

Estimating recent and historical effective population size of marine and freshwater sticklebacks

Xueyun Feng^{1,2}, Ari Löytynoja², Juha Merilä^{1,3}

¹*Organismal and Evolutionary Biology Programme, University of Helsinki, FI-00014 University of Helsinki, Finland*

²*Institute of Biotechnology, University of Helsinki, FI-00014 University of Helsinki, Finland*

³*Area of Ecology and Biodiversity, School of Biological Sciences, University of Hong Kong, Hong Kong SAR*

Correspondence: xueyung.feng@helsinki.fi

Running head: Stickleback Effective Population Size

Abstract

Effective population size (N_e) is a quantity of central importance in evolutionary biology and population genetics, but often notoriously challenging to estimate. Analyses of N_e are further complicated by the many interpretations of the concept and the alternative approaches to quantify N_e utilising widely different properties of the data. On the other hand, alternative methods are informative for different time scales such that a set of complementary methods should allow piecing together the entire continuum of N_e from a few generations before the present to the distant past. To test this in practice, we inferred the continuum of N_e for 45 nine-spined stickleback populations (*Pungitius pungitius*) using whole-genome data. We found that the marine populations had the largest historical and recent N_e , followed by coastal and other freshwater populations. We identified the impact of both recent and historical gene flow on the N_e estimates obtained from different methods and found that simple summary statistics are informative in comprehending the events in the very recent past. Overall, our analyses showed that the coalescence-based trajectories of N_e in the recent past and the LD-based estimates of near-contemporary N_e are incongruent, though in some cases the incongruence might be explained by specific demographic events. Despite still lacking accuracy and resolution for the very recent past, the sequentially Markovian coalescent-based methods seem to provide the most meaningful interpretation of the real-life N_e varying across time.

Keywords: admixture, demographic history, effective population size, nine-spined stickleback, *Pungitius*

Introduction

By quantifying the magnitude of genetic drift and inbreeding in real-world populations, the concept of effective population size (N_e) has many applications in evolutionary (Charlesworth, 2009; Charlesworth & Charlesworth, 2010) and conservation biology (Frankham et al., 2010; Allendorf et al., 2012). In evolutionary biology, N_e is informative about the efficacy of selection, mutation, and gene flow as systematic evolutionary forces (Charlesworth & Charlesworth, 2010). In conservation biology, N_e is informative of the evolutionary potential and fitness of populations (Frankham et al., 2010; Hare et al., 2011; R. S. Waples, 2022). Numerous approaches have been developed to estimate N_e from genetic data (J. Wang, 2005; Palstra & Ruzzante, 2008; Luikart et al., 2010; Gilbert & Whitlock, 2015; Santiago et al., 2020). Some of these approaches seek to estimate the dynamics of historical N_e (reviewed in Beichman et al., 2018), whereas others are designed to estimate the contemporary N_e (Waples & Do, 2010). Despite the progress, estimation of contemporary N_e is still challenging, especially for large populations (Marandel et al., 2019) where the detection of signs of drift and inbreeding requires very large sample sizes (Waples et al., 2016; Marandel et al., 2019). As a result, obtaining accurate estimates of contemporary N_e for marine organisms with large population sizes is considered to be next to impossible (Hare et al., 2011; Marandel et al., 2019).

All methods used to estimate N_e from genetic data make assumptions, and violation of these assumptions may lead to errors and biases (Beichman et al., 2018; Nadachowska-Brzyska et al., 2022). For instance, most N_e estimation approaches assume populations to be closed (e.g., Waples & Do, 2010; but see also Santiago et al., 2020), while in real life, most populations are affected by at least some level of migration. The methods based on sequential Markovian Coalescent (Li & Durbin, 2011; Schiffels & Durbin, 2014) commonly used to estimate dynamics of the historical N_e are not immune to the effects of gene flow either, and changes in gene flow can yield N_e trajectories that mimic changes in population size (Beichman et al., 2018). In a simulation study of human demographic history, Hawks (2017) showed that even a small fraction of introgression or gene flow in the distant past can have a visible effect on the inference of N_e . Although studied theoretically and with simulations, empirical studies on the effects of introgression on the dynamics of historical and recent N_e are still lacking.

Studies of genetic variability in fish have shown that freshwater populations harbour less genetic variation than marine populations (e.g., Ward et al., 1994; DeWoody & Avise, 2000; DeFaveri & Merilä, 2015; Kivikoski et al., 2023), indicating smaller long-term N_e (N_e^{LT}) for freshwater than for marine fish. While this makes intuitive sense due to the limited size and fragmentation of freshwater habitats as compared to more continuous marine environments, the conjecture suffers from two problems. First, due to the difficulties in estimation of contemporary N_e for large marine populations (cf. Waples et al., 2016; Marandel et al., 2019), rigorous comparisons of marine and freshwater populations are rare (e.g., DeFaveri & Merilä, 2015). The recent development in Linkage Disequilibrium (LD) -based methods for the estimation of temporal changes in N_e might provide a solution for this, and the new methods have been shown to provide robust estimates even for populations of relatively large N_e (Santiago et al., 2020). Second, since species' life history traits have been implicated to affect their genetic diversities (Romiguier et al., 2014), any comparison of N_e

derived from genetic data should ideally control for the life history differences between marine and freshwater environments. One way to resolve this is to study euryhaline species which can reproduce with similar life histories both in freshwater and marine habitats.

The nine-spined stickleback (*Pungitius pungitius*) is a small euryhaline teleost fish with a circumpolar distribution. What makes it particularly interesting in the context of N_e estimation is its presence in both open marine habitats as well as in landlocked freshwater habitats. Previous studies have also found evidence for ancient and potentially ongoing introgression in some parts of its distribution range (Guo et al., 2019; Feng et al., 2022; Y. Wang et al., 2023). Together these two properties make the species well-suited to explore the dynamics of N_e in different habitats and the effects of gene flow on N_e estimates.

The aims of this study were three-fold. First, by making use of a large collection ($n = 888$) of nine-spined stickleback whole genome sequences, we reconstructed the dynamics of both recent and historical N_e in 45 marine and freshwater populations to assess whether the estimates of near-contemporary N_e (N_e^{NC} ; 1-200 generations ago) and recent N_e (N_e^R ; a few hundred to thousand generations ago) are higher in marine than in freshwater populations. Second, using the information on levels of introgression between two divergent nine-spined stickleback lineages, we studied whether introgression has influenced the temporal dynamics of N_e in the populations affected by it. Third, given the difficulty of estimating contemporary N_e in large marine populations, we explored whether N_e^{NC} or N_e^R estimates could provide reasonable proxies of contemporary N_e . For that, we compared the different estimates of N_e across populations, focusing specifically on the small isolated pond populations for which contemporary N_e can most reliably be estimated.

Methods

Study populations and sampling

The samples used in this study were collected in accordance with the national legislation of the countries concerned. The data used in this study originate from 45 *P. pungitius* populations covering much of the species distribution area in Eurasia, North America and the Far East (Table S1). Of these populations, 12 were from marine and nine from coastal freshwater populations with connection (or recent connection) to the sea. Of the true freshwater populations, eleven were from closed ponds (surface area < 4 ha), ten from lakes, two from rivers and one from a man-made drainage ditch (Table S1). The fish from all populations were collected during the species' local breeding season using either minnow traps or beach seine. Sampled fish were euthanized with MS-222 and preserved in ethanol until DNA extractions.

DNA extractions and sequencing

The genomic DNA was extracted from alcohol-preserved fin clips either the conventional phenol-chloroform method as described in Sambrook & Russell (2006) or by salting-out

method (Miller et al., 1988, Bruford et al., 1998). From the resulting DNA, sequencing libraries were prepared with insert size of 300-350 bp, and 150-bp paired-end reads were generated using Illumina HiSeq 2500/4000 instruments. The library preparations and sequencing procedures were performed at two different facilities: the Beijing Genomics Institute (Hong Kong SAR, China) and the DNA Sequencing and Genomics Laboratory at the University of Helsinki (Helsinki, Finland).

