

Understanding population structure in an evolutionary context: population-specific F_{ST} and pairwise F_{ST}

Suichi Kitada,^{1,*} Reiichiro Nakamichi,² and Hirohisa Kishino^{3,4}

¹Tokyo University of Marine Science and Technology, Tokyo 108-8477, Japan,

²Japan Fisheries Research and Education Agency, Yokohama 236-8648, Japan,

³Graduate School of Agriculture and Life Sciences, The University of Tokyo, Tokyo 113-8657, Japan, and

⁴The Research Institute of Evolutionary Biology, Tokyo 138-0098, Japan

*Corresponding author: Tokyo University of Marine Science and Technology, 4-5-7 Konan, Minato-ku, Tokyo 108-8477, Japan. Email: kitada@kaiyodai.ac.jp

Abstract

Populations are shaped by their history. It is crucial to interpret population structure in an evolutionary context. Pairwise F_{ST} measures population structure, whereas population-specific F_{ST} measures deviation from the ancestral population. To understand the current population structure and a population's history of range expansion, we propose a representation method that overlays population-specific F_{ST} estimates on a sampling location map, and on an unrooted neighbor-joining tree and a multi-dimensional scaling plot inferred from a pairwise F_{ST} distance matrix. We examined the usefulness of our procedure using simulations that mimicked population colonization from an ancestral population and by analyzing published human, Atlantic cod, and wild poplar data. Our results demonstrated that population-specific F_{ST} values identify the source population and trace the evolutionary history of its derived populations. Conversely, pairwise F_{ST} values represent the current population structure. By integrating the results of both estimators, we obtained a new picture of the population structure that incorporates evolutionary history. The generalized least squares estimate of genome-wide population-specific F_{ST} indicated that the wild poplar population expanded its distribution to the north, where daylight hours are long in summer, to coastal areas with abundant rainfall, and to the south where summers are dry. Genomic data highlight the power of the bias-corrected moment estimators of F_{ST} , whether global, pairwise, or population-specific, that provide unbiased estimates of F_{ST} . All F_{ST} moment estimators described in this paper have reasonable processing times and are useful in population genomics studies.

Keywords: demographic history; genetic diversity; migration; population structure; range expansion

Introduction

Quantifying genetic relationships among populations is of substantial interest in population biology, ecology, and human genetics (Weir and Hill 2002). Appropriate estimates of population structure are the basis of our understanding of biology and biological applications, which vary from evolutionary and conservation studies to association mapping and forensic identification (Weir and Hill 2002). For such objectives, Wright's F_{ST} (Wright 1951) is commonly used to quantify the genetic divergence of populations, and there have been many informative reviews of F_{ST} estimators (e.g., Balloux and Lugon-Moulin 2002; Weir and Hill 2002; Rousset 2004, 2007; Beaumont 2005; Excoffier 2007; Holsinger and Weir 2009; Gaggiotti and Foll 2010; Bhatia et al. 2013). The traditional F_{ST} estimators have been defined as the ratio of the between-population variance to the total variance in allele frequencies (Wright 1965; Cockerham 1969, 1973; Weir and Cockerham 1984; Balloux and Lugon-Moulin 2002; Excoffier 2007; Holsinger and Weir 2009). An alternative approach for estimating population differentiation is to use population-specific F_{ST} estimators (Balding and Nichols 1995; Nicholson et al. 2002; Weir and Hill 2002; Weir et al. 2005; Gaggiotti and Foll 2010; Weir and

Goudet 2017). Model-based Bayesian approaches, based on beta and/or Dirichlet distributions, for estimating population-specific F_{ST} have been proposed (Balding and Nichols 1995; Nicholson et al. 2002; Falush et al. 2003; Beaumont and Balding 2004). In addition to model-based methods, moment estimators of population-specific F_{ST} have been derived (Weir and Hill 2002; Weir and Goudet 2017). A large number of approaches exist for estimating F_{ST} that have different underlying assumptions (global, pairwise, or population-specific F_{ST}) and the framework used, such as frequentist and/or Bayesian. There have been many comprehensive reviews of traditional and population-specific F_{ST} estimators, as indicated above, most of which were written from the viewpoint of theoretical issues. Conversely, there has been no formal comparative study to describe their differences in terms of the evolutionary scenarios that best describe the data. Although these issues are well understood among statistical geneticists and theoretical population geneticists, empirical researchers, particularly those working in non-model organisms, could benefit from studies that address the problem.

In practice, the F_{ST} value estimated from a set of population samples is called the global F_{ST} , which measures population

Received: January 20, 2021. Accepted: August 27, 2021

© The Author(s) 2021. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

differentiation overall populations (e.g., Pérez-Lezaun et al. 1997). Additionally, F_{ST} values between pairs of population samples (pairwise F_{ST} , Reynolds et al. 1983) are routinely used to estimate population structure in human genetics (Pérez-Lezaun et al. 1997; Ramachandran et al. 2005), conservation biology (Palsbøll et al. 2007; Schwartz et al. 2007), and evolutionary biology and ecology (e.g., Hemmer-Hansen et al. 2013a; Therkildsen et al. 2013a; Geraldès et al. 2014; McKown et al. 2014a; Jorde et al. 2015; Rougemont et al. 2020). Divergent selection across an environmental gradient can influence population structure (Nosil et al. 2009; Orsini et al. 2013), and researchers have examined geographic distance and habitat differences between populations as explanatory variables that affect population structure estimated based on genome-wide (average over loci) pairwise F_{ST} values (e.g., Rousset 1997; Bradbury and Bentzen 2007; Petrou et al. 2014; Jorde et al. 2015; Kitada et al. 2017). To identify the adaptive divergence of a gene among populations, locus-population-specific F_{ST} was developed using empirical Bayes (Beaumont and Balding 2004) and full Bayesian methods (BayeScan) (Foll and Gaggiotti 2008). The Bayesian methods have been applied in many cases across various species to identify outlier single-nucleotide polymorphisms (SNPs) (e.g., Limborg et al. 2012; Therkildsen et al. 2013a; Geraldès et al. 2014). However, genome-wide population-specific F_{ST} is new among biologists. Despite the expected usefulness of genome-wide population-specific F_{ST} in evolutionary biology (Weir and Goudet 2017), applications have been sparse to date (e.g., Nicholson et al. 2002; Weir et al. 2005; Foll and Gaggiotti 2006; Buckleton et al. 2016; Rougemont et al. 2020).

Traditional F_{ST} estimators were originally developed to estimate F_{ST} over a metapopulation (global F_{ST}) based on a set of population samples (Cockerham 1969; Nei and Chesser 1983; Weir and Cockerham 1984; Excoffier 2007; Rousset 2007). In this study, we use Nei and Chesser's (1983) bias-corrected G_{ST} moment estimator (hereafter, NC83) as a pairwise F_{ST} estimator (Supplementary Note). The pairwise F_{ST} between populations (i, j) is defined as:

$$pwF_{ST}^{ij} = \frac{H_T^{ij} - H_S^{ij}}{H_T^{ij}} = 1 - \frac{H_S^{ij}}{H_T^{ij}} \quad (1)$$

where H_T^{ij} is total heterozygosity over all populations and H_S^{ij} is within-population heterozygosity.

We apply Weir and Goudet's (2017) bias-corrected moment estimator of population-specific F_{ST} (hereafter, WG) (Supplementary Note). When only allele frequencies are used, the WG population-specific F_{ST} at a locus is defined as:

$$psF_{ST}^i = \frac{M_W^i - M^B}{1 - M^B}$$

where M_W^i is the within-population matching of two distinct alleles of population i and M^B is the between-population-pair matching average over pairs of populations i, i' . M^B is homozygosity over pairs of populations. We represent heterozygosity over all pairs of populations as $1 - M^B = H_B$, and $1 - M_W^i = H_{Si}$. Therefore:

$$psF_{ST}^i = \frac{H_B - H_{Si}}{H_B} = 1 - \frac{H_{Si}}{H_B} \quad (2)$$

This formulation is reasonable because WG population-specific F_{ST} uses "allele matching, equivalent to homozygosity

and complementary to heterozygosity as used by Nei (1973), rather than components of variance" (Weir and Goudet 2017). H_B is heterozygosity for all pairs of populations, whereas H_T^{ij} in Equation (1) is heterozygosity for the pair of populations. Equation (2) shows that WG population-specific F_{ST} measures population-specific genetic diversity (H_{Si}) under the framework of the relatedness of individuals and identifies the population with the greatest genetic diversity as the ancestral or oldest population. Because populations close to the ancestral population have had more opportunities for mutations than recently founded populations (Liu et al. 2006), they are likely to have the highest heterozygosity and low values of population-specific F_{ST} . Thus, WG population-specific F_{ST} works to infer evolutionary history through genetic diversity in terms of heterozygosity under the assumption that the ancestral population had the highest genetic diversity. By combining population-specific and pairwise F_{ST} estimates, we can infer the present population structure, which reflects evolutionary history.

In this study, our objective is to demonstrate to empirical population geneticists and biologists how the two types of genome-wide F_{ST} estimators can be combined to help elucidate the population structure (pairwise F_{ST}) in the evolutionary context (population-specific F_{ST}). In our approach, the current population structure is represented by an unrooted neighbor-joining (NJ) tree (Saitou and Nei 1987) and a multi-dimensional scaling (MDS) plot based on pairwise F_{ST} values, with population history being inferred by overlaying population-specific F_{ST} values on the population structure. The colors of the populations (names and/or sampling points) based on the WG genome-wide population-specific F_{ST} values enable the inference of the historical order of population colonization. We also present a representation on a geographical map, which is useful for visually understanding population history in a distribution range.

First, we examine the usefulness of our procedure using stepping-stone simulations that mimic population colonization from a single ancestral population for five scenarios of population range expansion. We then apply our approach to three datasets of human, Atlantic cod (*Gadus morhua*), and wild poplar (*Populus trichocarpa*). Human evolutionary history, migration, and population structure have been particularly well studied (e.g., Diamond 1997; Rosenberg et al. 2002; Ramachandran et al. 2005; Liu et al. 2006; Pickrell and Pritchard 2012; Lipson et al. 2013; Hellenthal et al. 2014; Rutherford 2016; Nielsen et al. 2017). These patterns are well known by statistical/theoretical population geneticists and biologists; therefore, testing our integrative F_{ST} analysis on this dataset could provide a good example of the usefulness of this practical approach. Although dense human SNP datasets are currently available, we used microsatellite data for illustrative purposes, because the quality of human microsatellites has been examined thoroughly, and they are fundamentally neutral (Kanitz et al. 2018). Also, microsatellites are highly polymorphic and have more information at a locus than SNPs (Schlötterer 2004).

