect was also found by Coates *et al*.

659  (2009) when comparing SNPs and microsatellites, who postulated that it was driven by the

660  differences in mutational processes of the markers. The second major difference between RADseq

661  and microsatellite results was that clusters inferred by DAPC from the RADseq data were

662  considerably more distinct compared to the full microsatellite dataset, emphasising the fine scale

663  structure in the data (which is particularly apparent in the northern Finnish populations). We ruled

664  out the possibility of these differences being caused by the reduction in number of populations, their

665  spatial uniformity or number of individuals per population used in RADseq by creating two partial

666  microsatellite datasets and comparing these to results from the RADseq-SNPs. Differences between

667  marker types were consistently reproducible whether full or partial microsatellite datasets were used

668  in the analyses.

669

670  It is also worth noting that the number of populations or the number of samples per population had

671  no apparent impact on IBD and DAPC results between the microsatellite datasets. This is in contrast

672  to predictions of patchy sampling of IBD made by Schwartz and McKelvey (2009), perhaps

673  because of the strong population structure in *C. carassius,* and likelihood that a sufficiently

674  informative number of populations was included even in the reduced datasets.

675

676    SNP loci provide several advantages over microsatellites additional to those highlighted here. SNPs

677    are more densely and evenly distributed across the genome (Xing *et al.* 2005) and have been shown

678    to display lower error rates during genotyping (Montgomery *et al.* 2005). For example, Morin *et al.*

679    (2009a) showed that HW proportions are very sensitive to microsatellite genotyping errors. SNPs

680    also lend themselves to a plethora of evolutionary applications, including the identification of

681    outlier loci (Hohenlohe *et al.* 2012) or small regions of introgression in the genome (Hohenlohe *et*

682    *al.* 2013). Lastly, SNPs are also much less susceptible to homoplasy than microsatellites (Morin *et*

683    *al.* 2004). Van Oppen *et al.* (2000) found evidence of homoplasy in 10 out of 13 microsatellite loci,

684    which had accumulated in approximately 700,000 years and Cornuet *et al.* (2010) show that such

685    homoplasy makes microsatellites unreliable and error prone when used in DIYABC for inference

686    over long time scales. For these reasons, SNPs have a clear advantage over microsatellites for the

687    purposes of characterising population divergence over long time scales. This may explain why

688    preliminary microsatellite analyses in DIYABC showed insufficient power to identify a most likely

689    colonisation scenario.

690

691    *Conclusions*

692    We have identified the most likely routes of post-glacial colonisation in *C. carassius*, which deviate

693    from the general patterns observed in other European freshwater fishes. This has resulted in two,

694    previously-unidentified major lineages in Europe, which future broad-scale monitoring and

695    conservation strategies should take into account.

696

697    Although our RADseq sampling design included only 17.6% of samples included in the full

698    microsatellite dataset this was sufficient to produce a robust phylogeography in agreement with the

699    microsatellite dataset, and emphasised the fine scale structure among populations. We therefore

700   conclude that, if made to choose between the comprehensively sampled microsatellite approach or

701   the RADseq approach with fewer samples but many more loci, the RADseq approach presents the

702   better option for the phylogeography of *C. carassius,* with the huge number of SNP loci

703   overcoming the limitations imposed by reduced sample number. We also predict that this will hold

704   true for systems with similar genetic characteristics to ours, *i.e*. strong population structure

705   characterised by IBD.

706

722

723   **References**

724   Andrews S (2010) *FastQC: a quality control tool for high throughput sequence data*. Babraham

725   Bioinformatics. Available at www.bioinformatics.babraham.ac.uk/projects/fastqc.

726   Arkhipov SA, Ehlers J, Johnson RG, Wright HE Jr (1995) Glacial drainage towards the

727   Mediterranean during the Middle and Late Pleistocene. *Boreas*, **24**, 196–206.

728   Balkenhol N, Landguth EL (2011) Simulation modelling in landscape genetics: on the need to go

729   further. *Molecular ecology*, **20**, 667–670.

730   Bernatchez L, Wilson CC (1998) Comparative phylogeography of Nearctic and Palearctic fishes.

731   *Molecular ecology*, **7**, 431–452.

732   Bianco P (1990) Potential role of the palaeohistory of the Mediterranean and Paratethys basins on

733   the early dispersal of Euro-Mediterranean freshwater fishes. *Ichthyological exploration of*

734   *freshwaters*, **1**, 167 - 184

735   Björck S (1995) A review of the history of the Baltic Sea, 13.0-8.0 ka BP. *Quaternary*

736   *international: the journal of the International Union for Quaternary Research*, **27**, 19–40.

737   Bohlen J, Šlechtová V, Bogutskaya N, Freyhof J (2006) Across Siberia and over Europe:

738   Phylogenetic relationships of the freshwater fish genus Rhodeus in Europe and the phylogenetic

739   position of *R. sericeus* from the River Amur. *Molecular phylogenetics and evolution*, **40**, 856–865.

740   Bohlen J, ŠLechtová V, Doadrio I, Ráb P (2007) Low mitochondrial divergence indicates a rapid

741   expansion across Europe in the weather loach, *Misgurnus fossilis* (L.). *Journal of fish biology*, **71**,

742   186–194.

743   Bănărescu P (1990) *Zoogeography of Fresh Waters. Vol. 1. General Distribution and Dispersal of*

744   *Freshwater Animals*. Aula-Verlag, Wiesbaden.

745 Bănărescu P (1992) *Zoogeography of fresh waters. Vol. 2. Distribution and dispersal of freshwater*

746 *animals in North America and Eurasia*. Aula-Verlag, Wiesbaden.

747 Bradbury IR, Hamilton LC, Dempson B *et al.* (2015) Transatlantic secondary contact in Atlantic

748 Salmon, comparing microsatellites, a single nucleotide polymorphism array and restriction-site

749 associated DNA sequencing for the resolution of complex spatial structure. *Molecular Ecology* **24**,

750 5130–5144.

751 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013) Stacks: an analysis tool set

752 for population genomics. *Molecular ecology*, **22**, 3124–3140.

753 Coates BS, Sumerford DV, Miller NJ *et al.* (2009) Comparative performance of single nucleotide

754 polymorphism and microsatellite markers for population genetic analysis. *The Journal of heredity*,

755 **100**, 556–564.

756 Copp GH (1991) Typology of aquatic habitats in the great ouse, a small regulated lowland river.

757 *Regulated Rivers: Research & Management*, **6**, 125–134.

758 Copp G, Sayer C (2010) *Norfolk Biodiversity Action Plan–Local Species Action Plan for Crucian*

759 *Carp (Carassius carassius). Norfolk Biodiversity Partnership Reference: LS⁄ 3*. Fisheries &

760 Aquaculture Science, Lowestoft.

761 Copp G, Tarkan S, Godard M, Edmonds N, Wesley K (2010) Preliminary assessment of feral

762 goldfish impacts on ponds, with particular reference to native crucian carp. *Aquatic invasions*, **5**,

763 413–422.

764 Copp GH, Černý J, Kováč V (2008) Growth and morphology of an endangered native freshwater

765 fish, crucian carp *Carassius carassius*, in an English ornamental pond. *Aquatic conservation:*

766 *marine and freshwater ecosystems*, **18**, 32–43.

767  Cornuet JM, Luikart G (1996) Description and power analysis of two tests for detecting recent

768  population bottlenecks from allele frequency data. *Genetics*, **144**, 2001–2014.

769  Cornuet JM, Ravigne V, Estoup A (2010) Inference on population history and model checking

770  using DNA sequence and microsatellite data with the software DIYABC (v1.0). *BMC*

771  *Bioinformatics* **11**, 401.

772  Cornuet J-M, Pudlo P, Veyssier J *et al.* (2014) DIYABC v2.0: a software to make approximate

773  Bayesian computation inferences about population history using single nucleotide polymorphism,

774  DNA sequence and microsatellite data. *Bioinformatics*, **30,** 1187-1189.

775  Costedoat C, Gilles A (2009) Quaternary pattern of freshwater fishes in Europe: comparative

776  phylogeography and conservation perspective. *The Open Conservation Biology Journal*, **3**, 36-48.

777  Culling MA, Janko K, Boron A *et al.* (2006) European colonization by the spined loach (*Cobitis*

778  *taenia*) from Ponto-Caspian refugia based on mitochondrial DNA variation. *Molecular ecology*, **15**,

779  173–190.

780  Darriba D, Taboada GL, Doallo R, Posada D (2012) jModelTest 2: more models, new heuristics

781  and parallel computing. Nature Methods **9**, 772.

782  Davey JW, Cezard T, Fuentes-Utrilla P *et al.* (2012) Special features of RAD Sequencing data:

783  implications for genotyping. *Molecular ecology*, **22**, 3151–3164.

784  Deinhardt M (2013) The invasive potential of Prussian carp in Finland under the light of a novel

785  semi-clonal reproductive mechanism, Masters thesis, University of Jyväskylä.

786  Dowling TE, Tibbets CA, Minckley WL, Smith GR, McEachran JD (2002) Evolutionary

787  Relationships of the Plagopterins (Teleostei: Cyprinidae) from Cytochrome b Sequences. *Copeia*,

788  **2002**, 665–678.

789 Drummond AJ, Suchard MA, Xie D, Rambaut A (2012) Bayesian phylogenetics with BEAUti and

790 the BEAST 1.7. *Molecular biology and evolution*, **29**, 1969–1973.

791 Durand JD, Persat H, Bouvet Y (1999) Phylogeography and postglacial dispersion of the chub

792 (*Leuciscus cephalus*) in Europe. *Molecular ecology*, **8**, 989–997.

793 Eaton DAR (2014) PyRAD: assembly of de novo RADseq loci for phylogenetic analyses.

794 *Bioinformatics*, **30**, 1844–1849.

795 Emerson KJ, Merz CR, Catchen JM *et al.* (2010) Resolving postglacial phylogeography using high-

796 throughput sequencing. *Proceedings of the National Academy of Sciences*, **107**, 16196–16200.

797 Epperson BK, Mcrae BH, Scribner K *et al.* (2010) Utility of computer simulations in landscape

798 genetics. *Molecular ecology*, **19**, 3549–3564.

799 Frankham R, Briscoe D,  Ballou J (2002) *Introduction to Conservation Genetics*. Cambridge

800 University Press, Cambridge, UK.

801 Fraser DJ, Bernatchez L (2001) Adaptive evolutionary conservation: towards a unified concept for

802 defining conservation units. *Molecular ecology*, **10**, 2741–2752.

803 Freyhof J, Brooks E (2011) European red list of freshwater fishes. Luxembourg: Publications Office

804 of the European Union.

805 Gibbard P, Head MJ (2009) The Definition of the Quaternary System/Era and the Pleistocene

806 Series/Epoch. *Quaternaire*, **20**, 125–133.

807 Gibbard PL, Rose J, Bridgland DR (1988) The history of the great northwest European rivers

808 during the past three million years [and discussion]. *Philosophical transactions of the Royal Society*

809 *of London. Series B, Biological sciences*, **318**, 559–602.

810 Goudet J (1995) FSTAT, a computer program to calculate F-Statistics. *Heredity*, **86**, 485-486

811    Grosswald MG (1980) Late Weichselian ice sheet of Northern Eurasia. *Quaternary Research*, **13**,

812    1–32.

813    Gum B, Gross R, Geist J (2009) Conservation genetics and management implications for European

814    grayling, *Thymallus thymallus*: synthesis of phylogeography and population genetics. *Fisheries*

815    *management and ecology*, **16**, 37–51.

816    Guindon S, Gascuel O (2003) A simple, fast, and accurate algorithm to estimate large phylogenies

817    by maximum likelihood. Systematic Biology, **52**, 696–704.

818    Hasegawa M, Kishino H, Yano T (1985) Dating of the human-ape splitting by a molecular clock of

819    mitochondrial DNA. *Journal of molecular evolution*, **22**, 160–174.

820    Hauser L, Baird M, Hilborn R, Seeb LW, Seeb JE (2011) An empirical comparison of SNPs and

821    microsatellites for parentage and kinship assignment in a wild sockeye salmon (*Oncorhynchus*

822    *nerka*) population. *Molecular ecology resources*, **11 Suppl 1**, 150–161.

823    Hazewinkel M (Ed.) (1994) *Encyclopeidia of Mathematics (set)*. Kluwer, Dordrecht, Netherlands.

824    Henkel CV, Dirks RP, Jansen HJ *et al.* (2012) Comparison of the exomes of common carp

825    (*Cyprinus carpio*) and zebrafish (*Danio rerio*). *Zebrafish*, **9**, 59–67.

826    Hess JE, Matala AP, Narum SR (2011) Comparison of SNPs and microsatellites for fine-scale

827    application of genetic stock identification of Chinook salmon in the Columbia River Basin.

828    *Molecular ecology resources*, **11**, 137–149.

829    Hewitt GM (1999) Post-glacial re-colonization of European biota. *Biological journal of the Linnean*

830    *Society. Linnean Society of London*, **68**, 87–112.

831    Hoban SM, Gaggiotti OE, Bertorelle G (2013) The number of markers and samples needed for

832    detecting bottlenecks under realistic scenarios, with and without recovery: a simulation-based study.

833    *Molecular ecology*, **22**, 3444–3450.

834   Holopainen IJ, Aho J, Vornanen M, Huuskonen H (1997) Phenotypic plasticity and predator effects

835   on morphology and physiology of crucian carp in nature and in the. *Journal of fish biology*, **50**,

836   781–798.