Data processing

The short-read data were mapped to the latest nine-spined stickleback reference genome (Kivikoski et al., 2021) using the Burrows-Wheeler Aligner v.0.7.17 (BWA MEM algorithm; Li, 2013) and its default parameters. Duplicate reads were marked with samtools v.1.7 (Li et al., 2009) and variant calling was performed with the Genome Analysis Toolkit (GATK) v.3.6.0 and v.4.0.1.2 (McKenna et al., 2010) following the GATK Best Practices workflows. In more detail, RealignerTargetCreator and IndelRealigner (from v.3.6.0) tools were applied to realign reads around indels, HaplotypeCaller was used to call variants for each individual (parameters set as -stand emit conf 3, -stand call cof 10, - GQB (10,50), variant index type linear and variant index parameter 128000), and finally GenotypeGVCFs was used to jointly call the variants for all the samples using its default parameters. Binary SNPs were extracted with bcftools v.1.7 (Danecek et al., 2021) excluding sites located within identified repetitive sequences (Varadharajan et al., 2019) and negative mappability mask regions combining the identified repeats and unmappable regions (Kivikoski et al., 2021). Sites showing low (<8x) or too high (>25x) average coverage, low (<20) genotype quality, low (<30) quality score and more than 25% missing data were filtered out using vcftools v.0.1.5 (Danecek et al., 2011). Data from the known sex chromosomes (LG12) were removed from further analysis. For details about the subsequent filtering of the dataset used in different analyses, see Table S2.

Analyses of linkage disequilibrium and genetic relatedness

The magnitude of linkage disequilibrium (LD) and its decay are informative on N_e , level of inbreeding, and migration (Flint-Garcia et al., 2003). Hence, we characterised LD patterns in all populations by estimating the squared correlation coefficient r^2 between each pair of SNPs with PopLDdecay (Zhang et al., 2019) with its default settings. We restricted the analysis to the largest linkage group LG4 and used LG1 for cross-validation (SNP Set 2 in Table S2). The LD decay curve was plotted with R (R Core Team, 2020).

High levels of LD in a population may indicate (i) small N_e , (ii) increased inbreeding, and/or (iii) recent migration/admixture. To distinguish between these, we first estimated the average inbreeding coefficients (F_{IS}) for each population using vcftools --het. We then used ngsRelate v.2 (Hanghøj et al., 2019) to calculate the r_{xy} , the pairwise relatedness within populations (Hedrick & Lacy, 2015). As a measure of temporal gene flow, we examined the LD decay patterns. Within the same ecotype, populations showing atypical LD decay patterns were considered as potentially affected by temporal gene flow.

GONE analyses

We estimated the near-contemporary N_e (N_e^{NC}) using GONE (Santiago et al., 2020). This method utilises the LD patterns in the data and has been shown to be robust for time spans of 0-200 generations before present, even when N_e is relatively large (Santiago et al., 2020). The analyses were performed using the SNP Set1 (Table S2) along with a genetic map lifted-over from the reference genome version 6 (Varadharajan et al., 2019) to version 7 (Kivikoski et al., 2021). According to Santiago et al. (2020), the possible bias from recent gene flow can be mitigated by limiting the recombination fractions threshold (hc). Following that, we repeated the analyses using a hc of 0.01 and 0.05. For each population, twenty independent replicates were performed with the default settings. The harmonic mean of estimates for generations 15-50 before the present was taken as the estimate of near-contemporary N_e^{NC} .

MSMC2 analyses

MSMC2 (Malaspinas et al., 2016) was used to estimate the recent and historical N_e . As the method can (with a reasonable runtime) analyse at most eight haplotypes, we selected and utilised the four individuals with the highest sequencing coverage from each population. The input files were generated following Schiffels & Wang (2020), and along with mask files generated by bamCaller.py, the mappability masks (Kivikoski et al., 2021) were applied. Estimates were carried out with default settings, and the outputs were processed assuming mutation rate of 4.37×10^{-9} per site per generation (Zhang et al., 2023) and a generation length of two years (DeFaveri et al., 2014). To conduct bootstrap estimations, the input data were chopped into 1 Mb blocks and an artificial 400 Mb long genome was generated by random sampling with replacement. 20 artificial bootstrap datasets were generated using Multihetsep_bootstrap.py from msmc-tools (<https://github.com/stschiff/msmc-tools>) and analysed with the same settings as the original data. In all analyses, the first two time segments (which usually are untrustworthy; Schiffels & Durbin, 2014; Sellinger et al., 2021) were discarded and the value for the third most recent time segment was used as the estimate of N_e^R , the recent N_e .

Long-term N_e estimation

The average long-term N_e (N_e^{LT}) were estimated using the formula (Kimura, 1983) $N_e = \pi / (4\mu)$, where μ is the mutation rate, assumed to be 4.37×10^{-9} per site per generation. π was obtained from folded site frequency spectra (SFS), estimated for each population directly from the bam data with ANGSD v.0.921 (Malaspinas et al., 2016), using the R script from Walsh et al., (2022) modified to fit folded SFSs. Sites with more than 70% heterozygote counts were removed and the mappability masks (Kivikoski et al., 2021) were applied in data filtering. For details, see Table S2.

Results

Summary statistics show unexpected variation among freshwater populations

We found the genetic diversity to be generally higher in marine and coastal freshwater than in inland freshwater populations and decreasing together with the connectivity of the habitat class (Fig. 1a). However, the patterns of LD decay over physical distance showed marked differences within and between habitat classes (Fig. 1b). While the decay of LD is faster and shows lower average LD (r^2) in marine than in freshwater populations, the latter show considerable variation, some being similar to marine populations with low levels of LD and others containing very high levels of LD (Fig. 1b). Despite being estimated with the full set of filtered SNPs, the estimates of inbreeding coefficients (F_{IS}) are meaningful for comparison of differences across the study populations. We found the marine populations to have generally low F_{IS} , while those for the freshwater populations were highly variable (Fig. 1c). Since a few populations showed higher than average within-population variation in F_{IS} , we estimated the pairwise relatedness (r_{xy} , Hedrick et al. 2015) within each population. The relatedness showed an increase with the degree of isolation of the habitat, and unexpectedly, some freshwater samples were more closely related than the others ($r_{xy} > 0.5$ shown in yellow colour in Fig. 1d), forming deme-like sub-populations within the given site (Fig. S1). The initial analysis revealed that the samples from Lake Riiokjärvi, Finland, show an exceptionally high level of LD and strong patterns of inbreeding in comparison to other lake populations, and certain individuals within the population were more closely related than others.

Marine and freshwater populations show distinct near-contemporary demographic histories

We inferred the near-contemporary demographic histories using the LD-based method GONE (Santiago et al., 2020). Estimates of the near-contemporary N_e (N_e^{NC}) were generally smaller in freshwater than in marine populations (Fig. 2a), consistent with the observed variability in the level of LD and F_{IS} . The smallest N_e^{NC} were obtained for pond populations (Lund, Sweden and Pyöreälampi, Finland), and the largest N_e^{NC} for marine populations and two lake populations (Lake Floatingstone, Canada and Lake Ukonjärvi, Finland; Fig. 2). However, our sample from the Lake Floatingstone population is smaller than others (eight individuals vs. 20+ individuals from most other populations) and the high estimate may be unreliable (Santiago et al., 2020).

A closer look at the N_e^{NC} results revealed a small number of populations with sharp drops in the estimated N_e during the last 20 generations and abnormally low N_e (<10) around 100-120 generations ago (Fig. S3). Such patterns can be expected if the level of inbreeding or LD has recently increased due to migration, population admixture or population bottlenecks (Santiago et al., 2020). Consistent with this, a comparison of N_e^{NC} estimates under different recombination bin cut-offs (hc) revealed either inconsistent estimates or very large ranges of values – the typical symptoms of recent gene flow – for populations with abnormal LD, inbreeding and/or relatedness patterns (Fig. 2; Santiago et al., 2020). An extreme example

of this is FIN-RII whose N_e^{NC} estimates vary between 14 and 924,485 (Fig. 2): the same population is the clearest outlier in the LD decay analysis (Fig. 1b). On the other hand, a majority of populations showed steady estimates of N_e^{NC} under different values of hc , indicating a stable recent history with no gene flow (Fig. 2).