The Atlantic cod SNP were genotyped from mature fish samples collected from the North Atlantic from the northern range margin of the species in Greenland, Norway, and the Baltic Sea. Two ecotypes (migratory and stationary) that were able to interbreed were genetically differentiated (Hemmer-Hansen et al. 2013a; Berg et al. 2016). The inclusion of both types of data may improve the understanding of the demographic history of highly migratory marine fish. The wild poplar SNP data were collected from the American Pacific Northwest. Male poplar trees produce pollen and female trees produce small seeds with fine hairs,

which enable dispersal of this species by wind (Geraldès et al. 2014). The samples covered various regions over a range of 2500 km near the Canadian–US border in British Columbia (BC) at alt between 0 and 800 m, where the variations in photoperiod and temperature have a north-south cline, while the variations in temperature, rainfall, and drought have an east-west (coastal to inland) cline (Geraldès et al. 2014). Each sampling location was associated with 11 environmental and geographical parameters. The analysis of environmental variables and population-specific F_{ST} values may provide a good example for understanding the history of the range expansion of a wind-dispersed species. By analyzing different types of data with species-specific ecology and migration history, the usefulness of our approach may be identified to enable us to understand the current population structure in an evolutionary context.

Materials and methods

Population colonization simulations

To test the performance of our visual representation, we conducted simulations that mimicked the colonization of populations from a single ancestral population (population 1). We modeled five types of stepping-stone colonization: one, two, and three-directional population expansion; three-directional grid colonization from an edge; and eight-directional grid colonization from the center, with 24 demes (populations 2–25) (Figure 1, A–E). Our objective is to describe current population structure using an unrooted NJ tree and an MDS plot, and to infer population history by overlaying population-specific F_{ST} values on the population structure. When gene flow is limited between adjacent populations, as shown in our simulations, the estimated population structure corresponds to population history.

We set the effective population size of the ancestral population to $N_e = 10^4$ (twice the number of individuals in diploid organisms in a random mating population). At the beginning of colonization, 1% of N_e migrated into the adjacent vacant habitat once every 10 generations. For convenience, we considered one simulation cycle to be one generation. The effective population size of the newly derived population increased to the same size as the ancestral population ($N_e = 10^4$) after one generation, and the population exchanged 1% of N_e genes with adjacent population(s) in every generation. Like the ancestral population, 1% of N_e individuals migrated into the adjacent vacant habitat once every 10 generations. We simulated the allele frequencies of SNPs in the ancestral and 24 derived populations. We also examined the cases in which the effective population size of the ancestral population was 10 times greater ($N_e = 10^5$) than that of the newly derived population ($N_e = 10^4$).

We generated the initial allele frequencies in the ancestral population, q , at 100,000 neutral SNP loci from the predictive equilibrium distribution, $f(q) \propto q^{-1}(1-q)^{-1}$ (Wright 1931). Additionally, we introduced 10 newly derived SNPs to each existing population in each generation. When a new SNP emerged in a population, we set the initial allele frequency of the newly derived SNP to 0.01 in the population and 0 in the other populations. This mimicked new mutations that survived in the initial phase after their birth. We considered these 100,000 ancestral SNPs and newly derived SNPs to be “unobserved.” We changed the allele frequencies of these SNPs using random drift under a binomial distribution in every generation. The frequencies of the derived alleles decreased for many of the SNPs over simulation generations and lost their polymorphism. After 260 simulation

generations, we randomly selected SNPs that retained their polymorphism as “observed” SNPs. For grid colonization models (Figure 1, D–E), we randomly selected polymorphic SNPs after 100 simulation generations.

In our simulations, we selected 10,000 ancestral SNPs and 500 newly derived SNPs. Then, we generated 50 individuals for each population. We randomly generated the genotypes of these 10,500 SNPs for each individual following the allele frequencies in the population to which each individual belonged. Thus, we obtained “observed” genotypes of 1,250 individuals (= 50 individuals \times 25 populations) at 10,500 SNP loci. To examine the effect of generations on genetic diversity in newly derived SNPs, we also selected 9000 and 7000 ancestral SNPs, and 1000 and 3000 newly derived SNPs, respectively. We converted the simulated genotypes into Genepop format (Raymond and Rousset, 1995; Rousset, 2008). We then computed genome-wide population-specific and pairwise F_{ST} values between the 25 populations.

Visual representation of population structure and demographic history

We integrated genome-wide population-specific and pairwise F_{ST} estimates on a map of sampling locations on an NJ tree and MDS plot. We visualized the magnitude of the genome-wide population-specific F_{ST} values using a color gradient based on $\text{rgb}(1 - F_{ST, \cdot 0}, 0, F_{ST, \cdot 0})$, where $F_{ST, \cdot 0} = (F_{ST} - \min F_{ST}) / (\max F_{ST} - \min F_{ST})$. This conversion represents the standardized magnitude of a population-specific F_{ST} value at the sampling point, with colors between red (for the smallest F_{ST}) and blue (for the largest F_{ST}). We drew the F_{ST} map using the `sf` package in R, where we plotted sampling locations based on the longitudes (lon) and lat. We connected sampling points with genome-wide pairwise F_{ST} values smaller than a given value using yellow lines to visualize the connectivity between populations. Under the assumption of Wright’s island model at equilibrium between drift, mutation, and migration (Wright 1931), $F_{ST} \approx \frac{1}{4N_e m + 1}$, where N_e is the effective population size and m is the average rate of migration between populations (Slatkin 1987). For example, $F_{ST} = 0.02$ refers to $4N_e m \approx 49$ migrants per generation (see Whitlock and McCauley 1999; Waples and Gaggiotti 2006). The value was arbitrarily used in our case studies. We plotted the genome-wide population-specific F_{ST} values on a dot chart with standard errors estimated using Equation (S5) (Supplementary Note). We drew the NJ tree based on the distance matrix of the genome-wide pairwise F_{ST} values using the `nj` function in the R package `ape`. We performed MDS analysis on the pairwise F_{ST} distance matrix using the `cmdscale` function in R. We used the cumulative contribution ratio up to the k th axis ($j = 1, \dots, k, \dots, K$) as the explained variation measure, which we computed using the R function as $C_k = \sum_{j=1}^k \lambda_j / \sum_{j=1}^K \lambda_j$, where λ_j is the eigenvalue and $\lambda_j = 0$ if $\lambda_j < 0$. We colored the sampling locations on the F_{ST} maps, dot charts, NJ trees, and MDS plots using a color gradient of the magnitude of genome-wide population-specific F_{ST} values. We also examined a diverging color palette instead of blue to red to test the resolution using `RColorBrewer` on the F_{ST} maps.

Computing F_{ST} values

We converted the genotype data into Genepop format (Raymond and Rousset 1995; Rousset 2008) for implementation in the R package `FinePop2_ver.0.2`. We computed genome-wide pairwise F_{ST} values [NC83, Supplementary Equation (S3)] using the `pop_pairwiseFST` function in `FinePop2`. We calculated expected

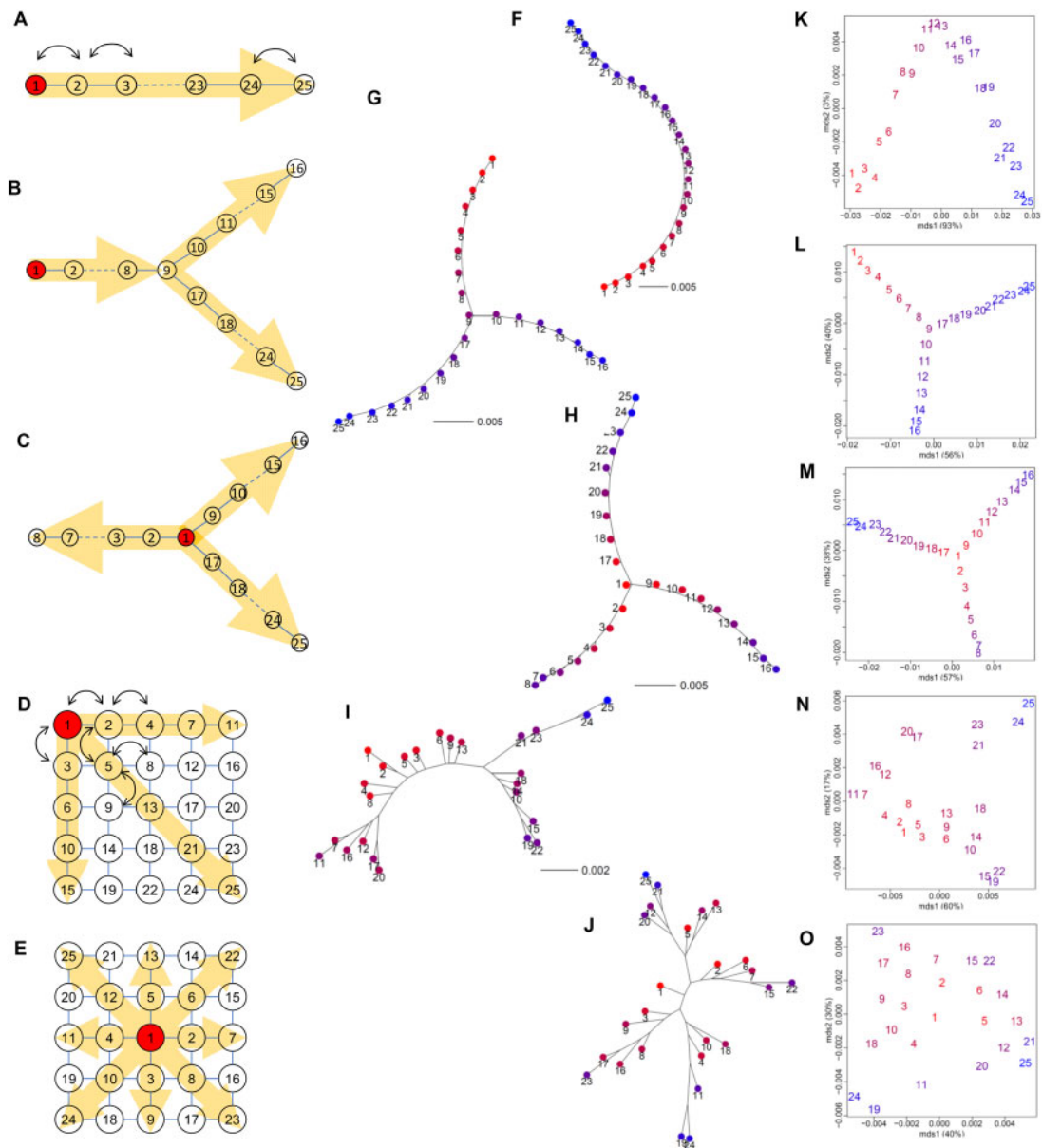


Figure 1 Results from population colonization simulations. Schematic diagrams of the models: (A) one, (B) two, (C) three-directional colonization, (D) three-directional grid colonization, and (E) eight-directional grid colonization. Population 1 in red is ancestral, and the yellow arrows indicate the direction of colonization. Lines show opportunities for migration. The effective population size of the newly derived population increased to the same size as the ancestral population ($N_e = 10^4$) after one simulation generation, and each population exchanged 1% of N_e genes with adjacent population(s) in every generation, as indicated by the arrows (see the text). Neighbor-joining (NJ) unrooted trees (F–J) and multi-dimensional scaling (MDS) plots (K–O) based on the pairwise F_{ST} distance matrix overlaid with population-specific F_{ST} values for each model. The color of each population shows the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

heterozygosity for each population using the read. GENEPOP function. We computed the genome-wide WG population-specific F_{ST} [Supplementary Equation (S4)] values using the `pop_specificFST` function. We applied Bayesian population-specific F_{ST} estimators on human data. We maximized Supplementary Equation (S7) and estimated the empirical Bayesian population-specific F_{ST} (Beaumont and Balding 2004) at each locus according to Supplementary Equation (S8). Then we averaged these values over all loci. For the full Bayesian model, we used `GESTE_ver. 2.0` (Foll and Gaggiotti 2006) to compute the genome-wide population-specific F_{ST} values. We examined the shrinkage effect of the Bayesian population-specific F_{ST} estimator

on inferring the ancestral population using a set of subsamples (37 populations) chosen from 51 populations.