837   Holopainen IJ, Oikari A (1992) Ecophysiological effects of temporary acidification. *Annales*

838   *Zoologici Fennici*, **29**, 29–38.

839   Holopainen IJ, Tonn WM, Paszkowski CA (1997) Tales of two fish: the dichotomous biology of

840   crucian carp (*Carassius carassius* (L.)) in northern Europe.

841   Hänfling B, Hellemans B, Volckaert F, Carvalho GR (2002) Late glacial history of the cold-adapted

842   freshwater fish *Cottus gobio*, revealed by microsatellites. *Molecular ecology,* **11,** 1717–1729.

843   Hänfling B, Bolton P, Harley M, Carvalho GR (2005) A molecular approach to detect hybridisation

844   between crucian carp (*Carassius carassius*) and non-indigenous carp species (*Carassius* spp. and

845   *Cyprinus carpio*). *Freshwater biology*, **50**, 403–417.

846   Hickley P, Chare S (2004) Fisheries for non-native species in England and Wales: angling or the

847   environment? *Fisheries management and ecology*, **11**, 203-212.

848   IUCN 2015. The IUCN Red List of Threatened Species. Version 2015-4, http://www.iucnredlist.org

849   . Downloaded on 19 August 2015.

850   Janko K, Culling MA, Ráb P, Kotlík P (2005) Ice age cloning--comparison of the Quaternary

851   evolutionary histories of sexual and clonal forms of spiny loaches (*Cobitis*; Teleostei) using the

852   analysis of mitochondrial DNA variation. *Molecular ecology*, **14**, 2991–3004.

853   Janson S, Wouters J, Bonow M, Svanberg I, Olsén KH (2014) Population genetic structure of

854   crucian carp (*Carassius carassius*) in man-made ponds and wild populations in Sweden.

855   *Aquaculture international: journal of the European Aquaculture Society*, **23**, 359-368.

856    Jombart T, Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-wide SNP data.

857    *Bioinformatics* , **27**, 3070–3071.

858    Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal components: a new

859    method for the analysis of genetically structured populations. *BMC genetics*, **11**, 94.

860    Kontula T, Väinölä R (2001) Postglacial colonization of Northern Europe by distinct

861    phylogeographic lineages of the bullhead, *Cottus gobio*. *Molecular ecology*, **10**, 1983–2002.

862    Koskinen MT, Ranta E, Piironen J *et al.* (2000) Genetic lineages and postglacial colonization of

863    grayling (*Thymallus thymallus*, Salmonidae) in Europe, as revealed by mitochondrial DNA

864    analyses. *Molecular ecology*, **9**, 1609–1624.

865    Kostecki R (2014) Stages of the Baltic Sea evolution in the geochemical record and radiocarbon

866    dating of sediment cores from the Arkona Basin. *Oceanological and Hydrobiological Studies*, **43**,

867    237–246.

868    Kotlík P, Berrebi P (2001) Phylogeography of the barbel (*Barbus barbus*) assessed by

869    mitochondrial DNA variation. *Molecular ecology*, **10**, 2177–2185.

870    Landguth EL, Fedy BC, Oyler-McCance SJ *et al.* (2012) Effects of sample size, number of markers,

871    and allelic richness on the detection of spatial genetic pattern. *Molecular ecology resources*, **12**,

872    276–284.

873    Larmuseau MHD, Freyhof J, Volckaert FAM, Van Houdt JKJ (2009) Matrilinear phylogeography

874    and demographical patterns of *Rutilus rutilus*: implications for taxonomy and conservation. *Journal

875    of fish biology*, **75**, 332–353.

876    Lelek A (1980) Threatened freshwater fishes of Europe. Council of Europe, Strasbourg, France.

877    Librado P, Rozas J (2009) DnaSP v5: a software for comprehensive analysis of DNA

878    polymorphism data. *Bioinformatics* , **25**, 1451–1452.

879   Luikart G, Cornuet J-M (2008) Empirical Evaluation of a Test for Identifying Recently

880   Bottlenecked Populations from Allele Frequency Data. *Conservation biology*, **12**, 228-237.

881   Lusk S, Hanel L, Luskova V (2004) Red List of the ichthyofauna of the Czech Republic:

882   Development and present status. *Folia Zoologica*, **53**, 215–226.

883   Lusk S, Lusková, V, Hanel L (2010) Alien fish species in the Czech Republic and their impact on

884   the native fish fauna. *Folia Zoology*, **59**, 57–72.

885   Mabuchi K, Senou H, Suzuki T, Nishida M (2005) Discovery of an ancient lineage of *Cyprinus*

886   *carpio* from Lake Biwa, central Japan, based on mtDNA sequence data, with reference to possible

887   multiple origins of koi. *Journal of Fish Biology,* **66**, 1516–1528.

888   Mezhzherin SV, Kokodii SV, Kulish AV, Verlatii DB, Fedorenko LV (2012) Hybridization of

889   crucian carp *Carassius carassius* (Linnaeus, 1758) in Ukrainian reservoirs and the genetic structure

890   of hybrids. *Cytology and genetics*, **46**, 28–35.

891   Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2006) Rapid and cost-effective

892   polymorphism identification and genotyping using restriction site associated DNA (RAD) markers.

893   *Genome research*, **17**, 240–248.

894   Moore A, Russell IC, Potter ECE (1990) The effects of intraperitoneally implanted dummy acoustic

895   transmitters on the behaviour and physiology of juvenile Atlantic salmon, *Salmo salar* L. *Journal of*

896   *fish biology*, **37**, 713–721.

897   Morin PA, Luikart G, Wayne RK, the SNP workshop group (2004) SNPs in ecology, evolution and

898   conservation. *Trends in ecology & evolution*, **19**, 208–216.

899   Morin PA, Pease VL, Hancock BL *et al.* (2010) Characterization of 42 single nucleotide

900   polymorphism (SNP) markers for the bowhead whale (*Balaena mysticetus*) for use in discriminating

901   populations. *Marine Mammal Science*, **26**, 716–732.

902 Mrakovčić M, Buj I, Mustafić P, Ćaleta M, Zanella D (2007) *Croatian Red List: Freshwater fish.*

903 Department of Zoology, Faculty of Science, Zagreb.

904 Navodaru I, Buijse AD, Staras M (2002) Effects of Hydrology and Water Quality on the Fish

905 Community in Danube Delta Lakes. *International review of hydrobiology*, **87**, 329–348.

906 Nei M (1987) Molecular evolutionary genetics. Columbia University Press, New York.

907 Nesbø CL, Fossheim T, Vøllestad LA, Jakobsen KS (1999) Genetic divergence and

908 phylogeographic relationships among European perch (*Perca fluviatilis*) populations reflect glacial

909 refugia and postglacial colonization. *Molecular ecology*, **8**, 1387–1404.

910 Van Oosterhout C, Hutchinson WF, Wills DPM, Shipley P (2004) micro-checker: software for

911 identifying and correcting genotyping errors in microsatellite data. *Molecular ecology notes*, **4**,

912 535–538.

913 Peery MZ, Kirby R, Reid BN *et al.* (2012) Reliability of genetic bottleneck tests for detecting recent

914 population declines. *Molecular ecology*, **21**, 3403–3418.

915 Petit RJ, El Mousadik A, Pons O (1998) Identifying Populations for Conservation on the Basis of

916 Genetic Markers. *Conservation biology: the journal of the Society for Conservation Biology*, **12**,

917 844–855.

918 Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus

919 genotype data. *Genetics*, **155**, 945–959.

920 R Core Team (2013) *R: a language and environment for statistical computing*, Vienna, Austria.

921 Ramachandran S, Deshpande O, Roseman CC *et al.* (2005) Support from the relationship of genetic

922 and geographic distance in human populations for a serial founder effect originating in Africa.

923 *Proceedings of the National Academy of Sciences of the United States of America*, **102**, 15942–

924 15947.

925   Rasmussen H (1959) Fish ponds and fish rearing. In: *Kulturhistoriskt lexikon för nordisk medeltid*

926   (eds: Andersson I, Granlund J)*,* pp. 307–309, Malmö, Allhem.

927   Reyjol Y, Hugueny B, Pont D *et al.* (2006) Patterns in species richness and endemism of European

928   freshwater fish. *Global ecology and biogeography*, **16**, 65–75.

929   Rylková K, Kalous L, Bohlen J, Lamatsch DK, Petrtýl M (2013) Phylogeny and biogeographic

930   history of the cyprinid fish genus Carassius (Teleostei: Cyprinidae) with focus on natural and

931   anthropogenic arrivals in Europe. *Aquaculture* , **380,** 13–20.

932   Salzburger W, Brandstätter A, Gilles A *et al.* (2003) Phylogeography of the vairone (Leuciscus

933   souffia, Risso 1826) in Central Europe. *Molecular ecology*, **12**, 2371–2386.

934   Savini D, Occhipinti-Ambrogi A, Marchini A *et al.* (2010) The top 27 animal alien species

935   introduced into Europe for aquaculture and related activities. *Journal of applied ichthyology*, **26**, 1–

936   7.

937   Sayer CD, Copp GH, Emson D *et al.* (2011) Towards the conservation of crucian carp Carassius

938   carassius: understanding the extent and causes of decline within part of its native English range.

939   *Journal of fish biology*, **79**, 1608–1624.

940   Schwartz MK, McKelvey KS (2009) Why sampling scheme matters: the effect of sampling scheme

941   on landscape genetic results. *Conservation genetics*, **10**, 441–452.

942   Simic, V, Simic S, Cirkovic M, Pantovic N (2009) *Preliminary red list of the fishes of Serbia*.

943   COMBAFF-First Conference on Conservation and Management of Balkan Freshwater Fishes.

944   Simon A, Gozlan RE, Robert Britton J, van Oosterhout C, Hänfling B (2014) Human induced

945   stepping-stone colonisation of an admixed founder population: the spread of topmouth gudgeon

946   (*Pseudorasbora parva*) in Europe. *Aquatic sciences*, **77**, 17–25.

947    Svanberg I, Bonow M, Olsén H (2013) Fish ponds in Scania, and Linnaeus's attempt promote

948    aquaculture in Sweden. In: *Svenska Linnésällskapets årsskrift* (eds David B, Gunnar D), pp. 85–

949    100. Svenska Linnésällskapet, Uppsala.

950    Svärdson G (1998) Plostglacial dispersal and reticulate evolution of Nordic Coregonids. *Nordic*

951    *journal of freshwater research*, **74**, 3–32.

952    Takada M, Tachihara K, Kon T *et al.* (2010) Biogeography and evolution of the *Carassius auratus*-

953    complex in East Asia. *BMC evolutionary biology*, **10**, 7.

954    Tamura K, Stecher G, Peterson D, Filipski A, Kumar S (2013) MEGA6: Molecular Evolutionary

955    Genetics Analysis version 6.0. *Molecular biology and evolution*, **30**, 2725–2729.

956    Tarkan AS, Cucherousset J, Zięba G, Godard MJ, Copp GH (2010) Growth and reproduction of

957    introduced goldfish *Carassius auratus* in small ponds of southeast England with and without native

958    crucian carp *Carassius carassius*. *Journal of applied ichthyology*, **26**, 102–108.

959    Tavaré S (1986) Some probabilistic and statistical problems in the analysis of DNA sequences. In:

960    *Lectures in mathematics in the life sciences* (ed Miura RM), pp. 57–86. American Mathematical

961    Society, Providence, RI.

962    Tsipas G, Tsiamis G, Vidalis K, Bourtzis K (2009) Genetic differentiation among Greek lake

963    populations of *Carassius gibelio* and *Cyprinus carpio carpio*. *Genetica* **136**, 491–500.

964    Urho L, Lehtonen H (2008) Fish species in Finland. Finnish Game and Fisheries Research Institute,

965    Helsinki.

966    Van Oppen MJH, Rico C, Turner GF, Hewitt GM (2000) Extensive Homoplasy, Nonstepwise

967    Mutations, and Shared Ancestral Polymorphism at a Complex Microsatellite Locus in Lake Malawi

968    Cichlids. *Molecular Biology and Evolution*, **17**, 489–498.

969    Weir B, Cockerham C (1984) Estimating F-statistics for the analysis of population structure.

970    *Evolution*, **38***, 1358-1370.*

971    Wheeler A (2000) Status of the crucian carp, Carassius carassius (L.), in the UK. *Fisheries*

972    *management and ecology*, **7**, 315–322.

973    Wolfram G, Mikschi E (2007) Rote Liste der Fische (Pisces) Österreichs. In: *Rote Liste gefährdeter*

974    *Tiere Österreichs, Teil 2*. Grüne Reihe des  Lebensministeriums Band 14/2. (ed Zulka K), pp. 61–

975    198. Böhlau-Verlag, Wien.

976    Wouters J, Janson S, Lusková V, Olsén KH (2012) Molecular identification of hybrids of the

977    invasive gibel carp *Carassius auratus gibelio* and crucian carp *Carassius carassius* in Swedish

978    waters. *Journal of fish biology*, **80**, 2595–2604.

979    Xu P, Zhang X, Wang X *et al.* (2014) Genome sequence and genetic diversity of the common carp,

980    *Cyprinus carpio*. *Nature genetics*, **46**, 1212-1219.

981

982    **Author contributions**

983    DLJ collected, samples, performed lab work, analysed data and wrote the manuscript. BH was

984    involved with conception of the project, advised on all steps of analyses and commented on the

985    manuscript. GHC was involved in the conception of the project, contributed samples and

986    commented on the manuscript. LJLH advised on the analysis and commented on the manuscript.