Historical N_e reveals aberrations in specific freshwater populations

We estimated historical N_e with MSMC2 and found all N_e trajectories showing declining trends until the Last Glacial Maximum (LGM, ca. 20,000 years ago; Fig. 3). After that, the N_e of marine populations has grown (Fig. 3 a) while that of most freshwater populations has continued to decrease (Fig. 3c-d). The resolution towards the very recent times (~100-1000 years ago) is lost for all but the smallest freshwater populations. Within ecotypes, specific populations show distinct – but predictable – deviations from the consensus history of N_e . On one hand, the marine population from Hokkaido, Japan (Fig. 3a), and the freshwater population from Lake Floating Stone, Canada (Fig. 3c), stand out within the ecotypes and thus reflect their geographic isolation and the regional differences in climate during the Last Glacial Period. On the other hand, specific coastal freshwater and pond populations show N_e trajectories similar to those of the European marine populations (Fig. 3b), indicating recent colonisation of the freshwater environment. In general, the demographic histories of freshwater populations are more variable than those of marine populations (Fig. 3). The patterns reflect well the sizes and connectivities of the habitats and the known differences in their regional geographic histories.

Different N_e estimation methods show poor correlation

The three effective population size estimates, N_e^R , N_e^{LT} and N_e^{NC} , were obtained using different methods and reflect alternative definitions of N_e , and their similarities and differences are of interest. Of the three estimators, N_e^{NC} is the most variable (Fig. 4a) within ecotypes, possibly because of its greater sensitivity to the violations of the expected closed population history. The estimates of N_e^R and N_e^{LT} differ in magnitude but show consistency across the ecotypes and populations (Fig. 4b and 4c).

As we excluded the last two time segments of the MSMC2 trajectories as untrustworthy, the ages of N_e^R estimates for different populations vary from 76.61 to 3346.65 generations ago. One would expect that the "near-contemporary" (N_e^{NC}) and "recent" (N_e^R) estimates of N_e would be more similar when the time periods for the estimates are temporally closer. That is not the case (Fig. 4d) and there is no correlation between the age of the N_e^R estimate and the ratio N_e^{NC}/N_e^R ($r = -0.228$, $p = 0.132$). This means that the differences between the two estimates are not simply explained by N_e^{NC} reflecting the situation a few tens of generations ago and N_e^R indicating, for different populations, the situation from a few hundred to several thousands of generations ago (see Fig. 3). The ratio N_e^{LT}/N_e^R correlates negatively with the age of the N_e^R estimate ($r = -0.811$, $p = 1.474 \times 10^{-11}$), but the age of the estimate also directly depends on the magnitude of N_e and the correlation may simply indicate that N_e^R and N_e^{LT} disagree more strongly for the small than large populations.

Historical introgression may differently affect alternative N_e estimates

The Baltic Sea nine-spined sticklebacks provide a model system to study the effect of introgression on N_e . We previously showed widespread genetic introgression within the area (Feng et al., 2022), this introgression process possibly continuing among the marine populations but having ceased in the isolated pond populations. The age of the introgression event is unknown but it must have happened after the Eastern lineage colonised the Fennoscandia (Feng et al., 2022).

The trajectories of historical N_e should reflect the past introgression events. The Baltic Sea marine populations show an increase of N_e around 10,000 years ago (Fig. 3) and similar trajectories are seen in coastal freshwater populations known to contain moderate levels of admixture. However, the Swedish freshwater populations with low levels of admixture (3-5%) in a previous study (Feng et al., 2022) show recent MSMC2 trajectories similar to the consensus pattern of the pond ecotype. The area around Umeå, Sweden, is known to have been isolated from the Baltic Sea around 10,000 years ago (Mobley et al., 2011), suggesting a negligible impact of low degree of admixture on historical N_e estimates.

As the introgression event is relatively old, the admixture proportion and the N_e^R show no correlation among the marine population (Table 1). However, the admixture proportion of the marine populations correlates negatively and positively with N_e^{NC} and N_e^{LT} , respectively (Table 1). In the admixed freshwater populations, the admixture proportions are low (median $\alpha = 0.029$) and, with the possible exception of N_e^R , no correlation with N_e estimates is seen (Table 1).

Discussion

The aim of this study was to test whether the alternative methods for the estimation of N_e based on different signals within genomic data are consistent and allow for inferring the full continuum of each population's history. Our results revealed that all studied populations had experienced glacial contractions and, while traces of post-glacial expansions were seen in the marine populations, the freshwater populations have continued to decrease in size. Both historical and near-contemporary N_e were the largest in the marine populations, followed by the coastal freshwaters and other freshwater populations. However, the three alternative N_e estimates showed inconsistencies (Fig. 4), likely reflecting fundamental differences in the data features used by the different methods. Our analyses indicated that the near-contemporary and long-term N_e estimates for the marine populations, obtained with GONE and from genetic diversity, respectively, are affected by genetic introgression. Methods modelling N_e across time are not immune to introgression either, but they can place the effects of introgression on the appropriate time period and, in our analyses, showed no noticeable bias in the most recent time segments. We discuss these issues in more detail in the following.

Summary statistics can reveal hidden within-population aberrations

The effective population size N_e has many definitions (Husemann et al., 2016; Waples, 2022) but they all basically return to the Wright-Fisher model (Fisher, 1931; Wright, 1931) and aim to determine the size of an idealised population behaving genetically in a similar fashion to the target population. This idealised population is closed and panmictic, and thus, many approaches to infer demographic history make the same assumptions (reviewed in Schraiber & Akey, 2015; Nadachowska-Brzyska et al., 2022). Population subdivision, gene flow and overlapping generations are common in natural populations (Patton et al., 2019) but violate the assumptions of the model and can impact the inferences (reviewed in Loog, 2020; Marchi et al., 2021). In non-model species, limited sample sizes and uncertainty of the true conditions of the study population, e.g. the connectivity of the habitat, is often a challenge.

We found that simple summary statistics can greatly help to elucidate the demographic history, explain the observed patterns and prevent misinterpretations. In our study, the Riikojärvi lake population (FIN-RII) from Northern Finland is one of the populations standing clearly out and shows the highest level of linkage disequilibrium (LD), increased relatedness and the most variable N_e estimates in GONE analyses with different recombination fraction parameters. These patterns are not compatible with the population's nucleotide diversity and the level of LD in other populations from similar localities. We hypothesise that the high level of LD in FIN-RII results from recent hybridization between the native fish and genetically distinct migrants. Given the isolation of the location, natural migration seems unlikely and the unrelated fish may have hitchhiked along with the commercially valuable fish stocked in the lake.

Estimates of near-contemporary N_e are affected by recent migration

It is important to note that the N_e estimates in this study should be considered as indicators of the relative demographic changes in each population, and, as such, highlight the impacts of different demographic processes (e.g., population bottlenecks and inbreeding) and climate changes. As the alternative estimates are based on different definitions of N_e and utilise different features of the data, they probably should not be expected to be fully congruent and none of them should be considered as the absolute estimate of N_e .

Our analysis with GONE suggests that the model is generally robust across very different-sized populations. Reassuringly, the estimates for marine populations are roughly aligned, suggesting that the estimation of near-contemporary demographic histories, previously limited to small populations (DeFaveri & Merilä, 2015; Marandel et al., 2019; Nadachowska-Brzyska et al., 2022), may be applied to them as well. However, we also observed deviations and highly unstable estimates for some populations, particularly for those whose connectivity with neighbouring populations have fluctuated and that may have received migration, as well as for populations showing substructure. Such violations of the

expected closed population history are known to lead to biased estimates of recent N_e (1–50 generations before present, GBP), showing both sharply decreased and exponentially increased N_e estimates (Santiago et al., 2020). Limiting the recombination between SNPs can mitigate some of these impacts (Santiago et al., 2020 and Fig. S3), but this may also cause distorted LD patterns and loss of informative SNPs in the analysis, leading to inaccurate estimates (Santiago et al., 2020).