Inferring environmental effects on the observed population structure

To infer the geography and environment that were experienced by the population range expansion, we regressed the genome-wide population-specific F_{ST} values on the geographical and environmental variables ($j = 1, \dots, s$):

$$ps_{ST}^i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_s x_{si} + \varepsilon_i, \quad \varepsilon \sim N(0, \Omega), \quad (i = 1, \dots, r), \quad (3)$$

where Ω is the variance matrix of population-specific F_{ST} . We correlated residuals because of the population structure; therefore, the effective sample size was lower than the actual sample size. In such circumstances, ordinary least squares overestimate the precision. To account for the correlation, we derived the components of the variance–covariance matrix of the population-specific F_{ST} estimator [Supplementary Equations (S5) and (S6)] for generalized least squares (GLS). We performed this analysis on the wild poplar dataset, for which 11 environmental/geographical parameters were available for each sampling location. We used the variance–covariance matrix for the components of the variance matrix Ω in Equation (3), and performed regression using the GLS function in FinePop2 v0.2.

Three empirical datasets

We retrieved the human microsatellite data used in Rosenberg et al. (2002) from <https://web.stanford.edu/group/rosenberglab/index.html>. We removed the Surui sample (Brazil) from the data because that population was reduced to 34 individuals in 1961 as a result of introduced diseases (Liu et al. 2006). We retained genotype data ($n = 1035$) of 377 microsatellite loci from 51 populations categorized into six groups, as in the original study: 6 populations from Africa, 12 from the Middle East and Europe, 9 from Central/South Asia, 18 from East Asia, 2 from Oceania, and 4 from America. We obtained the lon and lat of the sampling sites from Cann et al. (2002) (Supplementary Table S1).

We combined the Atlantic cod SNP genotype data of 924 SNPs common to 29 populations reported in Therkildsen et al. (2013a, 2013b) and 12 populations reported in Hemmer-Hansen et al. (2013a, 2013b). We compared the genotypes associated with each marker in samples that were identical in the two studies, that is, CAN08 and Western_Atlantic_2008, ISO02 and Iceland_migratory_2002, and ISCO2 and Iceland_stationary_2002, and standardized the gene codes. We removed *cgpGmo.S1035*, whose genotypes were inconsistent between the two studies. We also removed *cgpGmo.S1408* and *cgpGmo.S893*, for which the genotypes were missing in several population samples in Therkildsen et al. (2013b). For simplicity, we removed temporal replicates from the Norway migratory, Norway stationary, North Sea, and Baltic Sea samples. The final dataset consisted of genotype data ($n = 1065$) for 921 SNPs from 34 populations: 3 from Iceland, 25 from Greenland, 3 from Norway, and 1 each from Canada, the North Sea, and the Baltic Sea. All individuals in the samples were adults, and most were mature (Therkildsen et al. 2013a). We used the lon and lat of the sampling sites in Hemmer-Hansen et al. (2013a). For the data from Therkildsen et al. (2013a), we estimated approximate sampling points from the map of the original study and recorded the lon and lat (Supplementary Table S2).

We retrieved wild poplar SNP genotype data and environmental/geographical data from the original studies of McKown et al. (2014a, 2014b). The genotype data contained 29,355 SNPs of 3,518 genes of wild poplar ($n = 441$) collected from 25 drainage areas (McKown et al. 2014c). Details of the array development and selection of SNPs are provided in Geraldès et al. (2011, 2013). A breakdown of the 25 drainages (hereafter, populations) is as follows: 9 in northern British Columbia (NBC), 2 in inland British Columbia (IBC), 12 in southern British Columbia (SBC), and 2 in Oregon (ORE) (Geraldès et al. 2014). We combined the original names of the clusters and population numbers, and used them

as our population labels (NBC1, NBC3, ..., ORE30). We associated each sampling location with 11 environmental and geographical parameters: lat, lon, alt, longest day length (DAY), frost-free days (FFD), mean annual temperature (MAT), mean warmest month temperature (MWT), mean annual precipitation (MAP), mean summer precipitation (MSP), annual heat-moisture index (AHM), and summer heat-moisture index (SHM) (Supplementary Table S3). The AHM was calculated in the original study as $(MAT + 10)/(MAP/1000)$; a large AHM indicates extremely dry conditions.

Results

Simulations of population colonization

First, we examined the effect of the number of simulation generations on genetic diversity in newly derived SNPs using the eight-directional grid simulation (Figure 1E). WG population-specific F_{ST} correctly identified the ancestral population and traced the population history, and population structure reflected the population history regardless of the numbers of ancestral SNPs (9000 and 7000) and newly derived SNPs (1000 and 3000) selected after 100 simulation generations (Supplementary Figure S1). The result was consistent with the case that used 10,500 SNP loci (10,000 ancestral SNPs + 500 newly derived SNPs) (Figure 1, J and O and Supplementary Figure S2E). In the following analysis, we used the results based on 10,500 SNP loci to generate clearer results, even for limited numbers of simulation generations (260 and/or 100).

In the one-directional simulation (Figure 1A), population-specific F_{ST} correctly identified the ancestral population with the highest genetic diversity (Supplementary Figure S2A), and populations were located in order from 1 to 25 on the NJ tree (Figure 1F). The first axis of the MDS plot explained 93% of the variance of the pairwise F_{ST} distance matrix and indicated population expansion from populations 1 to 25 (Figure 1K). In the two-directional simulation (Figure 1B), our analysis correctly identified the ancestral population (Supplementary Figure S2B) and detected that populations were split at population 9 and expanded in two directions (Figure 1G), which was consistent with the simulation scenario. The first axis of the MDS plot identified population expansion from populations 1 to 25 and explained 56% of the variance of the pairwise F_{ST} distance matrix, whereas the second axis identified another manner of population expansion from populations 1 to 16 and explained 40% of the variation (Figure 1L). In the three-directional simulation (Figure 1C), the ancestral population was also correctly identified (Supplementary Figure S2C). It was closely located to the adjacent populations 2, 9, and 17, but correctly detected three directions (Figure 1H). The first axis of the MDS plot identified population expansion from population 1 to populations 16 and 25, and explained 57% of the variance of pairwise F_{ST} , whereas the second axis identified population expansion from population 1 to populations 8 and 16, and explained 38% of the variance (Figure 1M).

In the three-directional grid colonization model from an edge (Figure 1D), population-specific F_{ST} correctly identified the ancestral population (Supplementary Figure S2D) and pairwise F_{ST} detected that populations expanded in three directions (Figure 1I), which agreed with the simulation scenario. The first axis of the MDS plot identified population expansion from population 1 to other edge populations (populations 11, 15, and 25), and explained 60% of the variance of the pairwise F_{ST} distance matrix, whereas the second axis indicated genetic differentiation between populations 24 and 25, and 15, 19, and 22, and explained

17% of the variance (Figure 1N). In the eight-directional grid colonization model (Figure 1E), population-specific F_{ST} identified the ancestral population (Supplementary Figure S2E) and pairwise F_{ST} estimated that populations expanded in five directions from the center (Figure 1J). The first axis of the MDS plot identified vertical population expansion from population 1 to populations 24 and 25 and explained 40% of the variance of the pairwise F_{ST} distance matrix, and the second axis indicated horizontal population expansion from population 1 to populations 23 and 24, which explains 30% of the variance (Figure 1O). We obtained similar results in the cases in which the effective population size of the ancestral population was 10 times greater ($N_e = 10^5$) than that of the newly derived population ($N_e = 10^4$) (Supplementary Figure S3). We obtained very similar results from different data for more than 20 simulations (figures not shown).

Humans

The F_{ST} map (Figure 2A) shows integrated information from genome-wide population-specific and pairwise F_{ST} , which visualizes population structure with the migration and demographic history of human populations in terms of genetic diversity. Interestingly, Bantu Kenyans had the smallest F_{ST} value (shown

in red). Figure 2B ordered population-specific F_{ST} values from Africa to Central/South Asia, the Middle East, Europe, East Asia, Oceania, and America (Supplementary Table S4). As indicated by the sampling points connected by yellow lines with pairwise F_{ST} values below 0.02 (Figure 2A), genetic connectivity from Africa was low. Conversely, migration was substantial within Eurasia but much smaller than that inferred from Eurasia to Oceania and America. H_e was the highest in Africa, followed by the Middle East, Central/South Asia, Europe, and East Asia, but relatively small in Oceania and lowest in South America. The Karitiana in Brazil had the lowest H_e . The NJ tree (Figure 2C) integrated with population-specific F_{ST} values indicated that human populations originated from Bantu Kenyans and expanded to Europe through Mozabite, the Middle East, Central/South Asia, and East Asia. The Kalash was isolated from Europe/Middle East and Central/South Asia populations. Papuans/Melanesians and American populations diverged from between Central/South Asian and East Asian populations. The ordinal NJ tree of pairwise F_{ST} values divided the populations into five clusters: (1) Africa, (2) the Middle East, Europe, and Central/South Asia, (3) East Asia, (4) Oceania, and (5) America (Supplementary Figure S4). The first axis of the MDS plot highlighted differences between African and American