987    CDS and KHO provided samples and commented on manuscript.

988

989    **Data accessibility**

990    Genbank accession numbers for mtDNA sequences are provided in Table 1 of this manuscript. The

991    microsatellite data files for all DAPC analyses, mantel test matrices, mtDNA raw sequences,

992   sequence alignments, model testing outputs, tree files, RADseq loci catalog and VCF files used for

993   analyses and the DIYABC project files (containing all inputs, scenarios and parameter priors for

994   each analysis stage) have now been uploaded to Dryad (http://dx.doi.org/10.5061/dryad.t2j45). All

995   scripts used for DAPC, mantel tests and comparisons of RADseq and microsatellite datasets can be

996   found on GitHub (https://github.com/DanJeffries/Jeffries-et-al-2016-crucian-phylogeography). All

997   demultiplexed RADseq reads have been uploaded to the short read archive (Project accession:

998   SRP063043).

999

1000

1001  Figure 1. Population structure of *C. carassius* in Europe. a) Sampling locations (sites sampled with
1002  nuclear and mtDNA markers = red dots, mtDNA only = blue dots) and population cluster
1003  memberships from DAPC analysis. Pie chart size corresponds to microsatellite allelic richness. Pie
1004  chart colours for Danubian populations and RUS1 correspond to clusters in the broad scale DAPC
1005  analysis b)  and for all northern European populations colours correspond to clusters in the northern
1006  European DAPC analysis (mtDNA lineage 1 only) c). The Danube river catchment is shaded dark
1007  grey.

1008

1009  Figure 2. Maximum credibility tree calculated in BEAST for 100 *C. carassius cyt*b sequences. For
1010  the three maximally supported nodes, age is given above and the posterior probability distribution is
1011  given below, with 95% CI's represented by blue bars.

1012

1013  Figure 3. Comparison of DAPC results using a) RADseq dataset, b) M2 dataset and c) M3 dataset.
1014  Colours correspond between DAPC scatter plots and maps within but not between panels.

1015

1016  Figure 4. The postglacial recolonisation of *C. carassius* in Europe. Arrows represent the
1017  relationships between population pools used in DIYABC (grey circles) as inferred from Stage 1,
1018  scenario 9 (arrows outlined in black) and Stage 3, scenario 14d (arrows with no outline) analyses on
1019  RADseq data. Bottlenecks are represented by white-striped sections of arrows. Posterior time
1020  estimates in years for each demographic event are given in black, and estimates of Ne are given in
1021  blue. Blue diamonds represent ancestral populations inferred by DIYABC and the labels (a-f)
1022  correspond to their mention in the text. Hypothetical expansion fronts are represented by dashed
1023  contour lines and the Danube river catchment is shaded red. Hypothetical glacial refugia are
1024  represented by dashed blue circles (I - III). The blue dashed box (?) represents our inference that *C.*
1025  *carassius* expanded into central and perhaps northern Europe during the Riss-Würm interglacial,
1026  however we cannot estimate this range.

1027  SI Figure 1. DIYABC scenarios used in broad-scale analysis (Stage 1). See text for population
1028  poolings. See Table 3 for population poolings and prior parameter values.

1029

1030  SI Figure 2. All scenarios tested in stage 2 a) and stage 3 b) of DIYABC analysis. See Table 3 for
1031  population poolings and prior parameter values.

1032

1033  SI Figure 3. Filtering out merged ohnologs. a) Distribution of SNP locus coverage prior to
1034  removing loci that had observed heterozygosity higher than 0.5 in one or more population. b)
1035  Distribution of locus coverage after filtering, showing a loss of many high coverage loci and a
1036  reduction in mean SNP coverage. Note the loss of loci with high coverage.

1037

1038  SI Figure 4. Linear regressions for all samples a) *A*r against latitude; b) *A*r against longitude and for
1039  only samples in mtDNA lineage 1 c) *A*r against latitude; d) *A*r against longitude.

1040

1041  SI Figure 5. DAPC analysis of a) full microsatellite dataset (Excluding NOR2); for results used in
1042  Fig. 1) and b) Full RADseq dataset.

1043

1044  SI Figure 6. Isolation by distance a) in M1 dataset for mtDNA lineage 1 only (excluding NOR2), b)
1045  full RADseq dataset, c) M2 dataset and d) M3 dataset.

1046

1047  SI Figure 7. DAPC scatter plot for the 1000 SNP RADseq dataset used in the DIYABC analysis,
1048  showing the same population structure as inferred from the full RADseq dataset.

1049

1050  SI Figure 8. Broad scale DIYABC analysis (Stage 1) results. a) Direct approach (left) and Logistic
1051  regression (right) showing support for scenario 9. b) Model checking for scenario 9, showing that
1052  the observed data fall well within the cloud of datasets simulated from the posterior parameter
1053  distribution. c) Scenario 9 schematic.

1054

1055  SI Figure 9. Fine scale DIYABC analysis in northern Europe. a) Stage 2 - major topological
1056  variants of scenarios. Direct approach (top left) and Logistic regression (top right) showing support
1057  for scenario 14 and 13 respectively. Model checking (Middle) for scenario 14 (bottom), showing
1058  that the observed data fall well within the cloud of datasets simulated from the posterior parameter
1059  distribution. Note the model checking placed the observed data outside of the cloud of posterior
1060  datasets for scenario 13. b) Stage 3 - Minor scenario variants of scenario 14 from stage 2. Direct

1061    approach (top left), logistic regression (top right) and model checking (middle) all support scenario

1062    14d (bottom).

1063

1064    SI Figure 10. Comparison of spatial patterns of uniformity in geographic sampling regimes of the

1065    full M1 dataset locations (a, c) and the sampling location subset used in M2, M3, and RAD datasets

1066    (b,d). Estimates of G and L from true sampling locations are plotted using the black solid lines.

1067    Estimates of G and L from simulated locations based on random Poisson distribution is represented

1068    by the red dashed line. Grey shaded areas are the 95% confidence intervals around the random

1069    estimates. Both the G and L function estimates show that there is more clustering of sampling

1070    locations in the M1 dataset than in the M2, M3 and RAD subsets.

1071

1072    SI Figure 11. Change in a) number of RAD tags and b) average tag coverage for three individuals

1073    used in the preliminary Stacks tag mismatch parameter (M) tests.

1074

1075    SI Figure 12. Results of parameter tests for the Stacks module Populations. a) Number of SNP loci

1076    in final dataset for incrementing values of parameters –p, -r and –m; b) average coverage per SNP

1077    and per sample for the same parameter values; c) the number of loci which drop out in each

1078    population for each test value of the –p parameter

1079

1080

1081   Table 1. Location, number, genetic marker sampled, and accession numbers of samples and sequences used
1082   in the present study for microsatellite and mitochondrial DNA analyses. mtDNA sequence accession
1083   numbers can be found in SI table 2.

| Code | Accession | Location | Country | Drainage | Coordinates lat | long | Microsatellites | mtDNA | RADseq |
|---|---|---|---|---|---|---|---|---|---|
| GBR1 | | London | U.K. | U.K | 51.5 | 0.13 | 9 | | |
| GBR2 | | Reading | U.K. | U.K | 51.45 | -0.97 | 4 | | |
| GBR3 | | Norfolk | U.K. | U.K | 52.86 | 1.16 | 7 | | |
| GBR4 | | Norfolk | U.K. | U.K | 52.77 | 0.75 | 27 | | 9 |
| GBR5 | | Norfolk | U.K. | U.K | 52.77 | 0.76 | 14 | | |
| GBR6 | | Norfolk | U.K. | U.K | 52.54 | 0.93 | 29 | 3 | |
| GBR7 | | Norfolk | U.K. | U.K | 52.9 | 1.15 | 24 | 1 | 10 |
| GBR8 | | Hertfordshire | U.K. | U.K | 52.89 | 1.1 | 37 | 3 | 9 |
| GBR9 | | Norfolk | U.K. | U.K | 52.8 | 1.1 | 27 | | |
| GBR10 | | Norfolk | U.K. | U.K | 52.89 | 1.1 | 14 | | |
| GBR11 | | Norfolk | U.K. | U.K | 52.92 | 1.16 | 20 | | |
| BEL1 | | Bokrijk | Belgium | Scheldt River | 50.95 | 5.41 | 13 | 1 | |
| BEL2 | | Meer van Weerde | Belgium | Scheldt River | 50.97 | 4.48 | 12 | | |
| BEL3 | | Meer van Weerde | Belgium | Scheldt River | 50.97 | 4.48 | 8 | | |
| GER1* | | Kruegersee | Germany | Elbe River | 52.03 | 11.97 | | 3 | |
| GER2 | | Münster | Germany | Rhine River | 51.89 | 7.56 | 21 | 3 | |
| GER3 | | Bergheim | Germany | Danube River | 48.73 | 11.03 | 9 | 3 | |
| GER4 | | Bergheim | Germany | Danube River | 48.73 | 11.03 | 8 | 3 | |
| CZE1 | | Lužnice | Czech Republic | Danube River | 48.88 | 14.89 | 9 | 3 | |
| POL1 | | Sarnowo | Poland | Vistula River | 52.93 | 19.36 | 33 | | |
| POL2 | | Kikót-Wies | Poland | Vistula River | 52.9 | 19.12 | 34 | | |
| POL3 | | Tupadly | Poland | Vistula River | 52.74 | 19.3 | 17 | 3 | 10 |
| POL4 | | Orzysz | Poland | Vistula River | 53.83 | 22.02 | 13 | 3 | 10 |
| EST1 | | Tartu | Estonia | Baltic Sea | 58.39 | 26.72 | 5 | 3 | |
| EST2 | | Vehendi | Estonia | Baltic Sea | 58.39 | 26.72 | 5 | | |
| RUS4* | | Small lake, Velikaya river | Russia | Baltic Sea | 55.9 | 30.25 | 29 | 3 | |
| FIN1 | | Joensuu | Finland | Baltic Sea | 62.68 | 29.68 | 32 | 3 | |
| FIN2 | | Helsinki | Finland | Baltic Sea | 60.36 | 25.33 | 32 | | |
| FIN3 | | Jyväskylä | Finland | Baltic Sea | 62.26 | 25.76 | 37 | 3 | 10 |
| FIN4 | | Oulu | Finland | Baltic Sea | 65.01 | 25.47 | 7 | 3 | 8 |
| FIN5 | | Salo | Finalnd | Baltic Sea | 60.37 | 23.1 | 10 | 3 | |
| FIN6 | | Åland Island | Sweden | Baltic Sea | 60.36 | 19.85 | 8 | 3 | |
| SWE1 | | Gränbrydammen | Sweden | Baltic Sea | 59.87 | 17.67 | 25 | | |
| SWE2 | | Stordammen | Sweden | Baltic Sea | 59.8 | 17.71 | 21 | 3 | 10 |
| SWE3 | | Östhammar | Sweden | Baltic Sea | 60.26 | 18.38 | 27 | 3 | |
| SWE4 | | Umeå | Sweden | Baltic Sea | 63.71 | 20.41 | 9 | 3 | |
| SWE5 | | Kvicksund | Sweden | Baltic Sea | 59.45 | 16.32 | 9 | | |
| SWE7 | | Grillby | Sweden | Baltic Sea | 59.64 | 17.37 | 10 | | |
| SWE8 | | Skabersjo | Sweden | Baltic Sea | 55.55 | 13.15 | 19 | 3 | 10 |
| SWE9 | | Märsta | Sweden | Baltic Sea | 59.6 | 17.8 | 31 | 3 | |
| SWE10 | | Norrköping | Sweden | Baltic Sea | 58.56 | 16.27 | 29 | | 9 |
| SWE11 | | Gotland Island | Sweden | Baltic Sea | 57.85 | 18.79 | 11 | 3 | |
| NOR1 | | Oslo | Norway | North Sea | 60.05 | 9.94 | | 2 | |
| NOR2 | | Lake Prestvattnet, Tromsø | Norway | North Sea | 69.65 | 18.95 | 16 | | 9 |
| BLS | | | Belarus | Dnieper | 52.47 | 30.52 | 7 | 1 | |
| RUS1 | | Proran Lake | Russia | Don River | 47.46 | 40.47 | 10 | 3 | 9 |
| DEN1 | | Copenhagan | Denmark | Baltic Sea | 60.21 | 17.79 | 12 | | 10 |
| DEN2 | | Pederstrup | Denmark | Baltic Sea | 55.77 | 12.55 | 14 | | 8 |
| DEN3 | | Gammel Holte | Denmark | Baltic Sea | 56 | 12.5 | 14 | | |
| DEN4 | | Bornholm Island | Denmark | Baltic Sea | 55.17 | 14.86 | | | 5 |
| SWE12 | | Osterbybruk Mansion | Sweden | Baltic Sea | 55.73 | 12.34 | 14 | | 9 |
| SWE14 | | Wenngarn Castle | Sweden | Baltic Sea | 59.66 | 18.95 | 16 | | 9 |
| RUS2* | | Karma | Russia | Volga River | 52.9 | 58.4 | | 2 | |
| RUS3* | | Saygach'yedake | Russia | Volga River | 47.5 | 48.5 | | 4 | |
| TNO | | | Netherlands | North Sea | - | - | | 1 | |
| HUN1 | | Gödöllő | Hungary | Danube River | 47.61 | 19.36 | | 2 | 6 |
| HUN2 | | Vörösmocsár | Hungary | Danube River | 46.49 | 19.17 | | | |
| | | | | | | | **848** | **83** | **160** |
| | | | | | | | | Total number of fish = | 867 |