Coalescent-based inferences are not immune to historical gene flow

In contrast to many other methods for the estimation of N_e , the coalescent-based approaches inherently deal with the gene flow and population substructure. Individuals sampled from two different populations have a common ancestor – and coalesce – in the ancestral population from which the two present-day populations are derived. If the common ancestor is (e.g. due to a period of isolation) very far in the past, a panmictic Wright-Fisher population required to give an equally deep coalescence time may have to be very large in size, possibly orders of magnitude larger than the two present-day populations combined. As a result, coalescent-based estimates of N_e derived from the genetic diversity can be both theoretically accurate and completely useless for making practical N_e -based decisions. Exceptions to this are the sequentially Markovian coalescent-based approaches that model the coalescent rate across a segmented time in the past. The coalescent rate – reflecting the inverse of N_e – can vary and often provides meaningful estimates for the recent past despite phases of gene flow and population substructure deeper in the history. On the other hand, all deviations from the panmictic Wright-Fisher model have to be explained as changes in N_e and many historical bumps in N_e trajectories are in fact caused by ancestral population substructure.

Alternative demographic events all affecting the N_e estimates make the interpretation of the results challenging. We found an increase in N_e after the Last Glacial Period in all Baltic Sea marine populations and most of the coastal freshwater populations. This indeed may reflect the colonisation of new areas and population growth in improved climatic conditions, but it could as well result from genetic introgression between two evolutionary lineages (Feng et al., 2022). Simulations conducted by Hawks (2017) showed that even small amounts of introgression from a genetically divergent lineage can affect the inferred historical N_e and leave "waves" in the N_e trajectories of pairwise sequentially Markovian coalescent (PSMC). However, by removing introgressed variants, van der Valk et al., (2020) found no evidence that historical N_e trajectories were affected by introgression. Similarly, Beichman et al. (2017) concluded that MSMC is robust to low levels of admixture, such as the 1–4% of Neanderthal ancestry observed in modern humans (Green et al., 2010). However, given the higher amount of alien ancestry (13–30%, Feng et al., 2022) observed in the Baltic Sea marine populations and the putative post-glacial population expansion, it is difficult to untangle their impacts on the inferences and contributions to the increase in N_e over time. It is likely that both processes have contributed.

Impact of historical gene flow on alternative N_e estimates

Estimating the contemporary N_e for large populations remains a challenging task (DeFaveri & Merilä, 2015; J. Wang et al., 2016; Marandel et al., 2019; Nadachowska-Brzyska et al., 2022). In this study, we explored and compared different ways to estimate N_e across populations of highly different sizes. While none of N_e estimates is immune to the impact of population structure and gene flow, they can still provide valuable insights into the historical, recent and near-contemporary status of the studied populations. For example, the Japanese marine population sampled from a tide pool connected to Biwase Bay is known to have undergone recent inter-species introgression (Yamasaki et al., 2020). In our analyses, this population showed the highest N_e^{LT} and median level N_e^R , but it had a lower N_e^{NC} than the populations from the Baltic and White Sea. Despite their apparent relative incongruence, the results are consistent with the highly complex history of the population: N_e^{LT} reflects the mixed species ancestry, N_e^R the marine origin, and N_e^{NC} the very recent past in a shallow-water tide pool. Although anecdotal, the case well represents the differences behind the alternative estimates of N_e , their timescales and differences in sensitivity to demographic changes.

The three alternative N_e estimates greatly differ in the timescale they represent. The timing of N_e^{NC} is uniform for all populations (on average 32.5 GBP). However, the age of N_e^R (here, the third time segment estimated by MSMC2) varies across and within ecotypes, depending on the genetic diversity of the population. The local estimate of N_e is the inverse of the coalescent rate, and in small populations the rate is higher, providing more signal for the estimation. For the marine populations, the age of N_e^R is on average 1104.3 GBP (± 226.8), while it is 337.3 GBP (± 214.3) for the pond populations. On the other hand, the range of timescales for N_e^{LT} is even wider, as none of the populations has a common ancestor for all genetic variation in their current locality. If they would have, all the coalescent events would have happened and the MSMC2 trajectory would end as no information is available to see any deeper into the past. Given the very different time scales of the alternative methods, it is no surprise to observe variation and incongruencies between methods among different populations and ecotypes.

In addition to the different time scales, the underlying models of the three alternative N_e estimates are differently affected by historical introgression, demonstrated here by the incongruent correlations with the populations' admixture proportions. It is known that admixture can increase LD, leading to an underestimation of N_e by GONE (Saura et al., 2021). In contrast to that, introduced variants increase the genetic variation in the recipient population, thus leading to an increased N_e^{LT} . Interestingly, the N_e^R of the marine populations and all but the N_e^R for the pond populations were not correlated with admixture proportions. A plausible explanation is that the effects on N_e from the habitat itself are much larger than historical admixture in the pond populations, as the ponds are known to have different surface areas (Kivikoski et al., 2023) and different population histories across the areas, and the admixture proportions are in general low (median $\alpha = 0.029$). On the other hand, the ages of N_e^R for the marine populations are much younger than the recent secondary contacts (Fig. 3), and the admixture can correctly be taken into account in the more distant time segments.

We found that the N_e^{NC} and N_e^R are largely incongruent and the disagreement cannot be explained by the age of the N_e^R estimates. As the two estimates differed by two orders of magnitude for some populations, other populations having ratio N_e^{NC}/N_e^R close to one seems coincidental. N_e^{NC} can be underestimated because of admixture-induced high LD as well as recent fluctuations in demography (Waples, 2005) and subpopulation structure (Santiago et al., 2020). Confusingly, specific violations of the expected closed panmictic population, such as subpopulation structure created by gene flow, have opposite effects on N_e^{NC} and N_e^{LT} : sampling of a population mistakenly conducted across two populations would increase LD and strongly decrease N_e^{NC} , whereas the same error would increase the genetic variation and N_e^{LT} . Unless a population and its habitat had remained unchanged for a very long time, it seems impossible for the two estimates to truly agree with each other. On the other hand, comparing N_e^{NC} , N_e^R and N_e^{LT} is like comparing apples and oranges: sometimes such a comparison makes sense, often not. There are many definitions of N_e and they look and behave very differently (Frankham, 1995; Charlesworth, 2009; Husemann et al., 2016; Waples, 2022; Nadachowska-Brzyska et al., 2022).

Demographic history of nine-spined sticklebacks in Northern Europe

Pleistocene climatic oscillations have substantially affected the dynamics of numerous organisms. Previous studies have shown that glaciations and deglaciations are associated with population bottlenecks and expansions, and these have then shaped the genetic diversity and adaptive potential of contemporary populations (reviewed in Hewitt, 2004). Similarly to other species distributed across Europe (e.g. Backström et al., 2013; Liu et al., 2016), our analyses showed the European nine-spined sticklebacks to have undergone population contractions during the glacial period. After the Last Glacial Period (~11,000 years ago), the marine and freshwater ecotypes started to diverge, and the differing timings of the bottlenecks in various freshwater populations reflect their colonisation and the formation of freshwater habitat after the retreating glacial ice sheets. For instance, the pond populations from Umeå, Sweden, showed bottlenecks around 5,000-10,000 years ago, in line with the dating of the historical coastal line in the sampling area (10,000 years ago, Mobley et al. 2011). The White Sea pond populations are known to result from recent colonisations (Ziuganov & Zotin, 1995), and our results placed the colonisation bottleneck to 500 years ago (Figs. S3 and S4). On the other hand, the two river populations from the Baltic Sea coast (from Estonia and Latvia) are known to have similar proportions of Western Lineage ancestry as the nearby marine populations (Feng et al., 2022), and thus must have been established after the secondary contact between the two lineages.

Such fine details highlight the fundamental differences among the studied freshwater populations, particularly their ages and their origin from ancestral populations representing different pools of standing genetic variation. The latter is a major genetic resource for local adaptation (Barrett & Schluter, 2008), and the differences in access to the pool of standing genetic variation due to gene flow or historical demography can either constrain or facilitate local adaptation (e.g., Fang et al., 2021; Kempainen et al., 2021). In this respect, our results lay the groundwork for a deeper understanding of the role of ancestral polymorphism in the local adaptation of the nine-spined sticklebacks.