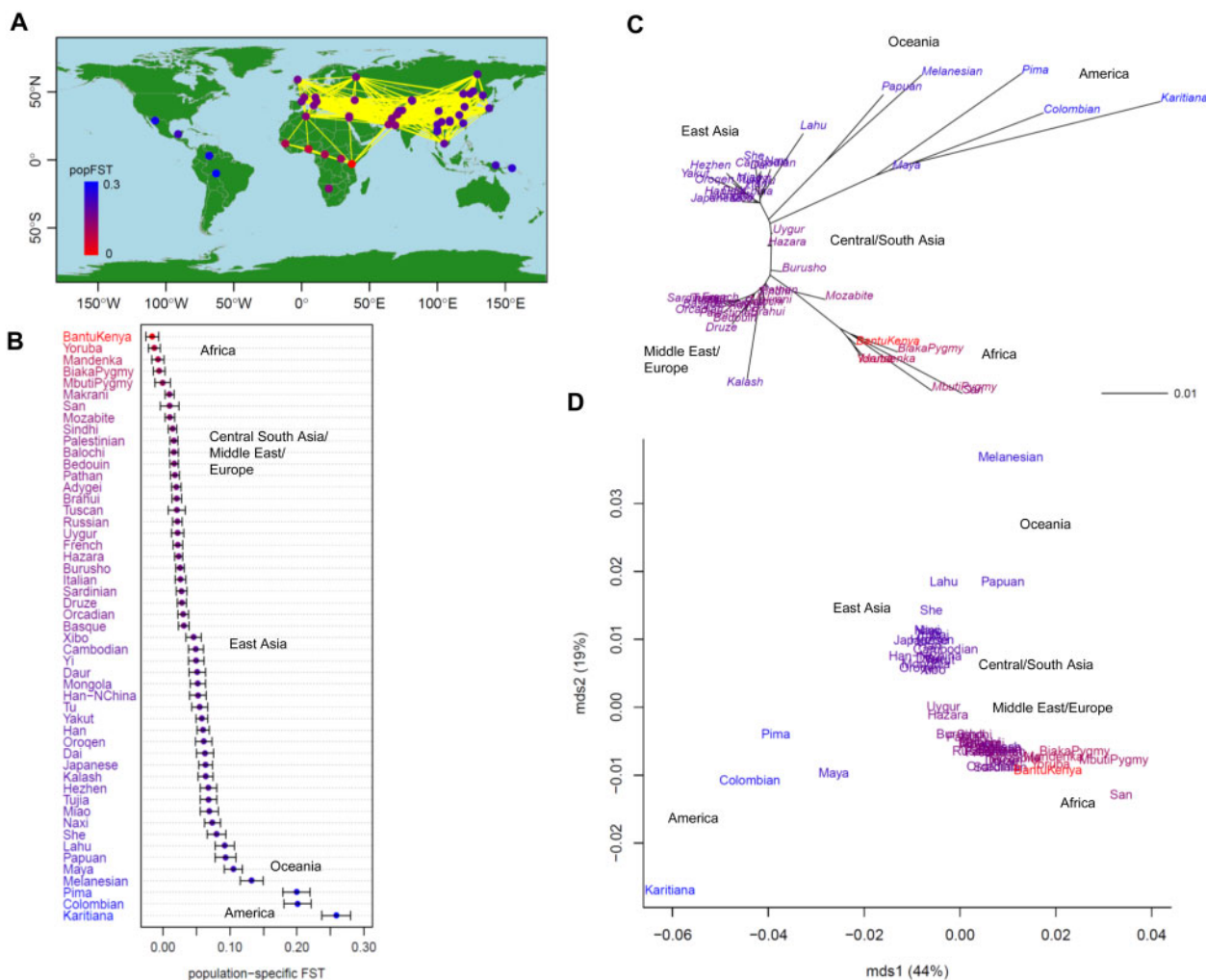


Figure 2 Population structure of 51 human populations ($n = 1035$; 377 microsatellites). (A) Map showing population connectivity with the magnitude of population-specific F_{ST} values. Populations connected by yellow lines are those with pairwise $F_{ST} < 0.02$. (B) Distribution of population-specific F_{ST} values $\pm 2 \times SE$. (C) Neighbor-joining (NJ) unrooted tree and (D) multi-dimensional scaling (MDS) based on pairwise F_{ST} overlaid with population-specific F_{ST} values on population labels. The color of each population indicates the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

populations and explained 44% of the variance of the pairwise F_{ST} distance matrix, whereas the second axis indicated genetic differentiation between Melanesian and Karitiana populations, and explained 19% of the variance (Figure 2D).

The Bayesian population-specific F_{ST} values estimated using the methods of Beaumont and Balding (2004) (empirical Bayes) and Foll and Gaggiotti (2006) (full Bayes) were nearly identical and the smallest F_{ST} values observed in the Middle East, Europe, and Central/South Asia (Supplementary Figure S5A, Supplementary Table S4). However, in African populations, they

were higher than the WG population-specific F_{ST} values (Figure S5B). Our F_{ST} map based on the empirical Bayesian population-specific F_{ST} values indicated that the Middle East, Europe, and Central/South Asia were centers of human origin (Figure 3A), which was consistent with that from the full Bayesian population-specific F_{ST} estimator (figure not shown). Our integrated NJ tree showed that the Hazara, Pakistan population was genetically closest to the human ancestors (Supplementary Figure S6A). The numbers of sampling locations of the 51 human populations were as follows: 6 from Africa, 12 from the Middle

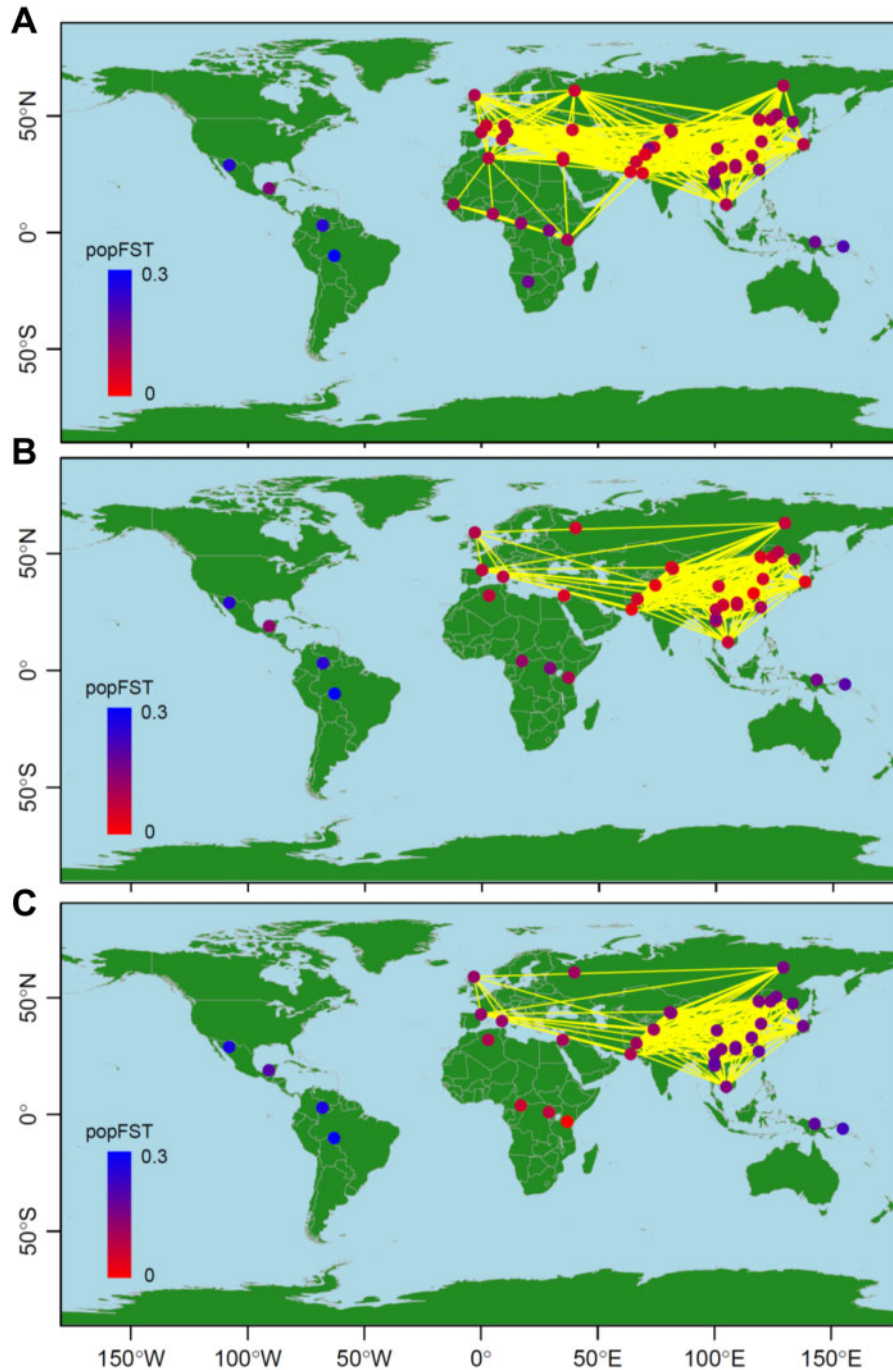


Figure 3 Population structure of humans based on Bayesian and moment estimators of population-specific F_{ST} . Results from the Bayesian population-specific F_{ST} estimator using (A) 51 samples and (B) 37 subsamples, and from (C) the WG population-specific F_{ST} moment estimator using 37 subsamples. The numbers of sampling locations of the subsamples were as follows: 3 from Africa, 6 from the Middle East/Europe, 9 from Central/South Asia, 18 from East Asia, 2 from Oceania, and 4 from America. Populations connected by yellow lines are those with pairwise $F_{ST} < 0.02$. The color of each population indicates the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

East/Europe, 9 from Central/South Asia, 18 from East Asia, 2 from Oceania, and 4 from America. When we used a subsample of the 37 human populations (3 from Africa, 6 from the Middle East/Europe, 4 from Central/South Asia, 18 from East Asia, 2 from Oceania, and 4 from America; Supplementary Table S5), the area with the highest population-specific F_{ST} values shifted toward Central/South Asia and East Asia (Figure 3B), whereas Bantu Kenyans had the smallest WG population-specific F_{ST} value (Figure 3C); this was consistent with the results from the full dataset (Figure 2A). The integrated NJ trees provided similar results (Supplementary Figure S6, B and C).