**Genbank mtDNA Sequences**

| Code | Accession | Reference | Country | Drainage |
|---|---|---|---|---|
| GER6 | DQ399917 | Kalous et al. (2007) | Germany | Baltic sea |
| GER6 | DQ399918 | Kalous et al. (2007) | Germany | Baltic sea |
| GER6 | DQ399919 | Kalous et al. (2007) | Germany | Baltic sea |
| GER7 | JN412540 | Rylková et al. (2013) | Germany | Hunte River |
| GER7 | JN412541 | Rylková et al. (2013) | Germany | Hunte River |
| GER7 | JN412542 | Rylková et al. (2013) | Germany | Hunte River |
| GER7 | JN412543 | Rylková et al. (2013) | Germany | Hunte River |
| GER8* | JN412537 | Rylková et al. (2013) | Germany | Lahn River |
| GER8* | JN412538 | Rylková et al. (2013) | Germany | Lahn River |
| CZE2 | GU991399 | Rylková et al. (2013) | Czech Republic | Elbe drainage |
| Milevsko | DQ399938 | Kalous et al. (2012) | Czech Republic | Elbe drainage |
| AUS1 | JN412534 | Rylková et al. (2013) | Austria | Danube river |
| AUS1 | JN412533 | Rylková et al. (2013) | Austria | Danube river |
| AUS2 | JN412535 | Rylková et al. (2013) | Austria | Danube river |
| AUS3 | JN412536 | Rylková et al. (2013) | Austria | Danube river |
| GBR12 | JN412539 | Rylková et al. (2013) | U.K. | U.K |

| GBR12 | GU991400 | Kalous et al. (2012) | U.K. | U.K |
|-------|----------|----------------------|------|-----|
| SWE15 | JN412545 | Rylková et al. (2013) | Sweden | Baltic sea |
| SWE16 | JN412544 | Rylková et al. (2013) | Sweden | Baltic sea |
| Ccarp1 | AB158807 | Mabuchi et al (2005) | Japan | - |
| Ccarp2 | DQ868875 | Tsipas et al. (2009) | Greece | - |
| Ccarp3 | KF574490 | Unpublished | India | - |

1084 † Also present

1085 * Location on Map (Fig. 1.a) is approximate

1086

1087 Table 2. Population pools, parameter priors used and median posterior parameter values inferred in the three
1088 stages of DIYABC analysis.

| Analysis stage | Population Pools | Scenarios tested | Parameter priors | Most likely Scenario | Median of posterior distributions of most likely scenario |
|---|---|---|---|---|---|
| 1 | **Pool 1** – GBR4, GBR7, GBR8, DEN1, DEN2, DEN3, FIN3, FIN4, POL3, POL4, SWE2, SWE8, SWE9, SWE10, SWE12, SWE14, NOR2 <br> **Pool 2** – DEN1, DEN2, DEN3 <br> **Pool 3** – FIN3, FIN4 | 1 - 11 | N1 = 10E+03 - 500E+03 <br> Nb1 = 10 - 100E+03 <br> N2 = 100 - 100E+03 <br> N3 = 100 - 200E+03 <br> t1 = 1E+03 - 1E+06 gens <br> t2 = 1E+03 - 3E+06 gens <br> ra = 0.001-0.999 <br> rb = 0.001-0.999 <br> rc = 0.001-0.999 <br> db = 10- 10E+03 gens | 9 | N1 =34700 <br> Nb1 =23700 <br> N2 =74900 <br> N3 =140000 <br> t1 =135000 <br> db =4460 <br><br> t2 =1090000 |
| 2 | **Pool 1** – GBR4, GBR7, GBR8 <br> **Pool 2** – DEN1, DEN2, DEN3 <br> **Pool 3** – FIN3, FIN4 | 12 - 16 | N1 = 10-4E+03 <br> N2 = 10 - 10E+03 <br> N3 = 10 - 20E+03 <br> N4 = 10 - 50E+03 <br> N5 = 10 - 20E+03 <br> N6 =10 - 400 <br> t1 = 100- 10E+03 gens <br> t1a = 100- 10E+03 gens <br> t2 =100- 10E+03 <br> t2a =100- 5E+03 gens <br> t2b = 500-20E+03 gens <br> t2c = 100 - 10E+03 gens <br> t2d = 100 - 10E+03 gens <br> t3 = 500 - 20E+03 gens <br> t3c =100 - 10E+03 gens <br> t3d =100 - 10E+03 gens <br> t4 =500 - 20E+03 gens <br> ra = 0.001-0.999 <br> rb = 0.001-0.999 | 14 | N1 =3670 <br> N2 =7520 <br> N3 =17400 <br> N4 =19400 <br> N5 =11800 <br> N6 =210 <br> t1 =6790 <br> t1a =2510 <br><br><br><br><br> t2d =6780 <br><br><br> t3d =8910 <br> t4 =12000 <br><br> rb =0.668 |
| 3 | **Pool 4** – POL3, POL4 <br> **Pool 5** – SWE2, SWE8, SWE9, SWE10, SWE12, SWE14 <br> **Pool 6** – NOR2 | 14a - 14f | N1 = 10-4E+03 <br> Nb1 = 10-10E+03 <br> N2 = 10 - 10E+03 <br> N3 = 10 - 20E+03 <br> Nb3 = 10-10E+03 <br> N4 = 10 - 50E+03 <br> N5 = 10 - 20E+03 <br> N6 =10 - 400 <br> Nb6 =10-10E+03 <br> t1 = 100- 10E+03 gens <br> t1a = 100- 10E+03 gens <br> t2d = 100 - 10E+03 gens <br> t3d = 100 - 10E+03 gens <br> t4 = 500 - 20E+03 gens <br> rb = 0.001-0.999 <br> da = 10 - 10E+03 gens <br> db = 10 - 10E+03 gens <br> dc = 10 - 10E+03 gens <br> dd = 10 - 10E+03 gens <br> de = 10 - 10E+03 gens | 14d | N1 =2390 <br> Nb1 =935 <br> N2 =8140 <br> N3 =9360 <br><br> N4 =17000 <br> N5 =11000 <br> N6 =138 <br><br> t1 =3750 <br> t1a =2460 <br> t2d =5900 <br> t3d =7970 <br> t4 =16800 <br> rb =0.619 <br><br><br> dc =9070 |

1089

1090

1091 Table 3. Summary statistics for M1, M2, M3 and RADseq datasets. RAD contains all RADseq data, M1
1092 contains all microsatellite data, M2 contains only microsatellite for the individuals used in the RADseq, and
1093 M3 contains all microsatellite data for all individuals that were available in populations that were used in
1094 RADseq.

| Dataset | Description | N samples | Mean N samples/pop | N. loci | Mean N.alleles/pop | Mean N.alleles/locus |
|---------|-------------|-----------|--------------------|---------|--------------------|----------------------|
| RAD | RADseq data only | 149 | 8.95 ± 1.4 | 13189 | 6723 | 2 |
| M1 | Full Microsatellite dataset | 848 | 17.2 ± 9.5 | 13 | 27 ± 8.8 | 7.6 |
| M2 | Microsatellites for RADseq samples only | 146 | 9.13 ± 0.8 | 13 | 24.4 ± 7.3 | 7.84 ± 5.1 |
| M3 | Microsatellites for all samples in populations used in RADseq | 313 | 19.6 ± 9.0 | 13 | 27.4 ± 8.1 | 11.23 ± 7.6 |

1095

1096

1097 Table 4. Pearson's product-moment correlation coefficients and paired t-tests comparing heterozygosities
1098 and $F_{ST}$s between M2, M3 and RADseq datasets. *** $P = {<}0.001$, ** $P = {<}\,0.005$, * $P = {<}\,0.05$.

| Heterozygosities (df = 18) | Pearsons correlation coefficient (t) | | |
|----------------------------|--------|----------|----------|
| | **M2** | 11.13*** | 3.85** |
| Paired T-tests | -2.4* | **M3** | 3.86** |
| | -9.71*** | -9.29*** | **RAD** |

| $F_{ST}$ (df = 105) | Pearsons correlation coefficient (t) | | |
|----------------------|--------|----------|----------|
| | **M2** | 46.26*** | 10.09*** |
| Paired T-tests | -6.21*** | **M3** | 9.05*** |
| | 13.74*** | 15.12*** | **RAD** |

1099

A) **RADseq -** 13,189 SNPs, n = 149.



DA eigenvalues

B) **M2 -** 13 Microsatellites, n = 146.



DA eigenvalues

B) **M3 -** 13 Microsatellites, n = 313.



DA eigenvalues

1   Comparing RADseq and microsatellites to infer
2   complex phylogeographic patterns, a real data
3   informed perspective in the Crucian carp, *Carassius*
4   *carassius,* L.

5

6   **Authors:** [1]Daniel L Jeffries, [2]Gordon H Copp, [1]Lori Lawson Handley, [3]K. Håkan Olsén, [4]Carl D
7   Sayer, [1]Bernd Hänfling

8

# Supporting Information

9

10   *Detecting hybrids*

11   *Methods*

12   In total we acquired tissue samples of 1354 fish from 72 populations. All samples were first

13   genotyped using multiplex 1 (SI table 1) which contained the 6 species diagnostic microsatellite

14   loci. These data were then analysed using the NewHybrids v. 1.1 (Anderson & Thompson 2002)

15   software package in order to determine whether each fish was *C. carassius, C. auratus, C. gibelio*

16   or a hybrid between any of these species.

17

18   NewHybrids uses allele frequencies to give a likelihood probability that an individual belongs to

19   one species or another, or if the individual belonged to one of several hybrid classes (F1, F2 or

20   backcross). Data from 20 *C. carassius* samples, which were confidently identified as pure from both

21   morphology and genotypes, and were not sympatric with non-native species, were included in each

22   analysis as baseline data. Priors were then added to the analyses specifying that these individuals

23   were indeed pure in order to give the software more power with which to assess allele frequencies

24   associated with *C. carassius*. To be sure to account for allele frequency differences between

25   different geographic regions, only pure individuals from regions neighbouring the hybrid population

26    were used. Individuals which had more than a 25% chance of being an F1 hybrid, F2 hybrid, or a

27    backcross were removed from population structure analyses and were not genotyped at the

28    additional 7 microsatellite loci (Multiplexes 2.1 and 2.2, SI table 1).

29

30    *Results*

31    Of the 1354 fish which were genotyped with microsatellites, 942 individuals across 55 populations

32    (86.7%) were identified as pure *C. carassius* using the first set of 6 species diagnostic loci in

33    NewHybrids analyses. 19 (1.8%) from 2 different populations were identified as *C. auratus*, 15 fish

34    (1.4%) from 4 populations were identified as. *C. gibelio* and 10 fish (0.93%) from two populations

35    were identified as *C. carpio*. NewHybrids identified 60 (5.5%) *C. carassius* x *C. auratus* hybrids,

36    25 (2.2%) *C. carassius* x *C. gibelio* hybrids, and 16 (1.5%) *C. carassius* x *C. carpio* hybrids. Of the

37    942 fish identified as pure *C. carassius*, 867 in existed in locations  (49 populations) where hybrids

38    or non-native species were not detected by microsatellite genotyping. To safeguard against cryptic

39    introgression which may produce erroneous results only these 867 pure *C. carassius* were used for

40    the main phylogeographic analyses and tests using either microsatellites, mtDNA or RADseq.

41

42    *RADseq data filtering and Stacks analysis parameter testing*

43    RADseq analyses were performed using only the first-end reads from the paired-end sequencing, as

44    coverage across the length of the second-end contigs was not consistent enough to call SNPs in all

45    individuals. For these first-end reads, raw data was first quality checked using FastQC (Andrews

46    2010), which assesses the per-base sequence quality and content of reads, and provides

47    comprehensive graphical outputs with which to assess the overall quality of raw sequencing data.

48    These analyses did not identify any individuals that had low overall sequence quality, therefore all

49    samples were retained for further analyses.

50

51     Preliminary analyses were also carried out using PyRAD (Eaton 2014), which allows for the

52     incorporation of allelic variants resulting from insertions and deletions. However, no significant

53     difference in the number of usable loci was shown. As Stacks provides more downstream

54     populations genetics facilities, this program was used for the final analyses.

55

56     Raw RADseq reads were first, demultiplexed using the "process_radtags" module distributed with

57     Stacks and our inline barcodes. Second, reads were filtered for any sequences containing Illumina

58     adapters or primers and trimmed to a length of 92 bp. Third, PCR duplicates introduced during

59     library preparation were removed using the "clone_filter" program (also distributed with Stacks).