Conclusions

We compared alternative approaches to quantify N_e using whole-genome data from 45 populations of nine-spined sticklebacks. Our analyses showed that the coalescence-based estimates of N_e in the recent past and the LD-based estimates of near-contemporary N_e are incongruent. The relative sizes of coalescence-based estimates of long-term N_e and of recent past for different populations were largely consistent but the estimates for a population could differ by an order of magnitude in size, the difference explained by the alternative approaches for incorporating the signals of past demographic events in the estimate of N_e . As a result, the sequentially Markovian coalescent-based methods seem to provide a meaningful interpretation of the real-life N_e , though they lose resolution for the very recent past for populations with large N_e . The LD-based estimates reflect a very different definition of N_e and could differ from the coalescent-based estimates by two orders of magnitude. The LD-based estimates can also be highly sensitive to violations of the expected closed population history and, in our analyses, the estimates differed by several orders of magnitude for seemingly comparable populations within the same ecotype. Nevertheless, our results suggest that when applied cautiously, the LD-based GONE can provide reasonable inferences of the very recent demographic history for panmictic marine populations for which the N_e estimation has traditionally been challenging.

Authors' contributions

J.M. started the project. X.F., J.M. and A.L. devised the research idea. X.F. performed the analyses with the help of A.L. X.F. and A.L. wrote the first draft and all authors participated in the writing of the final manuscript.

Acknowledgements

We thank Miinastiina Issakainen, Sami Karja, Laura Häkkinen and Kirsi Kähkönen for help in the laboratory; Kirsi Vestenius, Ari Savikko and Kare Koivisto for help with the information of fish transplantation in northern Finland; and numerous collaborators and colleagues for their help in obtaining the samples (listed in Acknowledgements of Feng et al., 2022). The advice and support from Paolo Momigliano, Petri Kempainen, and Mikko Kivikoski is gratefully acknowledged. Our research was supported by grants from the Academy Finland (129662, 134728 and 218343 to JM; 322681 to AL), China Scholarship Council (#201608520032 to XF), and Finnish Cultural Foundation (#00210295 to XF). Computational resources provided by the CSC-IT Center for Science, Finland, are acknowledged with gratitude.

Conflict of interest

The authors declare no competing interests.

Data availability statement

The whole-genome re-sequencing data have been published previously by Feng et al. (2022) and all the raw sequence data for this study can be accessed through European Nucleotide Archive (ENA) (<https://www.ebi.ac.uk/ena>) under accession code PRJEB39599. Other relevant data (e.g., filtered VCF files, input and output files) are available from the Zenodo Open Repository: <https://zenodo.org/record/xxxxxx>.

References

- Allendorf, F. W., Luikart, G. H., & Aitken, S. N. (2012). *Conservation and the genetics of populations*. John Wiley & Sons.
- Backström, N., Sætre, G.-P., & Ellegren, H. (2013). Inferring the demographic history of European Ficedula flycatcher populations. *BMC Evolutionary Biology*, 13(1), 2. <https://doi.org/10.1186/1471-2148-13-2>
- Barrett, R. D. H., & Schluter, D. (2008). Adaptation from standing genetic variation. *Trends in Ecology & Evolution*, 23(1), 38–44. <https://doi.org/10.1016/j.tree.2007.09.008>
- Beichman, A. C., Huerta-Sanchez, E., & Lohmueller, K. E. (2018). Using Genomic Data to Infer Historic Population Dynamics of Nonmodel Organisms. *Annual Review of Ecology, Evolution, and Systematics*, 49(1), 433–456. <https://doi.org/10.1146/annurev-ecolsys-110617-062431>
- Beichman, A. C., Phung, T. N., & Lohmueller, K. E. (2017). Comparison of Single Genome and Allele Frequency Data Reveals Discordant Demographic Histories. *G3 Genes|Genomes|Genetics*, 7(11), 3605–3620. <https://doi.org/10.1534/g3.117.300259>
- Bruford M W, Hanotte O, Brookfield J F Y, et al. (1998). Multilocus and single-locus DNA fingerprinting. *Molecular genetic analysis of populations: a practical approach*. 2: 287-336.
- Charlesworth, B. (2009). Effective population size and patterns of molecular evolution and variation. *Nature Reviews Genetics*, 10(3), Article 3. <https://doi.org/10.1038/nrg2526>
- Charlesworth, B., & Charlesworth, D. (2010). *Elements of Evolutionary Genetics*. Roberts

and Company Publishers.

- Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin, R., & 1000 Genomes Project Analysis Group. (2011). The variant call format and VCFtools. *Bioinformatics*, *27*(15), 2156–2158. <https://doi.org/10.1093/bioinformatics/btr330>
- Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., Whitwham, A., Keane, T., McCarthy, S. A., Davies, R. M., & Li, H. (2021). Twelve years of SAMtools and BCFtools. *GigaScience*, *10*(2), giab008. <https://doi.org/10.1093/gigascience/giab008>
- DeFaveri, J., & Merilä, J. (2015). Temporal Stability of Genetic Variability and Differentiation in the Three-Spined Stickleback (*Gasterosteus aculeatus*). *PLOS ONE*, *10*(4), e0123891. <https://doi.org/10.1371/journal.pone.0123891>
- DeFaveri, J., Shikano, T., & Merilä, J. (2014). Geographic Variation in Age Structure and Longevity in the Nine-Spined Stickleback (*Pungitius pungitius*). *PLOS ONE*, *9*(7), e102660. <https://doi.org/10.1371/journal.pone.0102660>
- DeWoody, J. A., & Avise, J. C. (2000). Microsatellite variation in marine, freshwater and anadromous fishes compared with other animals. *Journal of Fish Biology*, *56*(3), 461–473. <https://doi.org/10.1111/j.1095-8649.2000.tb00748.x>
- Fang, B., Kemppainen, P., Momigliano, P., & Merilä, J. (2021). Population Structure Limits Parallel Evolution in Sticklebacks. *Molecular Biology and Evolution*, *38*(10), 4205–4221. <https://doi.org/10.1093/molbev/msab144>
- Feng, X., Merilä, J., & Löytynoja, A. (2022). Complex population history affects admixture analyses in nine-spined sticklebacks. *Molecular Ecology*, *31*(20), 5386–5401. <https://doi.org/10.1111/mec.16651>
- Fisher, R. A. (1931). XVII.—The Distribution of Gene Ratios for Rare Mutations. *Proceedings of the Royal Society of Edinburgh*, *50*, 204–219. <https://doi.org/10.1017/S0370164600044886>
- Flint-Garcia, S. A., Thornsberry, J. M., & Buckler, E. S. (2003). Structure of Linkage

Disequilibrium in Plants. *Annual Review of Plant Biology*, 54(1), 357–374.

<https://doi.org/10.1146/annurev.arplant.54.031902.134907>

Frankham, R. (1995). Effective population size/adult population size ratios in wildlife: A review. *Genetics Research*, 66(2), 95–107.

<https://doi.org/10.1017/S0016672300034455>

Frankham, R., Ballou, J. D., & Briscoe, D. A. (2010). *Introduction to Conservation Genetics*. Cambridge University Press.

Gilbert, K. J., & Whitlock, M. C. (2015). Evaluating methods for estimating local effective population size with and without migration. *Evolution*, 69(8), 2154–2166.

<https://doi.org/10.1111/evo.12713>

Green, R. E., Krause, J., Briggs, A. W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li, H., Zhai, W., Fritz, M. H.-Y., Hansen, N. F., Durand, E. Y., Malaspinas, A.-S., Jensen, J. D., Marques-Bonet, T., Alkan, C., Prüfer, K., Meyer, M., Burbano, H. A., ... Pääbo, S. (2010). A Draft Sequence of the Neandertal Genome. *Science*, 328(5979),

710–722. <https://doi.org/10.1126/science.1188021>

Guo, B., Fang, B., Shikano, T., Momigliano, P., Wang, C., Kravchenko, A., & Merilä, J. (2019). A phylogenomic perspective on diversity, hybridization and evolutionary affinities in the stickleback genus *Pungitius*. *Molecular Ecology*, 28(17), 4046–4064.