Atlantic cod

The F_{ST} map (Figure 4A) visualizes the population structure, migration, and genetic diversity of the Atlantic cod populations. The Canadian population had the smallest population-specific F_{ST} value (shown in red) and the greatest H_e . H_e was also high in Greenland, low in other areas, and lowest in the Baltic Sea. Figure 4B shows the order of population-specific F_{ST} values from Canada to the Baltic sea (Supplementary Table S6). Greenland west coast populations (green in Supplementary Figure S7) generally had small population-specific F_{ST} values, whereas fjord populations (violet) had relatively higher population-specific F_{ST} values. The population-specific F_{ST} values were much higher for

populations in Iceland, Norway, and the North Sea, and were highest in the Baltic Sea. Based on pairwise F_{ST} values (< 0.02) between sampling points (Figure 4A), substantial migration was suggested between Greenland, Iceland, and Norway. Conversely, migration could be low from Canada to Greenland and from Iceland and Norway to the North and Baltic Seas. Our integrated NJ tree with population-specific F_{ST} values (Figure 4C) inferred that Atlantic cod originated from Canada, migrated to the west coast of Greenland, and then expanded their distribution to Iceland, Norway, the North Sea, and the Baltic Sea. According to the ordinal NJ tree of the pairwise F_{ST} distance matrix (Supplementary Figure S7), the populations were divided into four large clusters: (1) Canada; (2) Greenland west coast, (3) Greenland east coast, Iceland, and Norway; and (4) North and Baltic Seas. Fjord populations (in purple) formed a sub-cluster within the Greenland west coast, and migratory (orange) and stationary (magenta) ecotypes also formed a sub-cluster. The first axis of the MDS plot characterized the differentiation between Canadian and North Sea/Baltic Sea populations and explained 72% of the variance of the pairwise F_{ST} distance matrix, whereas the second axis highlighted the differentiation between Norwegian migratory populations and Canadian and North Sea/Baltic Sea populations, which explained 22% of the variance (Figure 4D).

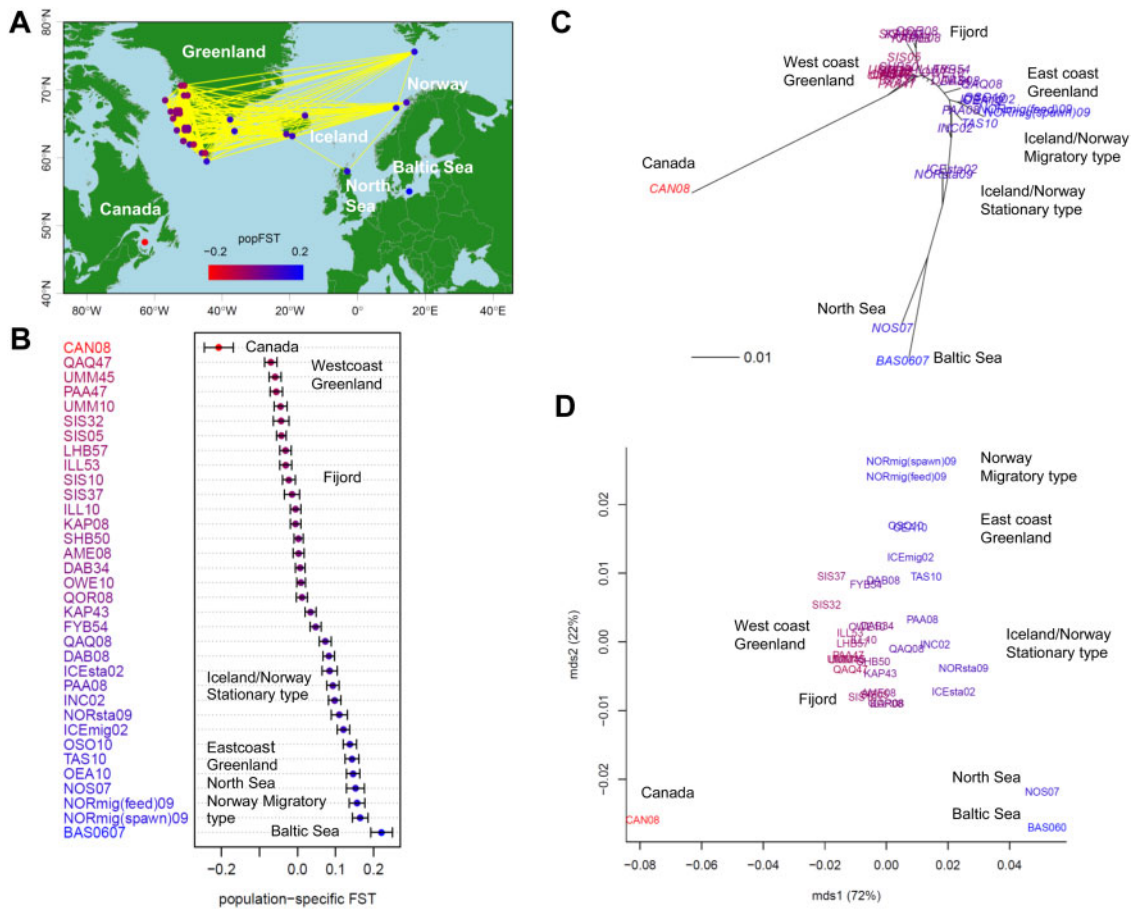


Figure 4 Population structure of 34 geographical samples of wild Atlantic cod ($n = 1065$; 921 SNPs). (A) Map showing population connectivity with the magnitude of population-specific F_{ST} values. Populations connected by yellow lines are those with pairwise $F_{ST} < 0.02$. (B) Distribution of population-specific F_{ST} values $\pm 2 \times SE$. (C) Neighbor-joining (NJ) unrooted tree and (D) multi-dimensional scaling (MDS) based on pairwise F_{ST} overlaid with population-specific F_{ST} values on population labels. The color of each population shows the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

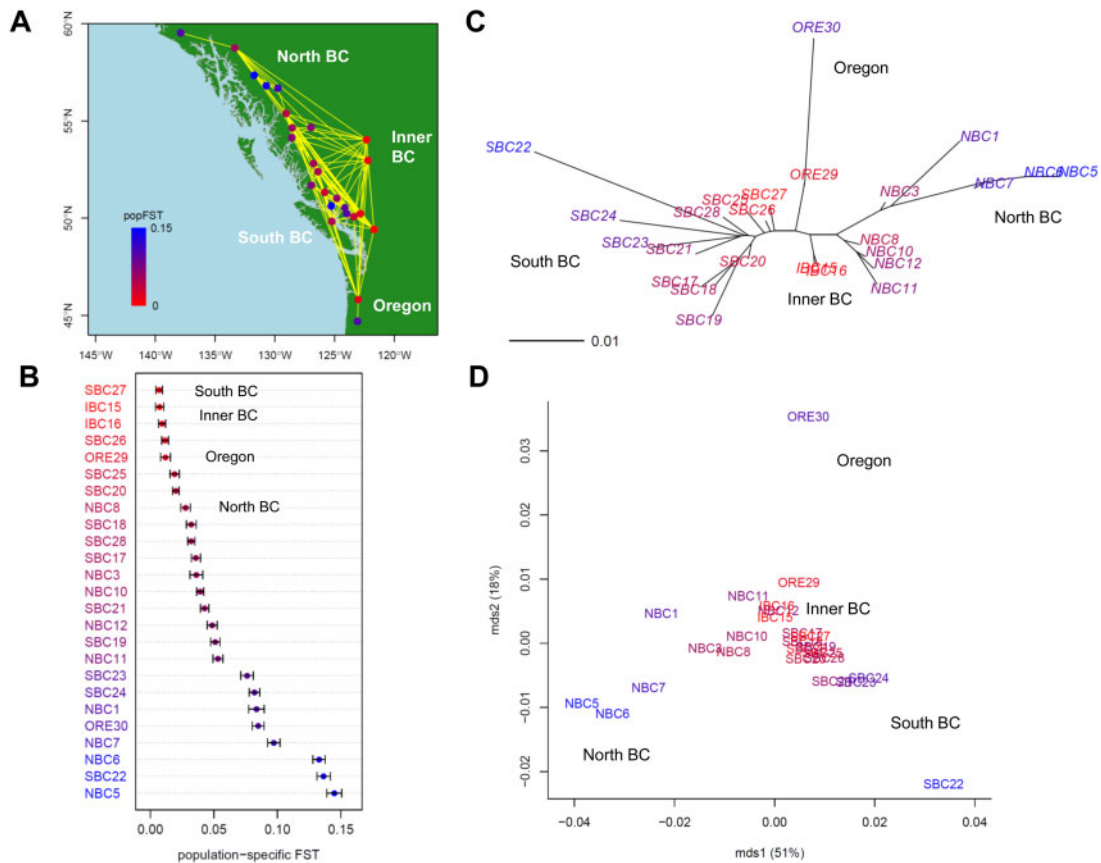


Figure 5 Population structure for 25 geographical samples of wild poplar ($n = 441$; 29,355 SNPs). (A) Map showing population connectivity with the magnitude of population-specific F_{ST} values. Populations connected by yellow lines are those with pairwise $F_{ST} < 0.02$. (B) Distribution of population-specific F_{ST} values $\pm 2 \times SE$. (C) Neighbor-joining (NJ) unrooted tree and (D) multi-dimensional scaling (MDS) based on pairwise F_{ST} overlaid with population-specific F_{ST} values on population labels. The color of each population indicates the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

Wild poplar

The F_{ST} map (Figure 5A) indicated that population-specific F_{ST} values were lowest in SBC27 and inner British Columbia (IBC15 and IBC16) (shown in red, Supplementary Figure S8). The sampling points connected by yellow lines (pairwise $F_{ST} < 0.02$) indicated migration between all populations. H_e was highest in SBC27, IBC15, and IBC16, and lowest in NBC5. Figure 5B shows that samples collected from areas close to the SBC coast had higher population-specific F_{ST} values than other SBC samples (Supplementary Table S7). The NBC samples had population-specific F_{ST} values similar to those of SBC. Among the NBC samples, NBC8 had the smallest population-specific F_{ST} , and NBC5 had the highest value, followed by NBC6 and NBC7. The pairwise F_{ST} NJ tree integrated with population-specific F_{ST} values (Figure 5C) suggested that wild poplar originated from around SBC27 and interior BC, and expanded in three directions: to the southern coast of BC, NBC and south-western Alaska, and ORE. The ordinal NJ tree based on the pairwise F_{ST} distant matrix divided populations into four large clusters: (1) IBC, (2) SBC, (3) NBC, and (4) ORE (Supplementary Figure S8). The population represented by sample ORE30 was isolated from ORE29. The first axis of the MDS plot characterized the differentiation between southern and northern populations and explained the 51% variance of the pairwise F_{ST} distance matrix, whereas the second axis characterized the southernmost ORE30 population, and explained the 18% variation (Figure 5D).

Table 1 Regression of genome-wide population-specific F_{ST} of 25 wild poplar populations on environmental variables

Variable	Estimate	SE	Z	P
DAY	0.0489	0.0164	2.99	0.003**
MAT	-0.0088	0.0086	-1.03	0.305
MAP	0.0001	0.0000	2.79	0.005**
SHM	0.0022	0.0009	2.38	0.018*

DAY: longest day length (h); MAT: mean annual temperature ($^{\circ}C$); MAP: mean annual precipitation (mm); SHM: summer heat-moisture index.