60     Finally, preliminary tests of parameter values for each module of the de novo stacks pipeline were

61     performed in order to identify "optimal" parameter values (i.e. where loci number and read depth

62     were stable) for use in the final Stacks analysis. These tests were carried out for 5 sets of 3

63     randomly chosen individuals from the RADseq dataset and, for each test, all non-test parameters

64     were kept as default.  In the ustacks module, which groups identical reads into stacks and then

65     stacks into loci, Parameters M and m were tested (See Catchen et al. 2013 for detailed description

66     of parameters). M values were increased in increments of 2 from 0 to 10. The efficiency of ustacks

67     in finding real loci was then examined with simple counts of the number of constructed loci at each

68     M parameter value and the read coverage of these loci. The expectation was that, at low parameter

69     values, divergent alleles (percentage divergence > M) at a locus will not merge (under-merging),

70     thus increasing the number of loci overall and decreasing the average coverage. In contrast high

71     parameter values could cause over-merging of paralogous loci and have the opposite effects on the

72     number of loci and coverage (Catchen et al. 2013). SI Fig. 11 shows the outputs for a single subset

73     of *C. carassius* samples, which was typical of all 5 subsets tried. In ustacks, an 'm' parameter value

74     of zero (minimum of 0 reads required to form a stack) resulted in a very large number of tags

75     (49000-54000) as expected. Likely due to many single reads containing sequencing error being

76     called as loci. The number of loci decreased by approximately 3000 – 4000 tags in the samples

77    tested at a required read depth of 2 (approx. 50,000), after which further increases in 'm' resulted in

78    small decreases in the number of tags. This likely reflects merging of paralogous loci, or low

79    coverage loci. Mean coverage across all loci within an individual of course reflected the 'm'

80    parameter increase, jumping initially from approx. 16 reads per locus with zero read depth required,

81    to 20-35 at a minimum required depth of two reads. On the basis of these results we chose an m = 8,

82    to ensure high power for SNP calling.

83

84    Incrementing over values of 'M' again met our expectations, with the number of loci dropping

85    significantly as the 'M' parameter was increased from zero to 2 mismatches allowed, and then

86    dropping more slowly with higher mismatch allowance. These further drops may again be allowing

87    for paralog merging between loci. The mean coverage of loci behaved as expected, with higher

88    mismatch allowance, more divergent reads can be added to existing stacks, inflating coverage for

89    those loci. On the basis of these results M=2 was chosen for final analyses.

90    Parameter tests were also performed for the cstacks parameter N, which is responsible for setting

91    the maximum mismatch threshold allowed between homologous loci among individuals in the locus

92    catalog. First, ustacks was run using chosen "optimal" parameters to obtain the inputs necessary for

93    cstacks. Cstacks was then run separately on each of the 5 sample subsets with values of N between

94    0 – 10, with increments of 2.

95

96    Finally, we tested three core parameters in the Populations module of Stacks, -m which is analogous

97    to the parameter of the same name in the ustacks module, -r, which specifies the number of

98    individuals within a give population that a locus must be present in, and –p which specifies the

99    number of populations that a locus must be present in (above the –r threshold) for it to be retained

100   in the final dataset (SI Fig. 12). –p was tested for values of between 13 – 19 populations, -r was

101   tested for values between 0.5 – 1.0 and –m was tested for values between 1-8 however, a the dataset

102   had previously been filtered at previous stages for loci present with a depth of 8 reads or higher, the

103    tests of –m in the populations stage were redundant.

104

105    *Final running parameters used*

106    For all parameter tests, the optimal values were taken to be those where the rate of change in either

107    RAD tag number, or coverage began to decrease. In ustacks, a maximum of two mismatches were

108    allowed between alleles at a given locus (M=2) and at least eight identical reads per stack (m=8)

109    were required. Default values were used for all other parameters. ustacks also called SNPs within

110    individuals at each locus. The cstacks module was then used to merge loci across individuals into a

111    catalog, where N=2 mismatches were allowed between individuals at a given locus. Individuals

112    were then searched against this catalog using Sstacks to determine their genotype at each catalog

113    locus. For the Populations module, optimal values were chosen so that loci that were shared

114    between at least 70% of individuals in each population (-r = 0.7), allowing loci to drop out in one or

115    two individuals in a population for reasons of low DNA sample quality or low coverage. Loci must

116    have also been present in 17 of the 19 populations (-p = 17), and have read depth of at least 8 (-m 8)

117    in each individual.

118

119    *DAPC & Running parameters*

120    *Methods*

121    Population structure was examined using Discriminant Analyses of Principal Components (DAPC,

122    (Jombart *et al.* 2010)) in adegenet. Similar to the more commonly used program, STRUCTURE

123    (Pritchard et al. 2000), DAPC is an individual-based approach that uses Principal Components

124    Analysis (PCA) to transform population genetic data and Discriminant Analysis (DA) to identify

125    clusters. The number of clusters is assessed using the K-means method, which is also used in

126    STRUCTURE (Pritchard *et al.* 2000). Unlike STRUCTURE, DAPC does not assume underlying

127    population genetics models such as Hardy-Weinberg Equilibrium (Jombart *et al.* 2010) and is

128    therefore more suitable for analysing C. carassius since populations are often bottlenecked

129    (Hänfling *et al.* 2005). An additional benefit of DAPC is that it maximizes between-group variation,

130    while minimizing variation within groups, allowing for optimal discrimination of between-

131    population structure (Jombart *et al.* 2010).

132

133    *Results*

134    For the full microsatellite dataset (M1), BIC scores indicated that between 11 and 19 genetic

135    clusters (**Error! Reference source not found.**) would be an appropriate model of the variation in the

136    data. We therefore chose 11 clusters to use in the discriminant analysis, retaining 8 principal

137    components as recommended by the spline interpolation a-scores (**Error! Reference source not**

138    **found.**c) and we kept 2 linear discriminants for plotting (**Error! Reference source not found.**b).

139

140    Three major lineages were found, one located in the Danube, one in the Don, and one spread across

141    northern Europe. However the large amount of divergence between them masked the population

142    structure present in northern Europe. We therefore subsetted the data, separating NEU populations

143    from RUS1, GER3, GER4, CZE1 (and SWE9, which was an outlier within NEU, **Error! Reference**

144    **source not found.**b) and reanalysed them with DAPC in order to better infer fine population structure

145    between them.

146

147    For the RADseq dataset, BIC scores suggested between 9 and 14 genetic clusters, similar to the

148    range inferred in the microsatellite data, we therefore chose 9 clusters to take forward in the

149    analysis. As recommended by spline interpolation, we retained 7 principal components and we kept

150    2 of the linear discriminants from the subsequent discriminant analysis

151

152    *Assessment of spatial uniformity of sampling locations*

153    *Methods*

154    In order to assess the geographic uniformity of the sampling regimes in each data subset, we used

155    two measures of spatial patterns. The nearest neighbour distance distribution function (G), measures

156    the distance of each sampling location to its nearest neighbour (Ripley 1991). The L-function is a

157    transformation (for ease of interpretation) of Ripley's K-function (Ripley 1991), which measures

158    the number of sampling locations within a given radius from each point. K has the advantage of

159    assessing the uniformity of the sampling regime over multiple scales, as opposed to only measuring

160    distances between closest neighbours as with G. In both cases, the estimates of G or K from our

161    sampling locations were compared against random Poisson distributions, which would represent

162    uniformly spaced sampling locations. 5% and 95% confidence thresholds for these Poisson

163    distributions were also calculated to allow us to determine whether our sampling regimes

164    significantly deviated from random (p <0.05). These calculations were performed using the Gest

165    and Lest functions (for G and L respectively) in the package "spatstats" in R (Baddeley & Turner

166    2005).

167

168    *Results*

169    Both methods used for the assessment of geographic uniformity of sampling locations shows that

170    the M1 dataset locations are more patchily distributed than those of the M2, M3 and RAD datasets

171    (**Error! Reference source not found.**).

172

173   *Additional discussion*

174   *Population structure in northwest Europe*

175   An intriguing result lies in the genetic similarity between populations in England with those in

176   Belgium and Germany. *C. carassius* has been designated as native to England, however this status

177   has been contentious in the past (Maitland 1972). Under the assumption that it is native, and

178   considering the observed diversity and divergence times between populations across mainland

179   Europe, we would expect to see stronger population structure between English and continental

180   Europe, which have been separated for approximately 7800 years (Coles 2000). Given the observed

181   diversity between populations across mainland Europe, which, according to DIYABC analysis, has

182   arisen relatively recently. Clearly further examination of this issue is warranted and molecular data

183   would be a value addition to the current evidence, which is predominantly anecdotal.

184

185   SI table 1. Microsatellite loci used, grouped by their combinations in multiplex reactions. Multiplex primer
186   mix ratios for PCR were chosen so as to give even peak strengths when analysing PCR products. Allele size
187   ranges are those present in C. carassius for all 43 putatively pure crucian populations.

| Locus | Multiplex # | Primer mix Ratios* | # Alleles | Allele size range | Ho | GenBank Accession no. | Reference |
|---|---|---|---|---|---|---|---|
| GF1 | 1 | 0.1 | 1 | 299 | 0 | U35614 | Zheng et al. 1995 |
| GF17 | 1 | 0.1 | 2 | 182-186 | 0.024 | U35616 | Zheng et al. 1995 |
| GF29 | 1 | 0.2 | 8 | 191-226 | 0.348 | U35618 | Zheng et al. 1995 |
| J7 | 1 | 0.07 | 10 | 202-228 | 0.109 | AY115095 | Yue & Orban 2002 |
| MFW2 | 1 | 0.1 | 1 | 161 | 0 | - | Croojimans et al. 1997 |
| Ca07 | 1 | 0.2 | 9 | 122-140 | 0.286 | D85428 | Yue & Orban 2004 |
| TE Buffer | 1 | 0.23 | | | | | |
| J69 | 2.1 | 0.4 | 14 | 213-241 | 0.404 | AY115106 | Yue & Orban 2002 |
| HJLY17 | 2.1 | 0.1 | 9 | 152-168 | 0.223 | DQ378986 | Zhi-Ying et al. 2006 |
| HJLY35 | 2.1 | 0.1 | 18 | 261-307 | 0.377 | DQ403242 | Zhi-Ying et al. 2006 |
| TE Buffer | 2.1 | 0.4 | | | | | |
| J20 | 2.2 | 0.2 | 9 | 171-218 | 0.149 | AY115099 | Yue & Orban 2002 |
| J58 | 2.2 | 0.1 | 14 | 119-147 | 0.398 | - | Yue & Orban 2002 |
| MFW7 | 2.2 | 0.35 | 25 | 160-206 | 0.464 | - | Croojimans et al. 1997 |
| MFW17 | 2.2 | 0.35 | 26 | 185-262 | 0.41 | - | Croojimans et al. 1997 |

188   * All primers used at 10mM per ul concentration, diluted in ddH20 from 100mM per ul stock

189     SI table 2. Genbank accession numbers for the mtDNA sequences used in this study.

190

| Sample code | Accession number |
| --- | --- |
| FIN5_01 | KT630314 |
| FIN5_02 | KT630315 |
| FIN5_03 | KT630316 |
| EST1_02 | KT630317 |
| GER1_01 | KT630318 |
| EST1_01 | KT630319 |
| GER1_03 | KT630320 |
| FIN6_01 | KT630321 |
| FIN6_02 | KT630322 |
| FIN6_03 | KT630323 |
| BEL1_03 | KT630324 |
| EST1_03 | KT630325 |
| GER2_02 | KT630326 |
| GER2_03 | KT630327 |
| GER4_02 | KT630328 |
| NOR1_01 | KT630329 |
| NOR1_02 | KT630330 |
| SWE11_01 | KT630331 |
| SWE11_02 | KT630332 |
| SWE11_03 | KT630333 |
| RUS2_02 | KT630334 |
| RUS4_01 | KT630335 |
| RUS4_03 | KT630336 |
| FIN1_01 | KT630337 |
| FIN1_02 | KT630338 |
| FIN1_03 | KT630339 |
| FIN4_01 | KT630340 |
| FIN4_02 | KT630341 |
| FIN4_03 | KT630342 |
| POL4_01 | KT630343 |
| POL4_02 | KT630344 |
| POL4_03 | KT630345 |
| RUS1_01 | KT630346 |
| RUS1_02 | KT630347 |
| RUS1_03 | KT630348 |
| SWE8_01 | KT630349 |
| SWE8_02 | KT630350 |
| SWE8_03 | KT630351 |
| POL3_01 | KT630352 |
| POL3_02 | KT630353 |
| POL3_03 | KT630354 |
| SWE4_01 | KT630355 |
| SWE4_02 | KT630356 |
| SWE4_03 | KT630357 |
| RUS3_01 | KT630358 |
| RUS3_03 | KT630359 |
| RUS3_04 | KT630360 |
| RUS2_01 | KT630361 |
| RUS4_02 | KT630362 |
| BLS_03 | KT630363 |
| RUS3_02 | KT630364 |
| SWE3_01 | KT630365 |
| SWE3_02 | KT630366 |
| SWE3_03 | KT630367 |
| SWE2_01 | KT630368 |
| SWE2_02 | KT630369 |
| SWE2_03 | KT630370 |
| SWE9_01 | KT630371 |
| SWE9_02 | KT630372 |
| SWE9_03 | KT630373 |
| GBR7_01 | KT630374 |
| GBR6_01 | KT630375 |
| GBR8_01 | KT630376 |
| GBR8_02 | KT630377 |
| GBR8_03 | KT630378 |
| GBR6_02 | KT630379 |
| GBR6_03 | KT630380 |
| CZE1_01 | KT630381 |
| CZE1_02 | KT630382 |
| CZE1_03 | KT630383 |
| GER4_01 | KT630384 |
| GER4_03 | KT630385 |
| GER1_02 | KT630386 |
| GER2_01 | KT630387 |
| FIN3_01 | KT630388 |
| FIN3_02 | KT630389 |
| FIN3_03 | KT630390 |
| HUN1_02 | KT630391 |
| GER3_01 | KT630392 |
| GER3_02 | KT630393 |
| GER3_03 | KT630394 |