<https://doi.org/10.1111/mec.15204>

Hanghøj, K., Moltke, I., Andersen, P. A., Manica, A., & Korneliussen, T. S. (2019). Fast and accurate relatedness estimation from high-throughput sequencing data in the presence of inbreeding. *GigaScience*, 8(5), giz034.

<https://doi.org/10.1093/gigascience/giz034>

Hare, M. P., Nunney, L., Schwartz, M. K., Ruzzante, D. E., Burford, M., Waples, R. S., Ruegg, K., & Palstra, F. (2011). Understanding and Estimating Effective Population Size for Practical Application in Marine Species Management. *Conservation Biology*,

25(3), 438–449. <https://doi.org/10.1111/j.1523-1739.2010.01637.x>

Hawks, J. (2017). Introgression Makes Waves in Inferred Histories of Effective Population

- Size. *Human Biology*, 89(1), 67–80.
- Hedrick, P. W., & Lacy, R. C. (2015). Measuring Relatedness between Inbred Individuals. *Journal of Heredity*, 106(1), 20–25. <https://doi.org/10.1093/jhered/esu072>
- Hewitt, G. M. (2004). Genetic consequences of climatic oscillations in the Quaternary. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 359(1442), 183–195.
- Husemann, M., Zachos, F. E., Paxton, R. J., & Habel, J. C. (2016). Effective population size in ecology and evolution. *Heredity*, 117(4), 191–192. <https://doi.org/10.1038/hdy.2016.75>
- Kemppainen, P., Li, Z., Rastas, P., Löytynoja, A., Fang, B., Yang, J., Guo, B., Shikano, T., & Merilä, J. (2021). Genetic population structure constrains local adaptation in sticklebacks. *Molecular Ecology*, 30(9), 1946–1961. <https://doi.org/10.1111/mec.15808>
- Kimura, M. (1983). *The neutral theory of molecular evolution*. Cambridge University Press.
- Kivikoski, M., Feng, X., Löytynoja, A., Momigliano, P., & Merilä, J. (2023). Determinants of genetic diversity in sticklebacks (p. 2023.03.17.533073). *bioRxiv*. <https://doi.org/10.1101/2023.03.17.533073>
- Kivikoski, M., Rastas, P., Löytynoja, A., & Merilä, J. (2021). Automated improvement of stickleback reference genome assemblies with Lep-Anchor software. *Molecular Ecology Resources*, 21(6), 2166–2176. <https://doi.org/10.1111/1755-0998.13404>
- Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *ArXiv Preprint ArXiv:1303.3997*.
- Li, H., & Durbin, R. (2011). Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357), Article 7357. <https://doi.org/10.1038/nature10231>
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., & 1000 Genome Project Data Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079.

<https://doi.org/10.1093/bioinformatics/btp352>

Liu, S., Hansen, M. M., & Jacobsen, M. W. (2016). Region-wide and ecotype-specific differences in demographic histories of threespine stickleback populations, estimated from whole genome sequences. *Molecular Ecology*, *25*(20), 5187–5202.

<https://doi.org/10.1111/mec.13827>

Loog, L. (2020). Sometimes hidden but always there: The assumptions underlying genetic inference of demographic histories. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *376*(1816), 20190719. <https://doi.org/10.1098/rstb.2019.0719>

Luikart, G., Ryman, N., Tallmon, D. A., Schwartz, M. K., & Allendorf, F. W. (2010). Estimation of census and effective population sizes: The increasing usefulness of DNA-based approaches. *Conservation Genetics*, *11*(2), 355–373.

<https://doi.org/10.1007/s10592-010-0050-7>

Malaspinas, A.-S., Westaway, M. C., Muller, C., Sousa, V. C., Lao, O., Alves, I., Bergström, A., Athanasiadis, G., Cheng, J. Y., Crawford, J. E., Heupink, T. H., Macholdt, E., Peischl, S., Rasmussen, S., Schiffels, S., Subramanian, S., Wright, J. L., Albrechtsen, A., Barbieri, C., ... Willerslev, E. (2016). A genomic history of Aboriginal Australia. *Nature*, *538*(7624), Article 7624. <https://doi.org/10.1038/nature18299>

Marandel, F., Lorance, P., Berthelé, O., Trenkel, V. M., Waples, R. S., & Lamy, J.-B. (2019). Estimating effective population size of large marine populations, is it feasible? *Fish and Fisheries*, *20*(1), 189–198. <https://doi.org/10.1111/faf.12338>

Marchi, N., Schlichta, F., & Excoffier, L. (2021). Demographic inference. *Current Biology*, *31*(6), R276–R279. <https://doi.org/10.1016/j.cub.2021.01.053>

McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*, *20*(9), 1297–1303.

<https://doi.org/10.1101/gr.107524.110>

Miller, S. (1988). A simple salting-out procedure tissue for extracting DNA from human

nucleated cells. *Nucleic Acids Res.* 16: 221.

Mobley, K. B., Lussetti, D., Johansson, F., Englund, G., & Bokma, F. (2011). Morphological and genetic divergence in Swedish postglacial stickleback (*Pungitius pungitius*) populations. *BMC Evolutionary Biology*, 11(1), 287.

<https://doi.org/10.1186/1471-2148-11-287>

Nadachowska-Brzyska, K., Konczal, M., & Babik, W. (2022). Navigating the temporal continuum of effective population size. *Methods in Ecology and Evolution*, 13(1), 22–41. <https://doi.org/10.1111/2041-210X.13740>

Palstra, F. P., & Ruzzante, D. E. (2008). Genetic estimates of contemporary effective population size: What can they tell us about the importance of genetic stochasticity for wild population persistence? *Molecular Ecology*, 17(15), 3428–3447.

<https://doi.org/10.1111/j.1365-294X.2008.03842.x>

Patton, A. H., Margres, M. J., Stahlke, A. R., Hendricks, S., Lewallen, K., Hamede, R. K., Ruiz-Aravena, M., Ryder, O., McCallum, H. I., Jones, M. E., Hohenlohe, P. A., & Storfer, A. (2019). Contemporary Demographic Reconstruction Methods Are Robust to Genome Assembly Quality: A Case Study in Tasmanian Devils. *Molecular Biology and Evolution*, 36(12), 2906–2921. <https://doi.org/10.1093/molbev/msz191>

R Core Team. (2020). *R: A Language and Environment for Statistical Computing*.

Romiguier, J., Gayral, P., Ballenghien, M., Bernard, A., Cahais, V., Chenuil, A., Chiari, Y., Dernaï, R., Duret, L., Faivre, N., Loire, E., Lourenco, J. M., Nabholz, B., Roux, C., Tsagkogeorga, G., Weber, A. a.-T., Weinert, L. A., Belkhir, K., Bierne, N., ... Galtier, N. (2014). Comparative population genomics in animals uncovers the determinants of genetic diversity. *Nature*, 515(7526), Article 7526.

<https://doi.org/10.1038/nature13685>

Sambrook, J., & Russell, D. W. (2006). *The condensed protocols from molecular cloning: A laboratory manual*. Cold Spring Harbor Laboratory Press.

Santiago, E., Novo, I., Pardiñas, A. F., Saura, M., Wang, J., & Caballero, A. (2020). Recent Demographic History Inferred by High-Resolution Analysis of Linkage Disequilibrium.

Molecular Biology and Evolution, 37(12), 3642–3653.

<https://doi.org/10.1093/molbev/msaa169>

Saura, M., Caballero, A., Santiago, E., Fernández, A., Morales-González, E., Fernández, J.,

Cabaleiro, S., Millán, A., Martínez, P., Palaikostas, C., Kocour, M., Aslam, M. L.,

Houston, R. D., Prchal, M., Bargelloni, L., Tzokas, K., Haffray, P., Bruant, J.-S., &

Villanueva, B. (2021). Estimates of recent and historical effective population size in

turbot, seabream, seabass and carp selective breeding programmes. *Genetics*

Selection Evolution, 53(1), 85. <https://doi.org/10.1186/s12711-021-00680-9>

Schiffels, S., & Durbin, R. (2014). Inferring human population size and separation history

from multiple genome sequences. *Nature Genetics*, 46(8), Article 8.