* $p < 0.05$.
** $p < 0.01$.

To avoid multicollinearity, we excluded 7 out of 11 environmental variables that were significantly correlated with each other: lat, lon, alt, FFD, MWMT, MSP, and AHM (Supplementary Table S3). Our GLS of genome-wide population-specific F_{ST} values on the four environmental variables (DAY, MAT, MAP, and SHM) indicated that DAY, MAP, and SHM were significant (Table 1). All estimates were positive, which indicated that higher population-specific F_{ST} values were expected for longer DAY (longer daylight time), higher MAP (abundant rain), and higher SHM (dry summers), and these values might reflect the directions of population expansion. The scatter plot of DAY and SHM (with each population colored according to the population-specific F_{ST} value) (Figure 6A) suggested three directions of population range expansion: the wild poplar that might have originated from around

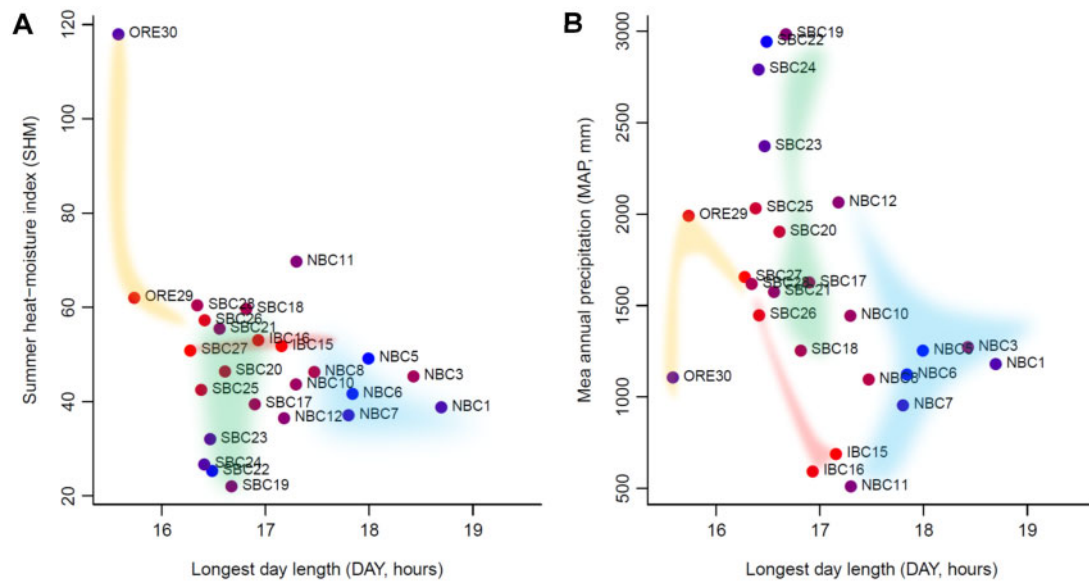


Figure 6 Population range expansion and key environmental variables. Longest day length vs (A) summer heat-moisture index and (B) mean annual precipitation for 25 geographical samples of wild poplar. The colored areas by the population clusters (see Supplementary Figure S8) show the inferred population expansion from IBC15, IBC16, and SBC27. The color of each population shows the magnitude of population-specific F_{ST} values between red (for the smallest F_{ST}) and blue (for the largest F_{ST}).

SBC27 and IBC15 expanded its distribution to NBC, where daylight hours are long in summer, as well as expanding to coastal SBC with its lower SHM (humid summer and abundant rainfall), and to the south (ORE29 and ORE30) with its higher SHM (dry summer). This was consistent in the scatter plot of DAY and MAP (Figure 6B).

Diverging color palette

RColorBrewer had 35 color palettes. Each palette had a minimum of eight colors. Two palettes had 12 colors (maximum) and nine had 11 colors. We chose a color palette with 10 colors (RdYlBu). The color gradient better identified the middle range of population-specific F_{ST} values, but failed to detect the ancestral population (Supplementary Figure S9).

CPU times

With an Intel Core i7-8650U CPU, 89.8 s of CPU time were required to compute the WG population-specific F_{ST} estimates and SEs of wild poplar (29,355 SNPs; 25 populations, $n = 441$). To obtain the pairwise F_{ST} (NC83) between all 300 ($= 25 \times 24/2$) population pairs, 120.7 s were required. CPU time is proportional to the number of SNPs analyzed.

Discussion

Genome-wide population-specific F_{ST} traced population history as reflected by genetic diversity

Our simulations demonstrated that the WG population-specific F_{ST} estimator identified the source population and traced the evolutionary history of its derived populations based on genetic diversity (heterozygosity estimated from each population). The NC83 pairwise F_{ST} estimator correctly estimated the current population structure. As explained in the Introduction, the population-specific F_{ST} estimator is a rescaling of expected heterozygosity, and we expect a linear relationship between expected heterozygosity and population F_{ST} . This shows that the

population-specific F_{ST} estimator implicitly assumes that populations closest to the ancestral population have the highest heterozygosity. In our three case studies, a linear relationship between H_e of each population ($= H_{Si}$) and psF_{ST}^i was evident (Supplementary Figure S10). The coefficient of determination, R^2 , was 0.91 for 51 human populations ($n = 1035$), 0.99 for 34 Atlantic cod populations ($n = 1065$), and 0.82 for 25 wild poplar populations ($n = 441$). The goodness of fit to the linear function should depend on the sample size (number of individuals). Our simulations evaluated the performance of the population-specific F_{ST} estimator for such cases. However, in populations that experienced extensive admixture events, heterozygosity was enhanced, whereas a bottleneck in the ancestral population reduced heterozygosity. In such cases, the population-specific F_{ST} estimator misidentifies the ancestral population.

In our analysis, the genome-wide WG population-specific F_{ST} values successfully illustrated human evolutionary history, and indicated that humans originated in Kenya, expanded from the Middle East into Europe and from Central/South Asia into East Asia, and then possibly migrated to Oceania and America (Figure 2). Kenya is located just below Ethiopia, where the earliest anatomically modern humans were found from fossils (Nielsen *et al.* 2017). Our results are also in good agreement with the highest levels of genetic diversity being detected in Africa (Rosenberg *et al.* 2002), the relationship uncovered between genetic and geographic distance (Ramachandran *et al.* 2005), the shortest colonization route from East Africa (Liu *et al.* 2006), and major migrations inferred from genomic data (Nielsen *et al.* 2017).

Our analysis indicated that Atlantic cod might originate in Canada (CAN08). Figure 4 suggested that the population expansion of Atlantic cod began by minimal gene flow from Canada. They might have first expanded to the west coast of Greenland before spreading to Iceland, the North Sea, Norway, and the Baltic Sea. This result was consistent with genomic evidence that Atlantic cod inhabit both sides of the Atlantic Ocean and evolved from a common evolutionary origin (Berg *et al.* 2017). The migratory ecotypes characterized by deeper and more offshore habitats

and long-distance migrations (Hemmer-Hansen et al. 2013a) may have played an important role in this expansion, as suggested in Figure 4C and Supplementary Figure S7. In our study, CAN08 had the highest H_e , which was lower in Iceland than in Greenland; this result implies that Icelandic populations were the descendants of colonists from Greenland, which in turn originated in Canada. The BAS0607 sample from the Baltic Sea had the highest population-specific F_{ST} and the lowest H_e , suggesting that Baltic cod is the newest population. This result agrees with the findings of a previous study, which identified Baltic cod as an example of a species subject to ongoing selection for reproductive success in a low salinity environment (Berg et al. 2015). In the Atlantic cod case study, CAN08 had the highest H_e and a very large negative population-specific F_{ST} value of -0.21 ± 0.019 compared with the maximum value of 0.22 ± 0.014 in BAS0607 (Supplementary Figure S10, Supplementary Table S6). The WG population-specific F_{ST} value can be negative (Weir and Goudet 2017). In the one and two-directional models of our simulations, the WG population-specific F_{ST} value was significantly negative in the ancestral population, whereas H_e was the largest (Supplementary Figure S2, A and B). Our consistent results between the simulations and Atlantic cod case study indicate that when gene flow from other populations into the source population is limited, a relatively large H_e (\hat{H}_{Si}) is maintained in the source population. In such cases with $\hat{H}_{Si} > \hat{H}_B$, Equation (2) produces negative values for population-specific F_{ST} .

Although the wild poplar samples used in this study might not cover the entire distribution range of the species, which extends from southern California to northern Alaska, Montana, and Idaho (Geraldes et al. 2013), the genome-wide population-specific F_{ST} values suggested three directions of population expansion of wild poplar: from SBC27 and IBC (IBC15, IBC16) to coastal British Columbia, southern ORE, and NBC (Figure 5). The largest population-specific F_{ST} value was found in the population with the smallest heterozygosity, SBC22, which may have resulted from a bottleneck (Geraldes et al. 2014).

Our continuous color gradient from red to blue successfully detected the ancestral population, whereas the R diverging color palettes, which had a limited number of colors (maximum of 12), better identified the middle range of population-specific F_{ST} values (Supplementary Figure S9). Both color gradients may be useful.

Genome-wide pairwise F_{ST} described current population structure

Our stepping-stone simulations did not account for long-range dispersal (Hallatschek and Fisher 2014). Human samples were collected from present populations, but they might reflect history from the Age of Discovery, when humans were travelling far beyond their native continents (Diamond 1997). Wild animals primarily move locally, but occasionally disperse over long distances (Hallatschek and Fisher 2014). The migratory ecotype of the Atlantic cod is characterized by its long-distance migration (Hemmer-Hansen et al. 2013a). Sea currents also play a role for passive transportation of fertilized eggs and juveniles. The pollen and small seeds with fine hairs of poplar trees enable long-range dispersal by wind (Geraldes et al. 2014). Our stepping-stone simulations were conducted with the assumption that dispersal was step-by-step, but long-range dispersals were not taken into consideration. When gene flow is limited between adjacent populations, as in our simulation scenarios, estimated population structure reflects population history.