191    SI table 3. Haplotype memberships for 101 Cytochrome B sequences used in Fig. 2.

| Lineage | Haplotype | N | Drainage (n populations) | Sample code |
|---|---|---|---|---|
| **1** | 1 | 3 | Baltic | FIN5 1-3 |
| | 2 | 1 | Baltic | EST1 2 |
| | 3 | 49 | Elbe(2), Baltic(9), Scheldt(1), Rhine(2), North sea(2), Vistula(6), Volga(4), Don(3), Danube(1), Hunte(4) | GER1 1,3, EST1 1, 3, SWE6 1 -3, BEL1 3 , GER2 2, 3, GER4 2, NOR 1, 2, SWE11 1-3, RUS2 2, RUS4 1, 3, FIN1 1-3, FIN4 1-3, POL4 1-3, RUS1 1-3, SWE8 1-3, POL5 1-3, SWE4 1-3, RUS3 1, 3, 4, CZE2 1, GER6 1 – 4, SWE14 1, SWE15 1 |
| | 4 | 1 | Volga | RUS2 1 |
| | 5 | 1 | Baltic | RUS4 2 |
| | 6 | 1 | Dnieper | BLS 3 |
| | 7 | 1 | Volga | RUS3 2 |
| | 8 | 3 | Baltic | SWE3 1-3 |
| | 9 | 2 | Baltic | SWE2 1 - 3 |
| | 10 | 3 | Baltic | SWE9 1-3 |
| | 11 | 13 | UK(4), Rhine(1), Baltic (2) | GBR7 1, GBR6 1-3, GBR8 1-3, NET 1, GER5 1-3, GBR12 1, 2 |
| | 12 | 3 | Baltic | FIN3 1-3 |
| **2** | 13 | 3 | Danube | GER4 1, 2, AUS3 1 |
| | 14 | 3 | Elbe(1), Rhine(1), Danube(1) | GER1 2, GER2 1, AUS2 1 |
| | 15 | 1 | Danube | CZE1 1 |
| | 16 | 1 | Danube | CZE1 2 |
| | 17 | 1 | Danube | CZE1 3 |
| | 18 | 2 | Danube | HUN 1, 2 |
| | 19 | 3 | Danube | GER3 1-3 |
| | 23 | 1 | Elbe | CZE2 2 |
| | 24 | 2 | Danube | AUS1 1, 2 |
| | 25 | 2 | Lahn | GER7 1, 2 |
| **Outgroup** | 20 | 1 | | Ccarp 1 |
| | 21 | 1 | | Ccarp 2 |
| | 22 | 1 | | Ccarp 3 |