<https://doi.org/10.1038/ng.3015>

Schiffels, S., & Wang, K. (2020). MSMC and MSMC2: The multiple sequentially markovian

coalescent. In *Statistical population genomics* (pp. 147–165). Humana.

Schraiber, J. G., & Akey, J. M. (2015). Methods and models for unravelling human

evolutionary history. *Nature Reviews Genetics*, 16(12), Article 12.

<https://doi.org/10.1038/nrg4005>

Sellinger, T. P. P., Abu-Awad, D., & Tellier, A. (2021). Limits and convergence properties of

the sequentially Markovian coalescent. *Molecular Ecology Resources*, 21(7),

2231–2248. <https://doi.org/10.1111/1755-0998.13416>

van der Valk, T., Gonda, C. M., Silegowa, H., Almanza, S., Sifuentes-Romero, I., Hart, T. B.,

Hart, J. A., Detwiler, K. M., & Guschanski, K. (2020). The Genome of the Endangered

Dryas Monkey Provides New Insights into the Evolutionary History of the Vervets.

Molecular Biology and Evolution, 37(1), 183–194.

<https://doi.org/10.1093/molbev/msz213>

Varadharajan, S., Rastas, P., Löytynoja, A., Matschiner, M., Calboli, F. C. F., Guo, B.,

Nederbragt, A. J., Jakobsen, K. S., & Merilä, J. (2019). A High-Quality Assembly of

the Nine-Spined Stickleback (*Pungitius pungitius*) Genome. *Genome Biology and*

Evolution, 11(11), 3291–3308. <https://doi.org/10.1093/gbe/evz240>

- Walsh, C. A. J., Momigliano, P., Boussarie, G., Robbins, W. D., Bonnin, L., Fauvelot, C., Kiszka, J. J., Mouillot, D., Vigliola, L., & Manel, S. (2022). Genomic insights into the historical and contemporary demographics of the grey reef shark. *Heredity*, *128*(4), Article 4. <https://doi.org/10.1038/s41437-022-00514-4>
- Wang, J. (2005). Estimation of effective population sizes from data on genetic markers. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *360*(1459), 1395–1409. <https://doi.org/10.1098/rstb.2005.1682>
- Wang, J., Santiago, E., & Caballero, A. (2016). Prediction and estimation of effective population size. *Heredity*, *117*(4), Article 4. <https://doi.org/10.1038/hdy.2016.43>
- Wang, Y., Wang, Y., Cheng, X., Ding, Y., Wang, C., Merilä, J., & Guo, B. (2023). Prevalent Introgression Underlies Convergent Evolution in the Diversification of Pungitius Sticklebacks. *Molecular Biology and Evolution*, *40*(2), msad026. <https://doi.org/10.1093/molbev/msad026>
- Waples, R. K., Larson, W. A., & Waples, R. S. (2016). Estimating contemporary effective population size in non-model species using linkage disequilibrium across thousands of loci. *Heredity*, *117*(4), Article 4. <https://doi.org/10.1038/hdy.2016.60>
- Waples, R. S. (2022). What Is N_e , Anyway? *Journal of Heredity*, *113*(4), 371–379. <https://doi.org/10.1093/jhered/esac023>
- Waples, R. S., & Do, C. (2010). Linkage disequilibrium estimates of contemporary N_e using highly variable genetic markers: A largely untapped resource for applied conservation and evolution. *Evolutionary Applications*, *3*(3), 244–262. <https://doi.org/10.1111/j.1752-4571.2009.00104.x>
- Ward, R. D., Woodwark, M., & Skibinski, D. O. F. (1994). A comparison of genetic diversity levels in marine, freshwater, and anadromous fishes. *Journal of Fish Biology*, *44*(2), 213–232. <https://doi.org/10.1111/j.1095-8649.1994.tb01200.x>
- Wright, S. (1931). Evolution in Mendelian Populations. *Genetics*, *16*(2), 97–159.
- Yamasaki, Y. Y., Kakioka, R., Takahashi, H., Toyoda, A., Nagano, A. J., Machida, Y., Møller, P. R., & Kitano, J. (2020). Genome-wide patterns of divergence and introgression

after secondary contact between *Pungitius* sticklebacks. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 375(1806), 20190548.

<https://doi.org/10.1098/rstb.2019.0548>

Zhang, C., Dong, S.-S., Xu, J.-Y., He, W.-M., & Yang, T.-L. (2019). PopLDdecay: A fast and effective tool for linkage disequilibrium decay analysis based on variant call format files. *Bioinformatics*, 35(10), 1786–1788.

<https://doi.org/10.1093/bioinformatics/bty875>

Zhang, C., Reid, K., Sands, A. F., Fraimout, A., Schierup, M. H., & Merilä, J. (2023). De novo mutation rates in sticklebacks (p. 2023.03.16.532904). *bioRxiv*.

<https://doi.org/10.1101/2023.03.16.532904>

Ziuganov, V. V., & Zotin, A. A. (1995). Pelvic Girdle Polymorphism and Reproductive Barriers in the Ninespine Stickleback *Pungitius Pungitius* (L.) From Northwest Russia.

Behaviour, 132(13–14), 1095–1105. <https://doi.org/10.1163/156853995X00478>

Figures and Table for manuscript:

Estimating recent and historical effective population size of marine and freshwater sticklebacks

Figures:

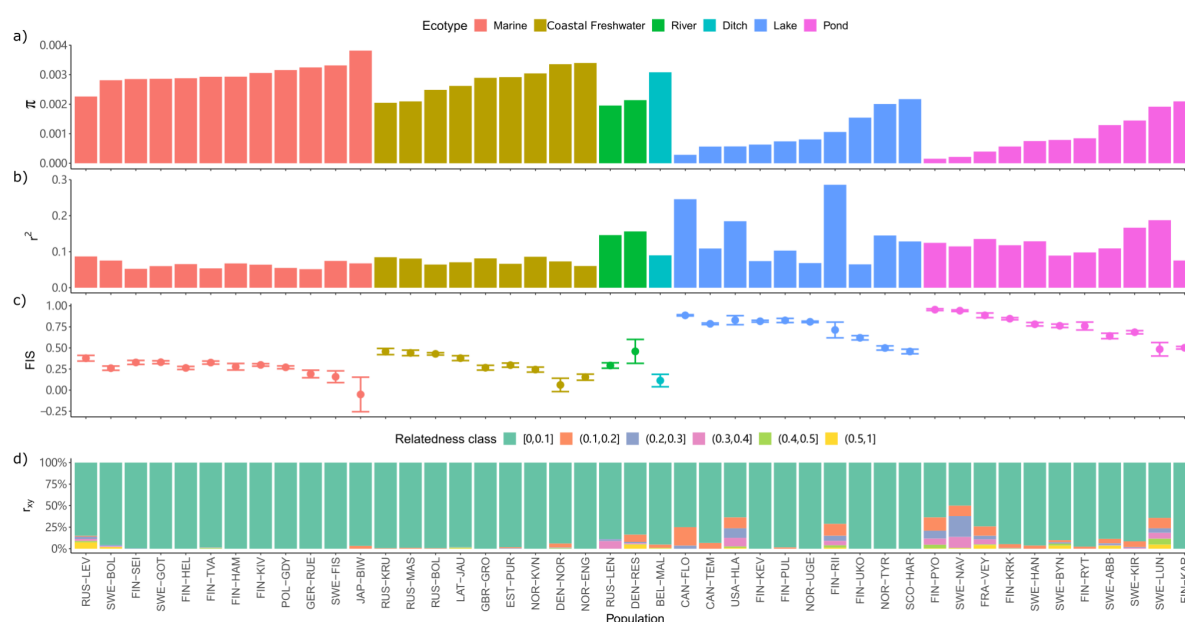


Fig. 1. Summary statistics for the studied nine-spined stickleback populations. a) Nucleotide diversity (π) across the autosomal chromosomes. b) LD calculated as the harmonic mean of r^2 for SNPs located 100-200 kb apart within LG4. c) Inbreeding coefficients (F_{IS}) with standard deviations. d) Relatedness (r_{xy}) for pairs of individuals within populations with colours representing the proportion of pairwise comparisons within a population falling in a specific relatedness class. For the LD decay curve for each population, see Fig. S2.