Genome-wide population-specific F_{ST} detects key environments that relate to population expansion

Our GLS of genome-wide population-specific F_{ST} values revealed that long daylight hours, abundant rainfall, and dry summer conditions are the key environmental factors that relate to the demographic history of wild poplar (Table 1). Wild poplar could have expanded its distribution by its fluffy seeds being blown away by the wind. In the NJ unrooted tree, the root of the populations cannot be fixed without out-group populations. There is no direct evidence, but we can infer from the population-specific F_{ST} values that wild poplar seems to have spread from SBC (SBC27) and IBC (IBC15, 16), and expanded its distribution to NBC, where daylight hours are long in summer, to coastal SBC with its rainy environment, and to southern ORE (ORE30) with its dry summer conditions (Figures 5 and 6). A previous study on wild poplar revealed that genes involved in drought response were identified as F_{ST} outliers using BayeScan (Foll and Gaggiotti 2008; Geraldes et al. 2014). The F_{ST} outlier test of Geraldes and colleagues also revealed that genes involved in the circadian rhythm and response to red/far-red light had high locus-specific global F_{ST} values. The first principal component of SNP allele frequencies of the poplar tree was significantly correlated with day length, and a previous enrichment analysis for population structuring uncovered genes related to circadian rhythm and photoperiod. The circadian clock pathway might play a central role in adaptation, and clinal variation in phenological traits might be an adaptive response to the north-south climatic variation of *P. trichocarpa* (Geraldes et al. 2014; McKown et al. 2014a). Our GLS analysis does not detect a locus-specific effect of environmental adaptation like genome scan methods, but detects key environmental variables that affected the population history through genome-wide population-specific F_{ST} . Our results confirmed the previous findings and the hypothesis of postglacial northward expansion of the poplar tree to refugia north of the southern margin of the ice sheet (Geraldes et al. 2014), showing the usefulness of applying the GLS estimate of genome-wide population-specific F_{ST} to infer environmental effects on the population expansion of species.

The two types of F_{ST} s can also be calculated for the genomic windows. By comparing the pairwise F_{ST} among genomic windows, it is possible to identify the genomic regions that are largely differentiated due to local adaptation by using the population branch statistics (Yi et al. 2010). The local principal component analysis using lostruct software (Li and Ralph 2019) and MDS (Fuller et al. 2020) also identify genomic regions linked to local adaptation. Likewise, by comparing population-specific F_{ST} among windows, it is possible to identify unique genomic regions for adaptation (Akey et al. 2002; Weir et al. 2005; Oleksyk et al. 2008).

Genome-wide F_{ST} moment estimators converge to their true means

Previous studies have suggested or indicated that the “ratio of averages” works better than the “average of ratios” as the number of independent SNPs increases (Reynolds et al. 1983; Weir and Cockerham 1984; Bhatia et al. 2013). Because “the combined ratio estimate (ratio of averages) is much less subject to the risk of bias than the separate estimate (average of ratios)” (Cockran 1977), scholars recommend using the “ratio of averages” estimators (Bhatia et al. 2013). To explicitly show the underlying mechanism, we used the observed heterozygosity of population *i* (\hat{H}_{Si}) as derived by Nei and Chesser (1983) (Supplementary Note). When the number of loci (*L*) increases, the average observed heterozygosity

over all loci converges to its expected value according to the law of large numbers as:

$$\frac{1}{L} \sum_{l=1}^L \left(1 - \sum_{u=1}^m \tilde{p}_{iu}^2\right) \rightarrow \frac{1}{L} \sum_{l=1}^L \left(1 - E\left[\sum_{u=1}^m \tilde{p}_{iu}^2\right]\right)$$

The observed gene diversity thus converges to the expected value:

$$\hat{H}_{Si} \left(1 - \frac{1}{n_i}\right) + \frac{\hat{H}_{oi}}{2n_i} \rightarrow H_{Si} \left(1 - \frac{1}{n_i}\right) + \frac{H_{oi}}{2n_i}$$

Similarly, \hat{H}_S and \hat{H}_T converge to their expected values. This example indicates that the numerators and denominators of bias-corrected F_{ST} moment estimators, whether global, pairwise, or population-specific, converge to their true means and provide unbiased estimates of F_{ST} in population genomics analyses with large numbers of SNPs.

Bayesian F_{ST} estimators measure the deviation from the average of the sampled populations

In the Bayesian framework, the population-specific F_{ST} is the coefficient of the genetic drift that represents the among-population variation of the allele frequencies at neutral loci from the allele frequencies of the ancestral population. The allele frequency of the ancestral population is assumed to be the among-population mean allele frequency. The analysis of human populations (Figure 3A) highlights the need to account for geographical heterogeneity in sampling fractions of populations. The unbiased estimate of the allele frequency of the ancestral population will be obtained as the weighted among-population average. The weights are inversely proportional to the sampling fractions.

The shrinkage effect on allele frequencies in Bayesian inference (Stein, 1956) may shift population-specific F_{ST} values toward the average of the entire population. Because of the shrinkage toward mean allele frequencies, the maximum likelihood and Bayesian estimators of locus-specific global F_{ST} improve the power to detect genes under environmental selection (Beaumont and Balding 2004; Foll and Gaggiotti 2008). An empirical Bayes genome-wide pairwise F_{ST} estimator (Kitada et al. 2007) is useful in cases involving a small number of polymorphic marker loci, particularly in high gene flow scenarios, but it suffers from the shrinkage effect when larger numbers of loci are used. The shrinkage of allele frequencies should affect inference in genome-wide population-specific F_{ST} , particularly in cases when samples (populations) were not representative of the populations.

Conclusions

The WG genome-wide population-specific F_{ST} moment estimator can identify the source population and trace the evolutionary history of the derived populations based on genetic diversity under the assumption that populations closest to the ancestral population have the highest heterozygosity. Conversely, the NC83 genome-wide pairwise F_{ST} moment estimator represents the current population structure. By integrating population-specific and pairwise estimates on F_{ST} maps, NJ trees, and MDS plots, we obtain a picture of population structure by incorporating evolutionary history. Our GLS analysis of genome-wide population-specific F_{ST} , which accounts for the

correlation between populations, provides insights into how a species expanded its distribution in different environments. Given a large number of loci, bias-corrected F_{ST} moment estimators—whether global, pairwise, or population-specific—provide unbiased estimates of F_{ST} supported by the law of large numbers. Genomic data highlight the usefulness of the bias-corrected moment estimators of F_{ST} .

Data availability

The authors affirm that all data necessary for confirming the conclusions of the article are present within the article, figures, tables, and supplementary material. The R codes for our representation method exemplified by the human data and simulations of population colonization used in this study are available in the Supplementary material at figshare: <https://doi.org/10.25387/g3.14813490>.

Acknowledgments

The authors express appreciation to editors Jeffrey Ross-Ibarra, Andrew Kern, Mark Beaumont, and Nicholas Barton for reviewing this manuscript and providing constructive comments. They also thank the reviewers for providing essential comments that significantly improved the earlier versions of the manuscript.

Funding

This study was supported by Japan Society for the Promotion of Science Grants-in-Aid for Scientific Research KAKENHI nos. 16H02788 and 19H04070 to HK, and 18K0578116 to SK.

Conflicts of interest

The authors declare that there is no conflict of interest.

Literature cited

- Akey JM, Zhang G, Zhang K, Jin L, Shriver MD. 2002. Interrogating a high-density SNP map for signatures of natural selection. *Genome Res.* 12:1805–1814. <https://doi.org/10.1101/gr.631202>
- Balding DJ, Nichols RA. 1995. A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity. *Genetica.* 96:3–12. <https://doi.org/10.1007/BF01441146>
- Balloux F, Lugon-Moulin N. 2002. The estimation of population differentiation with microsatellite markers. *Mol Ecol.* 11:155–165. <https://doi.org/10.1046/j.0962-1083.2001.01436.x>
- Beaumont MA. 2005. Adaptation and speciation: what can F_{ST} tell us? *Trends Ecol Evol.* 20:435–440. <https://doi.org/10.1016/j.tree.2005.05.017>
- Beaumont MA, Balding DJ. 2004. Identifying adaptive genetic divergence among populations from genome scans. *Mol Ecol.* 13: 969–980. <https://doi.org/10.1111/j.1365-294X.2004.02125.x>
- Berg PR, Jentoft S, Star B, Ring KH, Knutsen H, et al. 2015. Adaptation to low salinity promotes genomic divergence in Atlantic cod (*Gadus morhua* L.). *Genome Biol Evol.* 7:1644–1663. <https://doi.org/10.1093/gbe/evv093>
- Berg PR, Star B, Pampoulie C, Bradbury IR, Bentzen P, et al. 2017. Trans-oceanic genomic divergence of Atlantic cod ecotypes is associated with large inversions. *Heredity (Edinb).* 119:418–428. <https://doi.org/10.1038/hdy.2017.54>