192

193

194

195    SI table 4. Pairwise FST values calculated using the M1 dataset.

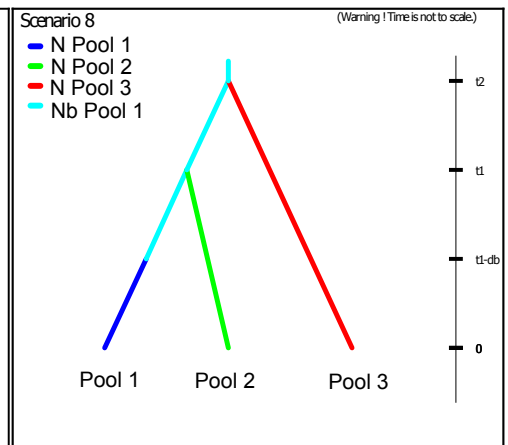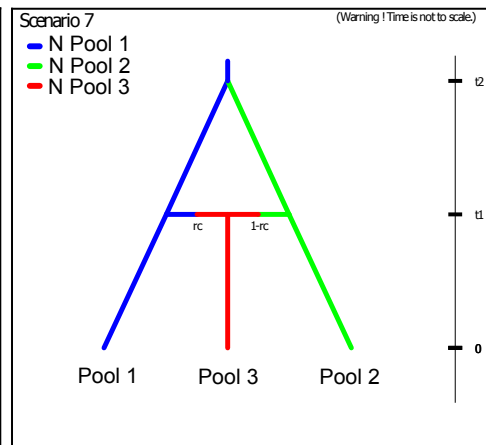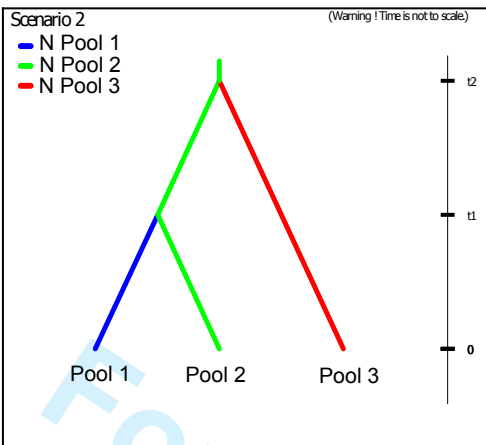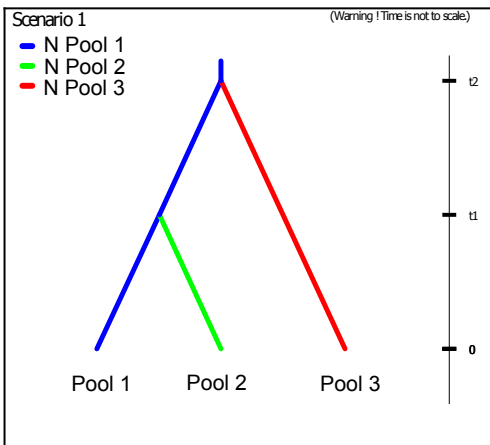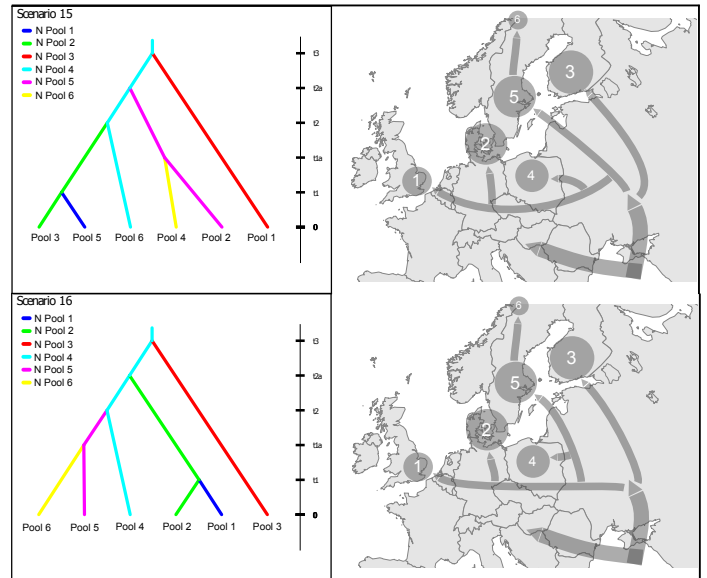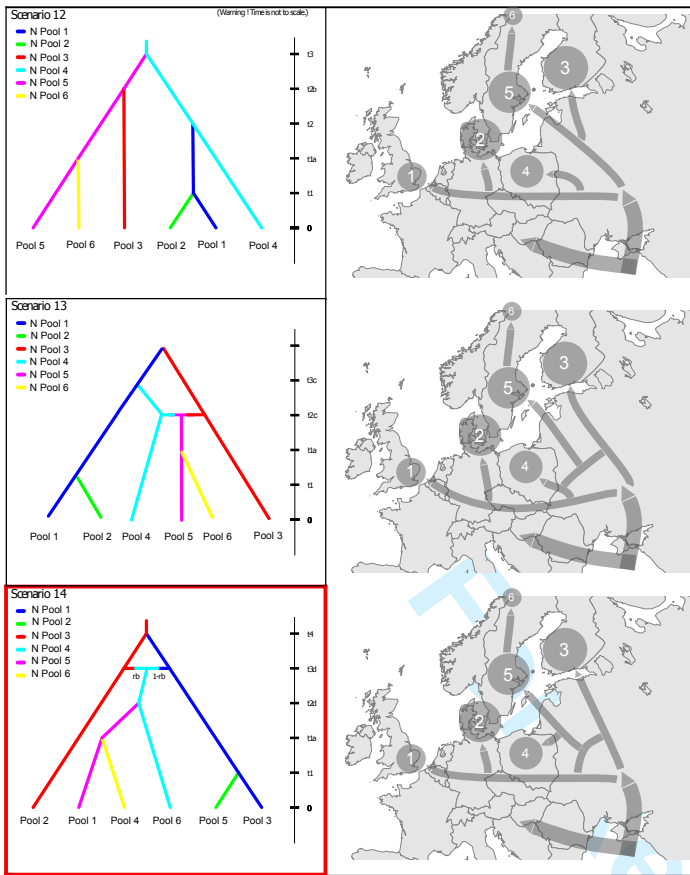| | GBR1 | GBR2 | GBR4 | BEL1 | BEL2 | BEL3 | FIN1 | RUS4* | FIN2 | CZE1 | GER2 | GER3 | GER4 | POL1 | POL2 | POL3 | POL4 | GBR7 | GBR3 | GBR8 | GBR9 | GBR11 | GBR5 | GBR6 | GBR10 | SWE4 | SWE3 | SWE5 | FIN6 | SWE7 | SWE2 | SWE1 | SWE9 | SWE10 | SWE11 | SWE8 | FIN5 | FIN3 | FIN4 | EST1 | EST2 | BLS | RUS1 | DEN1 | SWE12 | DEN2 | NOR2 | SWE14 | DEN3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GBR1 | | 0.307 | 0.531 | 0.312 | 0.198 | 0.346 | 0.785 | 0.472 | 0.407 | 0.604 | 0.256 | 0.613 | 0.628 | 0.226 | 0.291 | 0.342 | 0.368 | 0.436 | 0.364 | 0.378 | 0.518 | 0.317 | 0.517 | 0.302 | 0.376 | 0.479 | 0.444 | 0.419 | 0.458 | 0.542 | 0.591 | 0.404 | 0.839 | 0.548 | 0.793 | 0.39 | 0.428 | 0.72 | 0.596 | 0.628 | 0.526 | 0.491 | 0.623 | 0.319 | 0.626 | 0.261 | 0.768 | 0.457 | 0.233 |
| GBR2 | NS | | 0.67 | 0.316 | 0.247 | 0.366 | 0.783 | 0.482 | 0.446 | 0.588 | 0.332 | 0.6 | 0.618 | 0.266 | 0.309 | 0.357 | 0.378 | 0.611 | 0.535 | 0.562 | 0.716 | 0.381 | 0.651 | 0.451 | 0.501 | 0.476 | 0.478 | 0.443 | 0.518 | 0.566 | 0.594 | 0.396 | 0.853 | 0.616 | 0.826 | 0.454 | 0.444 | 0.725 | 0.572 | 0.645 | 0.522 | 0.459 | 0.59 | 0.346 | 0.664 | 0.357 | 0.864 | 0.454 | 0.268 |
| GBR4 | * | NS | | 0.588 | 0.445 | 0.532 | 0.774 | 0.498 | 0.327 | 0.69 | 0.267 | 0.708 | 0.716 | 0.19 | 0.325 | 0.315 | 0.484 | 0.15 | 0.401 | 0.288 | 0.223 | 0.248 | 0.185 | 0.432 | 0.145 | 0.508 | 0.41 | 0.433 | 0.422 | 0.543 | 0.57 | 0.402 | 0.817 | 0.506 | 0.774 | 0.501 | 0.439 | 0.717 | 0.601 | 0.663 | 0.497 | 0.488 | 0.683 | 0.472 | 0.648 | 0.362 | 0.627 | 0.525 | 0.312 |
| BEL1 | * | NS | * | | 0.065 | 0.023 | 0.732 | 0.479 | 0.427 | 0.601 | 0.253 | 0.609 | 0.617 | 0.284 | 0.293 | 0.359 | 0.347 | 0.512 | 0.36 | 0.442 | 0.523 | 0.295 | 0.502 | 0.291 | 0.436 | 0.449 | 0.447 | 0.446 | 0.483 | 0.524 | 0.586 | 0.412 | 0.8 | 0.583 | 0.75 | 0.47 | 0.436 | 0.696 | 0.569 | 0.614 | 0.481 | 0.467 | 0.608 | 0.363 | 0.569 | 0.362 | 0.73 | 0.462 | 0.283 |
| BEL2 | * | NS | * | NS | | 0 | 0.711 | 0.438 | 0.363 | 0.571 | 0.195 | 0.582 | 0.588 | 0.193 | 0.24 | 0.288 | 0.296 | 0.396 | 0.24 | 0.361 | 0.38 | 0.156 | 0.356 | 0.249 | 0.278 | 0.39 | 0.395 | 0.374 | 0.393 | 0.465 | 0.525 | 0.359 | 0.779 | 0.536 | 0.705 | 0.425 | 0.359 | 0.673 | 0.508 | 0.558 | 0.394 | 0.314 | 0.57 | 0.327 | 0.523 | 0.287 | 0.683 | 0.422 | 0.198 |
| BEL3 | NS | NS | * | NS | NS | | 0.724 | 0.447 | 0.382 | 0.563 | 0.204 | 0.573 | 0.581 | 0.232 | 0.249 | 0.303 | 0.296 | 0.472 | 0.306 | 0.423 | 0.482 | 0.215 | 0.439 | 0.279 | 0.353 | 0.407 | 0.412 | 0.381 | 0.418 | 0.474 | 0.54 | 0.368 | 0.807 | 0.561 | 0.731 | 0.462 | 0.369 | 0.686 | 0.521 | 0.577 | 0.41 | 0.39 | 0.559 | 0.352 | 0.534 | 0.34 | 0.738 | 0.428 | 0.233 |
| FIN1 | * | NS | * | * | * | * | | 0.498 | 0.537 | 0.742 | 0.586 | 0.746 | 0.745 | 0.513 | 0.475 | 0.508 | 0.532 | 0.745 | 0.761 | 0.738 | 0.797 | 0.695 | 0.763 | 0.718 | 0.737 | 0.419 | 0.515 | 0.532 | 0.587 | 0.627 | 0.55 | 0.437 | 0.75 | 0.642 | 0.756 | 0.685 | 0.56 | 0.569 | 0.43 | 0.521 | 0.456 | 0.487 | 0.717 | 0.632 | 0.697 | 0.666 | 0.676 | 0.485 | 0.591 |
| RUS4* | * | * | * | * | * | * | * | | 0.309 | 0.484 | 0.33 | 0.506 | 0.51 | 0.311 | 0.3 | 0.286 | 0.334 | 0.462 | 0.416 | 0.482 | 0.5 | 0.41 | 0.434 | 0.442 | 0.437 | 0.291 | 0.301 | 0.215 | 0.191 | 0.354 | 0.367 | 0.262 | 0.555 | 0.38 | 0.462 | 0.433 | 0.286 | 0.494 | 0.304 | 0.28 | 0.113 | 0.231 | 0.495 | 0.367 | 0.455 | 0.371 | 0.522 | 0.27 | 0.317 |
| FIN2 | * | NS | * | * | * | * | * | * | | 0.488 | 0.225 | 0.526 | 0.521 | 0.191 | 0.142 | 0.125 | 0.235 | 0.286 | 0.302 | 0.325 | 0.395 | 0.286 | 0.314 | 0.312 | 0.302 | 0.284 | 0.142 | 0.166 | 0.161 | 0.212 | 0.295 | 0.172 | 0.649 | 0.271 | 0.482 | 0.271 | 0.182 | 0.442 | 0.289 | 0.28 | 0.137 | 0.168 | 0.484 | 0.265 | 0.264 | 0.206 | 0.448 | 0.193 | 0.159 |
| CZE1 | NS | NS | * | * | * | * | * | * | * | | 0.38 | 0.342 | 0.364 | 0.43 | 0.421 | 0.364 | 0.462 | 0.573 | 0.546 | 0.572 | 0.672 | 0.596 | 0.637 | 0.555 | 0.571 | 0.471 | 0.444 | 0.347 | 0.408 | 0.445 | 0.587 | 0.456 | 0.791 | 0.555 | 0.615 | 0.448 | 0.395 | 0.69 | 0.535 | 0.479 | 0.402 | 0.388 | 0.477 | 0.384 | 0.484 | 0.44 | 0.677 | 0.418 | 0.408 |
| GER2 | * | * | * | * | * | NS | * | * | * | NS | | 0.379 | 0.381 | 0.146 | 0.189 | 0.181 | 0.232 | 0.142 | 0.111 | 0.113 | 0.269 | 0.177 | 0.226 | 0.139 | 0.168 | 0.263 | 0.256 | 0.2 | 0.186 | 0.275 | 0.39 | 0.226 | 0.654 | 0.355 | 0.507 | 0.207 | 0.22 | 0.552 | 0.351 | 0.358 | 0.228 | 0.237 | 0.458 | 0.168 | 0.337 | 0.146 | 0.453 | 0.299 | 0.128 |
| GER3 | NS | NS | * | * | * | NS | * | * | NS | * | NS | | 0.113 | 0.445 | 0.445 | 0.397 | 0.48 | 0.579 | 0.543 | 0.567 | 0.673 | 0.61 | 0.649 | 0.542 | 0.57 | 0.502 | 0.492 | 0.402 | 0.454 | 0.492 | 0.609 | 0.489 | 0.805 | 0.589 | 0.642 | 0.438 | 0.441 | 0.708 | 0.532 | 0.499 | 0.435 | 0.412 | 0.472 | 0.411 | 0.54 | 0.467 | 0.691 | 0.47 | 0.435 |
| GER4 | NS | NS | * | * | * | NS | * | * | NS | * | NS | NS | | 0.442 | 0.443 | 0.399 | 0.465 | 0.584 | 0.553 | 0.569 | 0.687 | 0.612 | 0.661 | 0.546 | 0.575 | 0.488 | 0.487 | 0.387 | 0.45 | 0.494 | 0.61 | 0.486 | 0.812 | 0.593 | 0.657 | 0.439 | 0.435 | 0.697 | 0.54 | 0.501 | 0.435 | 0.405 | 0.492 | 0.415 | 0.542 | 0.481 | 0.703 | 0.463 | 0.431 |
| POL1 | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.105 | 0.074 | 0.191 | 0.182 | 0.202 | 0.242 | 0.21 | 0.153 | 0.175 | 0.237 | 0.105 | 0.218 | 0.195 | 0.194 | 0.183 | 0.246 | 0.3 | 0.187 | 0.587 | 0.317 | 0.477 | 0.235 | 0.186 | 0.487 | 0.259 | 0.298 | 0.156 | 0.161 | 0.426 | 0.194 | 0.314 | 0.138 | 0.356 | 0.246 | 0.11 |
| POL2 | * | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.061 | 0.113 | 0.292 | 0.253 | 0.317 | 0.358 | 0.237 | 0.298 | 0.243 | 0.242 | 0.241 | 0.148 | 0.149 | 0.169 | 0.111 | 0.219 | 0.146 | 0.598 | 0.266 | 0.417 | 0.244 | 0.112 | 0.438 | 0.228 | 0.239 | 0.125 | 0.114 | 0.427 | 0.203 | 0.184 | 0.157 | 0.422 | 0.17 | 0.124 |
| POL3 | * | NS | * | * | NS | * | * | * | * | * | * | * | * | * | * | | 0.142 | 0.31 | 0.271 | 0.368 | 0.392 | 0.234 | 0.274 | 0.294 | 0.227 | 0.246 | 0.16 | 0.16 | 0.185 | 0.214 | 0.283 | 0.154 | 0.642 | 0.253 | 0.448 | 0.26 | 0.154 | 0.456 | 0.197 | 0.261 | 0.086 | 0.057 | 0.355 | 0.203 | 0.268 | 0.155 | 0.427 | 0.194 | 0.117 |
| POL4 | * | NS | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | | 0.416 | 0.301 | 0.418 | 0.491 | 0.323 | 0.413 | 0.281 | 0.358 | 0.263 | 0.285 | 0.184 | 0.246 | 0.24 | 0.34 | 0.22 | 0.69 | 0.391 | 0.547 | 0.344 | 0.211 | 0.446 | 0.269 | 0.269 | 0.204 | 0.177 | 0.464 | 0.266 | 0.286 | 0.261 | 0.53 | 0.257 | 0.202 |
| GBR7 | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | NS | | 0.153 | 0.072 | 0.364 | 0.164 | 0.244 | 0.286 | 0.134 | 0.497 | 0.405 | 0.388 | 0.391 | 0.514 | 0.529 | 0.321 | 0.8 | 0.452 | 0.74 | 0.355 | 0.426 | 0.685 | 0.542 | 0.63 | 0.424 | 0.406 | 0.608 | 0.37 | 0.606 | 0.277 | 0.637 | 0.499 | 0.279 |
| GBR3 | NS | NS | * | NS | NS | NS | * | * | NS | NS | NS | * | * | * | * | * | NS | NS | | 0.021 | 0.422 | 0.09 | 0.336 | 0.097 | 0.22 | 0.435 | 0.387 | 0.322 | 0.343 | 0.479 | 0.525 | 0.297 | 0.827 | 0.516 | 0.751 | 0.284 | 0.673 | 0.509 | 0.573 | 0.396 | 0.394 | 0.592 | 0.232 | 0.591 | 0.182 | 0.752 | 0.442 | 0.175 |
| GBR8 | * | * | * | * | * | * | * | * | NS | NS | * | * | * | * | * | * | * | NS | * | | 0.42 | 0.184 | 0.31 | 0.181 | 0.22 | 0.518 | 0.444 | 0.426 | 0.424 | 0.534 | 0.561 | 0.356 | 0.784 | 0.479 | 0.734 | 0.301 | 0.464 | 0.686 | 0.564 | 0.636 | 0.453 | 0.447 | 0.631 | 0.332 | 0.605 | 0.254 | 0.661 | 0.524 | 0.287 |
| GBR9 | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | | 0.205 | 0.021 | 0.38 | 0.159 | 0.577 | 0.483 | 0.528 | 0.517 | 0.661 | 0.621 | 0.458 | 0.841 | 0.528 | 0.814 | 0.61 | 0.529 | 0.728 | 0.651 | 0.723 | 0.519 | 0.495 | 0.652 | 0.553 | 0.751 | 0.504 | 0.757 | 0.608 | 0.395 |
| GBR11 | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | | 0.178 | 0.235 | 0.138 | 0.369 | 0.346 | 0.342 | 0.336 | 0.438 | 0.475 | 0.285 | 0.746 | 0.509 | 0.689 | 0.368 | 0.46 | 0.542 | 0.342 | 0.384 | 0.603 | 0.287 | 0.52 | 0.211 | 0.584 | 0.418 | 0.161 |
| GBR5 | NS | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | | 0.339 | 0.161 | 0.452 | 0.367 | 0.387 | 0.365 | 0.538 | 0.555 | 0.375 | 0.819 | 0.489 | 0.759 | 0.489 | 0.398 | 0.681 | 0.561 | 0.604 | 0.401 | 0.415 | 0.619 | 0.422 | 0.645 | 0.366 | 0.655 | 0.483 | 0.27 |
| GBR6 | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | | 0.278 | 0.452 | 0.39 | 0.358 | 0.36 | 0.463 | 0.533 | 0.366 | 0.773 | 0.474 | 0.686 | 0.293 | 0.387 | 0.634 | 0.513 | 0.519 | 0.398 | 0.413 | 0.599 | 0.272 | 0.511 | 0.235 | 0.666 | 0.429 | 0.228 |
| GBR10 | * | NS | * | * | * | * | * | * | NS | NS | * | * | * | * | * | NS | NS | * | * | * | * | * | * | NS | | 0.403 | 0.352 | 0.332 | 0.346 | 0.447 | 0.478 | 0.325 | 0.787 | 0.469 | 0.662 | 0.481 | 0.537 | 0.378 | 0.365 | 0.571 | 0.54 | 0.284 | 0.583 | 0.427 | 0.221 |
| SWE4 | * | NS | * | * | NS | * | * | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | | 0.233 | 0.235 | 0.192 | 0.363 | 0.329 | 0.176 | 0.652 | 0.357 | 0.578 | 0.458 | 0.224 | 0.436 | 0.224 | 0.294 | 0.132 | 0.222 | 0.47 | 0.351 | 0.473 | 0.378 | 0.507 | 0.294 | 0.25 |
| SWE3 | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | | 0.166 | 0.115 | 0.188 | 0.209 | 0.134 | 0.652 | 0.116 | 0.378 | 0.358 | 0.117 | 0.465 | 0.307 | 0.292 | 0.083 | 0.171 | 0.447 | 0.32 | 0.315 | 0.292 | 0.519 | 0.175 | 0.205 |
| SWE5 | NS | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | NS | | 0.084 | 0.168 | 0.248 | 0.115 | 0.625 | 0.168 | 0.404 | 0.337 | 0.103 | 0.462 | 0.235 | 0.214 | 0.064 | 0.113 | 0.378 | 0.274 | 0.26 | 0.271 | 0.513 | 0.137 | 0.191 |
| FIN6 | NS | NS | * | NS | NS | NS | * | * | NS | * | * | * | * | NS | NS | * | * | * | * | * | NS | NS | NS | NS | NS | * | * | NS | | 0.258 | 0.295 | 0.141 | 0.687 | 0.135 | 0.429 | 0.385 | 0.127 | 0.532 | 0.32 | 0.29 | 0.08 | 0.175 | 0.426 | 0.311 | 0.411 | 0.294 | 0.63 | 0.229 | 0.187 |
| SWE7 | NS | NS | * | * | NS | * | * | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | NS | NS | * | NS | | 0.205 | 0.141 | 0.77 | 0.279 | 0.501 | 0.406 | 0.12 | 0.555 | 0.37 | 0.36 | 0.253 | 0.195 | 0.463 | 0.362 | 0.15 | 0.345 | 0.641 | 0.201 | 0.297 |
| SWE2 | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.129 | 0.695 | 0.389 | 0.515 | 0.491 | 0.329 | 0.334 | 0.266 | 0.29 | 0.585 | 0.433 | 0.448 | 0.435 | 0.567 | 0.235 | 0.361 | | |
| SWE1 | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.589 | 0.318 | 0.439 | 0.281 | 0.136 | 0.389 | 0.193 | 0.205 | 0.108 | 0.157 | 0.489 | 0.2 | 0.29 | 0.201 | 0.368 | 0.168 | 0.174 |
| SWE9 | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.721 | 0.838 | 0.768 | 0.686 | 0.753 | 0.706 | 0.756 | 0.65 | 0.699 | 0.776 | 0.734 | 0.829 | 0.753 | 0.828 | 0.661 | 0.702 |
| SWE10 | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | | 0.444 | 0.437 | 0.176 | 0.578 | 0.477 | 0.374 | 0.331 | 0.558 | 0.419 | 0.407 | 0.386 | 0.16 | 0.3 | | |
| SWE11 | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | NS | | 0.64 | 0.36 | 0.701 | 0.642 | 0.636 | 0.381 | 0.442 | 0.62 | 0.605 | 0.64 | 0.642 | 0.851 | 0.378 | 0.545 |
| SWE8 | NS | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | | 0.387 | 0.604 | 0.429 | 0.472 | 0.342 | 0.309 | 0.532 | 0.181 | 0.441 | 0.159 | 0.638 | 0.363 | 0.233 |
| FIN5 | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | NS | * | NS | * | * | NS | * | NS | * | * | * | * | * | * | NS | | 0.48 | 0.288 | 0.265 | 0.118 | 0.155 | 0.405 | 0.314 | 0.271 | 0.295 | 0.561 | 0.178 | 0.199 |
| FIN3 | * | * | * | * | * | NS | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | | 0.402 | 0.357 | 0.43 | 0.423 | 0.687 | 0.566 | 0.606 | 0.605 | 0.657 | 0.464 | 0.528 |
| FIN4 | NS | NS | * | * | NS | * | * | * | NS | * | NS | * | * | * | * | * | * | NS | * | NS | NS | NS | * | * | NS | NS | NS | NS | * | * | * | * | * | NS | NS | * | NS | * | | 0.23 | 0.166 | 0.165 | 0.519 | 0.34 | 0.458 | 0.361 | 0.462 | 0.267 | 0.269 |
| EST1 | NS | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | | 0.158 | 0.191 | 0.474 | 0.39 | 0.461 | 0.452 | 0.674 | 0.234 | 0.34 |
| EST2 | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | NA | | 0.042 | 0.367 | 0.253 | 0.356 | 0.274 | 0.504 | 0.119 | 0.146 |
| BLS | NS | NS | * | * | NS | * | * | * | NS | * | NS | * | * | * | * | * | * | NS | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | NS | * | NS | NS | NA | | 0.364 | 0.245 | 0.269 | 0.265 | 0.406 | 0.138 | 0.181 |
| RUS1 | NS | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | NS | NS | NS | NS | NS | NS | NS | NS | NS | | 0.483 | 0.467 | 0.503 | 0.646 | 0.405 | 0.437 |
| DEN1 | NS | NS | * | * | * | * | * | * | NS | * | * | * | * | NS | * | * | * | NS | * | * | * | * | NS | * | * | * | * | * | NS | * | * | * | * | NS | NS | NS | NS | NS | NS | NA | NS | NS | NS | | 0.428 | 0.121 | 0.588 | 0.316 | 0.132 |
| SWE12 | NS | NS | * | * | NS | NS | * | * | NS | * | * | * | * | NS | * | * | * | NS | * | NS | * | NS | NS | NS | NS | NS | NS | NS | NS | NS | * | * | * | NS | NS | NS | NS | NS | NS | NA | NS | NS | NS | NS | | 0.431 | 0.723 | 0.265 | 0.359 |
| DEN2 | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | * | NS | * | NS | * | * | NA | * | * | * | NS | NS | | 0.611 | 0.335 | 0.099 |
| NOR2 | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | NA | * | * | * | NS | * | NS | | 0.532 | 0.541 |
| SWE14 | * | NS | * | * | * | * | * | * | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | * | * | * | * | NS | NS | * | NS | * | NS | NA | NS | NS | NS | NS | * | * | * | | 0.261 |
| DEN3 | * | NS | * | * | * | * | * | * | NS | NS | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | * | NS | NS | * | NS | * | * | NA | * | * | * | NS | * | * | NS | * | * | |