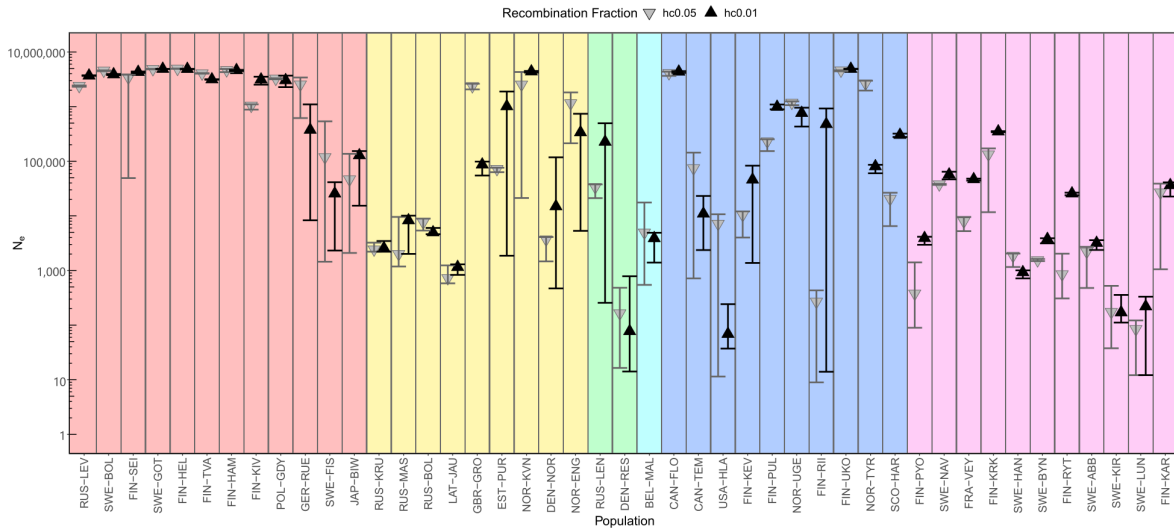


Fig. 2. The near-contemporary effective population size N_e^{NC} estimated with GONE. For each population, the mean and the range of N_e estimates obtained using a recombination fraction (hc) of 0.01 and 0.05 are shown. Populations are ordered as in Fig. 1. Background colours represent the ecotypes: marine in light red, coastal freshwater in light yellow, river in light green, ditch in light cyan, lake in light blue and pond in light pink. For the trajectories of each population, see Fig. S3.

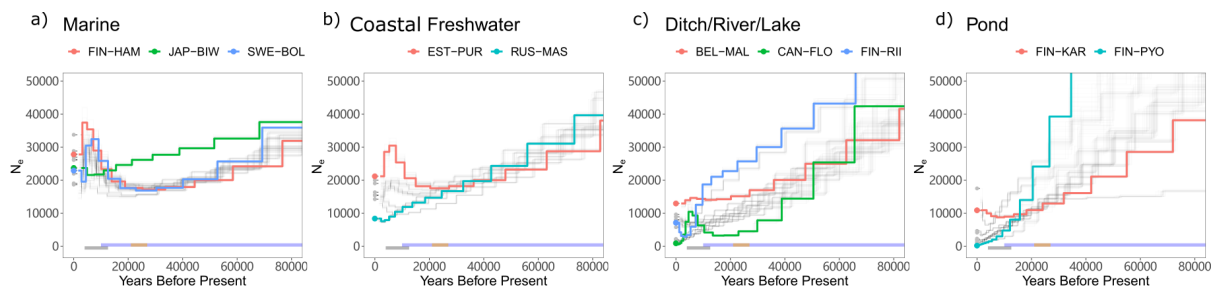


Fig. 3. Historical effective population sizes (N_e) for populations within different ecotypes inferred with MSMC2. a) Marine populations. b) Coastal and freshwater populations. c) Lake and stream populations. d) Pond populations. The x-axis shows the time in years before present based on generation time of two years and mutation rate of 4.37×10^{-9} per base pair per generation. The grey, blue and orange bars at the bottom indicate the times for the formation of the Baltic Sea, the Last Glacial Period and the Last Glacial Maximum, respectively. Within each panel, the thin lines correspond to the original inferences and 20 rounds of bootstrap replicates, and the bold lines show representative populations from each ecotype. The dots at $x=0$ represent the N_e^R estimates and are projected from the last trustworthy estimate (time segment three) for the corresponding population. The end points of the trajectory lines reflect the ages of the N_e^R estimates. For the trajectories of each population, see Fig. S4.

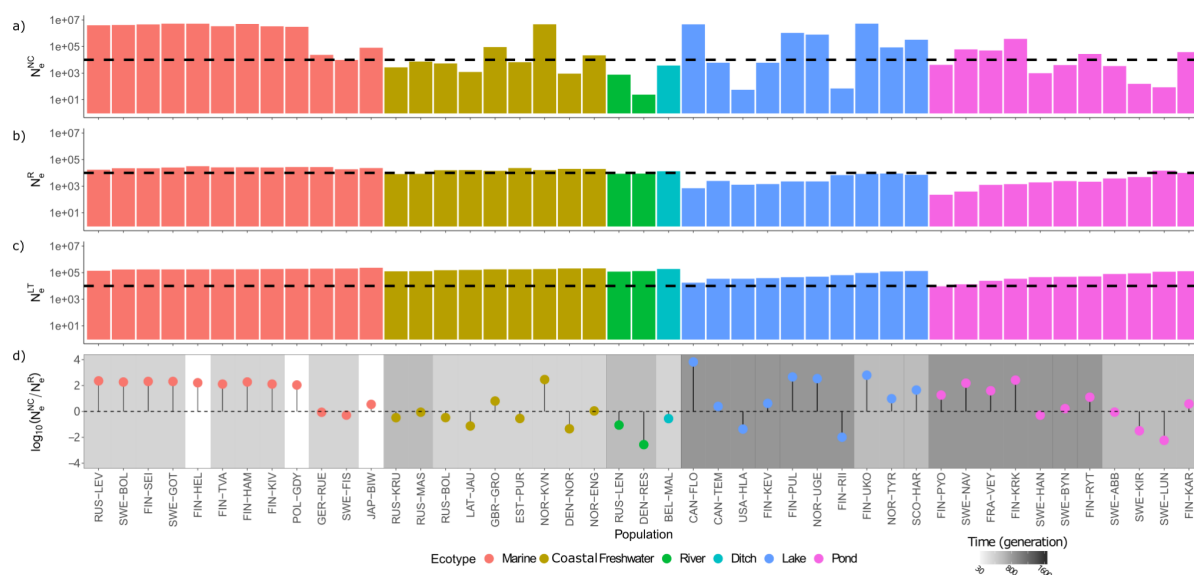


Fig. 4. Alternative N_e estimates for different populations. The colours indicate the ecotype and the dashed horizontal lines shows N_e of 1×10^4 . a) Near-contemporary N_e^{NC} estimated using GONE generations before the present. b) Recent N_e^R estimated with MSMC2. c) Long-term N_e^{LT} computed from genetic diversity as $N = \pi/(4\mu)$. d) The ratio N_e^{NC} / N_e^R with the background shading representing the age of the N_e^R estimate, with the darkest colours indicating an age approximately 100 years ago. Note the log-scale: the dashed line indicates the ratio of 1.

Table:

Table.1 Pearson product moment correlations between the different N_e estimates and the admixture proportions in marine and pond ecotypes. Ten marine and seven pond populations were included.

Method	Ecotype	r	p
N_e^{NC}	Marine	-0.843	0.002
N_e^R	Marine	-0.027	0.940
N_e^{LT}	Marine	0.805	0.005
N_e^{NC}	Pond	-0.269	0.560
N_e^R	Pond	0.770	0.043
N_e^{LT}	Pond	0.411	0.360