196

197

198    SI table 5. Pairwise FST values calculated using the RADseq dataset.

| | GBR8 | BEL1 | GBR4 | FIN3 | DEN1 | GBR7 | SWE12 | FIN4 | DEN2 | POL4 | RUS1 | SWE2 | SWE8 | SWE9 | SWE10 | NOR2 | POL3 | HUN2 | WEN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GBR8 | | 0.34971 | 0.35695 | 0.49475 | 0.223897 | 0.35613 | 0.406544 | 0.295019 | 0.293628 | 0.211876 | 0.38775 | 0.308973 | 0.273693 | 0.412263 | 0.321365 | 0.650207 | 0.146766 | 0.61801 | 0.397239 |
| BEL1 | | | 0.370425 | 0.390916 | 0.098308 | 0.381154 | 0.300836 | 0.225496 | 0.130398 | 0.152617 | 0.343954 | 0.22423 | 0.08032 | 0.326848 | 0.235947 | 0.522507 | 0.103445 | 0.597677 | 0.31111 |
| GBR4 | | | | 0.513779 | 0.231153 | 0.195241 | 0.423664 | 0.302246 | 0.316185 | 0.218776 | 0.392539 | 0.314975 | 0.284155 | 0.422534 | 0.331921 | 0.698989 | 0.149208 | 0.620806 | 0.412409 |
| FIN3 | | | | | 0.308284 | 0.517114 | 0.341754 | 0.198275 | 0.364426 | 0.222674 | 0.378729 | 0.27048 | 0.328488 | 0.331267 | 0.286862 | 0.562015 | 0.149991 | 0.614832 | 0.341565 |
| DEN1 | | | | | | 0.244594 | 0.239562 | 0.194342 | 0.06762 | 0.136982 | 0.356985 | 0.182005 | 0.085513 | 0.266461 | 0.190793 | 0.362014 | 0.102429 | 0.602815 | 0.237037 |
| GBR7 | | | | | | | 0.430574 | 0.31162 | 0.32391 | 0.229621 | 0.396753 | 0.319608 | 0.295939 | 0.433712 | 0.340292 | 0.692819 | 0.157339 | 0.621918 | 0.422803 |
| SWE12 | | | | | | | | 0.209406 | 0.282835 | 0.173199 | 0.363912 | 0.198857 | 0.259513 | 0.303204 | 0.211775 | 0.459576 | 0.122381 | 0.606576 | 0.250115 |
| FIN4 | | | | | | | | | 0.218225 | 0.142389 | 0.328888 | 0.154803 | 0.211809 | 0.203425 | 0.174944 | 0.316929 | 0.099233 | 0.586541 | 0.198636 |
| DEN2 | | | | | | | | | | 0.153556 | 0.362177 | 0.212179 | 0.101801 | 0.307702 | 0.222051 | 0.459347 | 0.108029 | 0.60623 | 0.284015 |
| POL4 | | | | | | | | | | | 0.321777 | 0.128672 | 0.150743 | 0.192186 | 0.138734 | 0.250273 | 0.073063 | 0.579543 | 0.161299 |
| RUS1 | | | | | | | | | | | | 0.341129 | 0.358602 | 0.368371 | 0.349288 | 0.396194 | 0.278006 | 0.516158 | 0.358584 |
| SWE2 | | | | | | | | | | | | | 0.19768 | 0.218326 | 0.145195 | 0.325228 | 0.094924 | 0.591579 | 0.151258 |
| SWE8 | | | | | | | | | | | | | | 0.289356 | 0.208013 | 0.401551 | 0.110433 | 0.604134 | 0.262799 |
| SWE9 | | | | | | | | | | | | | | | 0.257245 | 0.429544 | 0.136442 | 0.607715 | 0.29715 |
| SWE10 | | | | | | | | | | | | | | | | 0.350951 | 0.100275 | 0.598136 | 0.184722 |
| NOR2 | | | | | | | | | | | | | | | | | 0.165304 | 0.625602 | 0.426179 |
| POL3 | | | | | | | | | | | | | | | | | | 0.547371 | 0.111018 |
| HUN2 | | | | | | | | | | | | | | | | | | | 0.604399 |

199

Scenario 1 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

Pool 1  Pool 2  Pool 3

t2, t1, 0

Scenario 2 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

Pool 1  Pool 2  Pool 3

t2, t1, 0

Scenario 7 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

rc  1-rc

Pool 1  Pool 3  Pool 2

t2, t1, 0

Scenario 8 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3
- Nb Pool 1

Pool 1  Pool 2  Pool 3

t2, t1, t1-db, 0

Scenario 3 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

Pool 1  Pool 3  Pool 2

t2, t1, 0

Scenario 4 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

Pool 1  Pool 3  Pool 2

t2, t1, 0

Scenario 9 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3
- Nb Pool 1

Pool 1  Pool 2  Pool 3

t2, t1, t1-db, 0

Scenario 10 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3
- Nb Pool 1

Pool 1  Pool 2  Pool 3

t2, t2-db, t1, 0

Scenario 5 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

ra  1-ra

Pool 2  Pool 1  Pool 3

t2, t1, 0

Scenario 6 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3

1-rb  rb

Pool 3  Pool 2  Pool 1

t2, t1, 0

Scenario 11 (Warning ! Time is not to scale.)
- N Pool 1
- N Pool 2
- N Pool 3
- Nb Pool 1

Pool 3  Pool 1  Pool 2

t2, t2-db, t1, 0

Pool 1 - Northern Europe
Pool 2 - Don river
Pool 3 - Danube River

## a) Stage 2. NEU Major variants



Pool 1 - UK
Pool 2 - Denmark / S.Sweden
Pool 3 - Finland
Pool 4 - Poland
Pool 5 - Sweden
Pool 6 - Tromsø

## b) Stage 3. Scenario 14 Minor variants

a) Covereage per tag before Ho filtering

mean coverage = 29.072
n Tags = 18908

b) Covereage per tag after Ho filtering

mean coverage = 27.722
n Tags = 13189

c) Per-sample average coverage for all retained tags

a) All samples — Allelic Richness vs Latitude, $R^2 = -0.007$

b) All samples — Allelic Richness vs Longitude, $R^2 = 0.2887$ ***

c) mtDNA Lineage 1 only — Allelic Richness vs Latitude, $R^2 = -0.023$

d) mtDNA Lineage 1 only — Allelic Richness vs Longitude, $R^2 = 0.3156$ ***

a) M1, Lineage 1 only (excluding NOR2) — adj. R squared = 0.287 ***

b) RADseq data IBD — adj. R squared = 0.722***

c) M2 data IBD — adj. R squared = 0.455***

d) M3 data IBD — adj. R squared = 0.447***

a)



b)



c)

a) Stage 2. NEU major variants

b) Stage 3. NEU minor variants
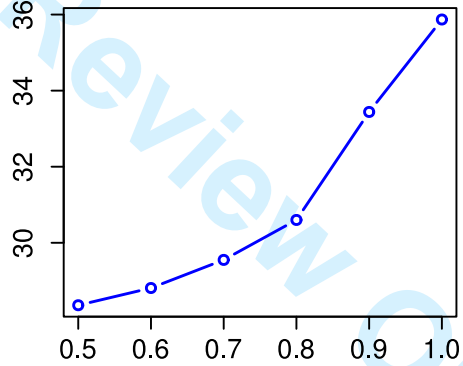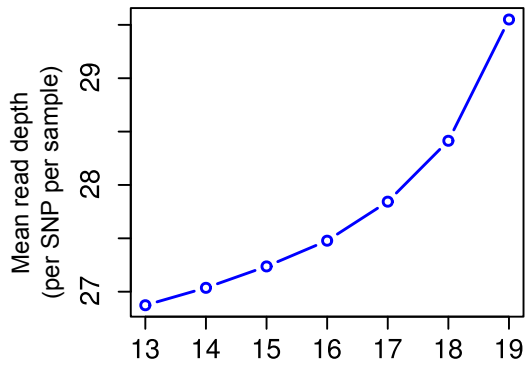
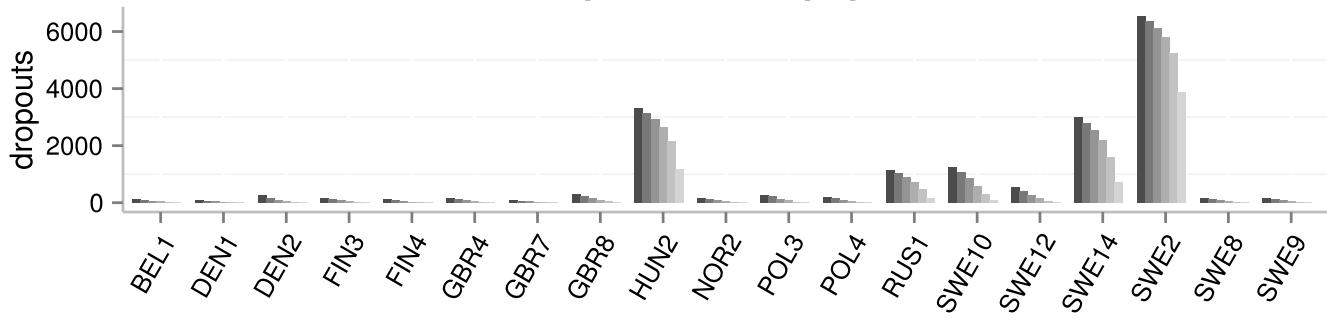# Spatial patterns in M2 vs M1 datasets



## M1 Gest

## M2 Gest

## M1 Lest

## M2 Lest

a) Change in number of tags with incrementing M value

b) change in average tag coverage with incrementing M value

a)

b)

c)

SNP dropout across populations