- Berg PR, Star B, Pampoulie C, Sodeland M, Barth JM, et al. 2016. Three chromosomal rearrangements promote genomic divergence between migratory and stationary ecotypes of Atlantic cod. *Sci Rep*. 6:23246. <https://doi.org/10.1038/srep23246>
- Bhatia G, Patterson N, Sankararaman S, Price AL. 2013. Estimating and interpreting F_{ST} : the impact of rare variants. *Genome Res*. 23:1514–1521. <http://www.genome.org/cgi/doi/10.1101/gr.154831.113>
- Bradbury IR, Bentzen P. 2007. Non-linear genetic isolation by distance: implications for dispersal estimation in anadromous and marine fish populations. *Mar Ecol Prog Ser*. 340:245–257. <https://doi.org/10.3354/meps340245>
- Buckleton J, Curran J, Goudet J, Taylor D, Thiery A, et al. 2016. Population-specific F_{ST} values for forensic STR markers: a worldwide survey. *Forensic Sci Int Genet*. 23:91–100. <https://doi.org/10.1016/j.fsigen.2016.03.004>
- Cann HM, De Toma C, Cazes L, Legrand MF, Morel V, et al. 2002. A human genome diversity cell line panel. *Science*. 296:261–262. <https://doi.org/10.1126/science.296.5566.261b>
- Cockerham CC. 1969. Variance of gene frequencies. *Evolution*. 23:72–84. <https://doi.org/10.1111/j.1558-5646.1969.tb03496.x>
- Cockerham CC. 1973. Analyses of gene frequencies. *Genetics*. 74:679–700. PubMed 17248636
- Cockran WG. 1977. *Sampling Techniques*. New York, NY: Wiley.
- Diamond J. 1997. *Guns, Germs and Steel: The Fates of Human Societies*. London, UK: Random House.
- Excoffier L. 2007. Analysis of population subdivision. In: DJ Balding, M Bishop and C Cannings, editors. *Handbook of Statistical Genetics*. 3rd ed. Chichester, UK: Wiley. p. 980–1020.
- Falush D, Stephens M, Pritchard JK. 2003. Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. *Genetics*. 164:1567–1587.
- Foll M, Gaggiotti O. 2006. Identifying the environmental factors that determine the genetic structure of populations. *Genetics*. 174:875–891. <https://doi.org/10.1534/genetics.106.059451>
- Foll M, Gaggiotti OE. 2008. A genome-scan method to identify selected loci appropriate for both dominant and codominant markers: a Bayesian perspective. *Genetics*. 180:977–993. <https://doi.org/10.1534/genetics.108.092221>
- Fuller ZL, Mocellin VJ, Morris LA, Cantin N, Shepherd J, et al. 2020. Population genetics of the coral *Acropora millepora*: toward genomic prediction of bleaching. *Science*. 369:eaba4674. <https://doi.org/10.1126/science.aba4674>
- Gaggiotti OE, Foll M. 2010. Quantifying population structure using the F-model. *Mol Ecol Resour*. 10:821–830. <https://doi.org/10.1111/j.1755-0998.2010.02873.x>
- Geraldes A, Difazio SP, Slavov GT, Ranjan P, Muchero W, et al. 2013. A 34K SNP genotyping array for *Populus trichocarpa*: design, application to the study of natural populations and transferability to other *Populus* species. *Mol Ecol Resour*. 13:306–323. <https://doi.org/10.1111/1755-0998.12056>
- Geraldes A, Farzaneh N, Grassa CJ, McKown AD, Guy RD, et al. 2014. Landscape genomics of *Populus trichocarpa* the role of hybridization limited gene flow and natural selection in shaping patterns of population structure. *Evolution*. 68:3260–3280. <https://doi.org/10.1111/evo.12497>
- Geraldes A, Pang J, Thiessen N, Cezard T, Moore R, et al. 2011. SNP discovery in black cottonwood (*Populus trichocarpa*) by population transcriptome resequencing. *Mol Ecol Resour*. 11 Suppl 1:81–92. <https://doi.org/10.1111/j.1755-0998.2010.02960.x>
- Hallatschek O, Fisher DS. 2014. Acceleration of evolutionary spread by long-range dispersal. *Proc Natl Acad Sci U S A*. 111:E4911–E4919. <https://doi.org/10.1073/pnas.1404663111>
- Hellenthal G, Busby GB, Band G, Wilson JF, Capelli C, et al. 2014. A genetic atlas of human admixture history. *Science*. 343:747–751. <https://doi.org/10.1126/science.1243518>
- Hemmer-Hansen J, Nielsen EE, Therkildsen NO, Taylor MI, Ogden R, et al.; FishPopTrace Consortium. 2013a. A genomic island linked to ecotype divergence in Atlantic cod. *Mol Ecol*. 22:2653–2667. <https://doi.org/10.1111/mec.12284>
- Hemmer-Hansen J, Nielsen EE, Therkildsen NO, Taylor MI, Ogden R, et al. 2013b. Data from: a genomic island linked to ecotype divergence in Atlantic cod. Dryad, Dataset. <https://doi.org/10.5061/dryad.9gf10>
- Holsinger KE, Weir BS. 2009. Genetics in geographically structured populations: defining, estimating and interpreting F_{ST} . *Nat Rev Genet*. 10:639–650. <https://doi.org/10.1038/nrg2611>
- Jorde PE, Søvik G, Westgaard JI, Albrechtsen J, André C, et al. 2015. Genetically distinct populations of northern shrimp, *Pandalus borealis*, in the North Atlantic: adaptation to different temperatures as an isolation factor. *Mol Ecol*. 24:1742–1757. <https://doi.org/10.1111/mec.13158>
- Kanitz R, Guillot EG, Antoniazza S, Neuenschwander S, Goudet J. 2018. Complex genetic patterns in human arise from a simple range-expansion model over continental landmasses. *PLoS One*. 13:e0192460. <https://doi.org/10.1371/journal.pone.0192460>
- Kitada S, Kitakado T, Kishino H. 2007. Empirical Bayes inference of pairwise F_{ST} and its distribution in the genome. *Genetics*. 177:861–873. <https://doi.org/10.1534/genetics.107.077263>
- Kitada S, Nakamichi R, Kishino H. 2017. The empirical Bayes estimators of fine-scale population structure in high gene flow species. *Mol Ecol Resour*. 17:1210–1222. <https://doi.org/10.1111/1755-0998.12663>
- Li H, Ralph P. 2019. Local PCA shows how the effect of population structure differs along the genome. *Genetics*. 211:289–304. <https://doi.org/10.1534/genetics.118.301747>
- Limborg MT, Helyar SJ, De Bruyn M, Taylor MI, Nielsen EE, et al.; FPT Consortium. 2012. Environmental selection on transcriptome-derived SNPs in a high gene flow marine fish, the Atlantic herring (*Clupea harengus*). *Mol Ecol*. 21:3686–3703. <https://doi.org/10.1111/j.1365-294X.2012.05639.x>
- Lipson M, Loh PR, Levin A, Reich D, Patterson N, et al. 2013. Efficient moment-based inference of admixture parameters and sources of gene flow. *Mol Biol Evol*. 30:1788–1802. <https://doi.org/10.1093/molbev/mst099>
- Liu H, Prugnolle F, Manica A, Balloux F. 2006. A geographically explicit genetic model of worldwide human-settlement history. *Am J Hum Genet*. 79:230–237. <https://doi.org/10.1086/505436>
- McKown AD, Guy RD, Klápště J, Geraldes A, Friedmann M, et al. 2014a. Geographical and environmental gradients shape phenotypic trait variation and genetic structure in *Populus trichocarpa*. *New Phytol*. 201:1263–1276. <https://doi.org/10.1111/nph.12601>
- McKown AD, Guy RD, Quamme L, Klápště J, Mantia JL, et al. 2014c. Association genetics, geography and ecophysiology link stomatal patterning in *Populus trichocarpa* with carbon gain and disease resistance trade-offs. *Mol Ecol*. 23:5771–5790. <https://doi.org/10.1111/mec.12969>
- McKown AD, Klápště J, Guy RD, Geraldes A, Porth I, et al. 2014b. Genome-wide association implicates numerous genes underlying ecological trait variation in natural populations of *Populus trichocarpa*. *New Phytol*. 203:535–553. <https://doi.org/10.1111/nph.12815>
- Nei M. 1973. Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci U S A*. 70:3321–3323. <https://doi.org/10.1073/pnas.70.12.3321>

- Nei M, Chesser RK. 1983. Estimation of fixation indices and gene diversities. *Ann Hum Genet.* 47:253–259. <https://doi.org/10.1111/j.1469-1809.1983.tb00993.x>
- Nicholson G, Smith AV, Jónsson F, Gústafsson Ó, Stefánsson K, et al. 2002. Assessing population differentiation and isolation from single-nucleotide polymorphism data. *J Royal Statistical Soc B.* 64: 695–715. <https://doi.org/10.1111/1467-9868.00357>
- Nielsen R, Akey JM, Jakobsson M, Pritchard JK, Tishkoff S, et al. 2017. Tracing the peopling of the world through genomics. *Nature.* 541: 302–310. <https://doi.org/10.1038/nature21347>
- Nosil P, Funk DJ, Ortiz-Barrientos D. 2009. Divergent selection and heterogeneous genomic divergence. *Mol Ecol.* 18:375–402. <https://doi.org/10.1111/j.1365-294X.2008.03946.x>
- Oleksyk TK, Zhao K, De La Vega FM, Gilbert DA, O'Brien SJ, et al. 2008. Identifying selected regions from heterozygosity and divergence using a light-coverage genomic dataset from two human populations. *PLoS One.* 3:e1712.
- Orsini L, Vanoverbeke J, Swillen I, Mergeay J, Meester LD. 2013. Drivers of population genetic differentiation in the wild: isolation by dispersal limitation, isolation by adaptation and isolation by colonization. *Mol Ecol.* 22:5983–5999. <https://doi.org/10.1111/mec.12561>
- Palsbøll PJ, Berube M, Allendorf FW. 2007. Identification of management units using population genetic data. *Trends Ecol Evol.* 22: 11–16. <https://doi.org/10.1016/j.tree.2006.09.003>
- Pérez-Lezaun A, Calafell F, Mateu E, Comas D, Ruiz-Pacheco R, et al. 1997. Microsatellite variation and the differentiation of modern humans. *Hum Genet.* 99:1–7. <https://doi.org/10.1007/s004390050299>
- Petrou EL, Seeb JE, Hauser L, Witteveen MJ, Templin WD, et al. 2014. Fine-scale sampling reveals distinct isolation by distance patterns in chum salmon (*Oncorhynchus keta*) populations occupying a glacially dynamic environment. *Conserv Genet.* 15:229–243. <https://doi.org/10.1007/s10592-013-0534-3>
- Pickrell J, Pritchard JK. 2012. Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genet.* 8: e1002967. <https://doi.org/10.1371/journal.pgen.1002967>
- Ramachandran S, Deshpande O, Roseman CC, Rosenberg NA, Feldman MW, et al. 2005. Support from the relationship of genetic and geographic distance in human populations for a serial founder effect originating in Africa. *Proc Natl Acad Sci U S A.* 102: 15942–15947. <https://doi.org/10.1073/pnas.0507611102>
- Raymond M, Rousset F. 1995. GENEPOP (version 1.2): population genetics software for exact tests and ecumenicism. *J Hered.* 86:248–249.
- Reynolds J, Weir BS, Cockerham CC. 1983. Estimation of the coancestry coefficient: basis for a short-term genetic distance. *Genetics.* 105:767–779. <https://doi.org/10.1093/genetics/105.3.767>
- Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, et al. 2002. Genetic structure of human populations. *Science.* 298:2381–2385. <https://doi.org/10.1126/science.1078311>
- Rougemont Q, Moore JS, Leroy T, Normandeau E, Rondeau EB, et al. 2020. Demographic history shaped geographical patterns of deleterious mutation load in a broadly distributed Pacific Salmon. *PLoS Genet.* 16:e1008348. <https://doi.org/10.1371/journal.pgen.1008348>
- Rousset F. 1997. Genetic differentiation and estimation of gene flow from F-statistics under isolation by distance. *Genetics.* 145: 1219–1228.
- Rousset F. 2004. *Genetic Structure and Selection in Subdivided Populations.* Princeton, NJ: Princeton University Press.
- Rousset F. 2007. Inferences from spatial population genetics. In: DJ Balding, M Bishop, C Cannings, editors. *Handbook of Statistical Genetics.* 3rd ed. Chichester, UK: Wiley. p. 945–979.
- Rousset F. 2008. Genepop'007: a complete reimplement of the Genepop software for Windows and Linux. *Mol Ecol Resour.* 8: 103–106. <https://doi.org/10.1111/j.1471-8286.2007.01931.x>
- Rutherford A. 2016. *A Brief History of Everyone Who Ever Lived: The Human Story Retold through Our Genes.* New York, NY: The Experiment.
- Saitou N, Nei M. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol.* 4:406–425. <https://doi.org/10.1093/oxfordjournals.molbev.a040454>
- Schlötterer C. 2004. The evolution of molecular markers—just a matter of fashion? *Nat Rev Genet.* 5:63–69. <https://doi.org/10.1038/nrg1249>
- Schwartz MK, Luikart G, Waples RS. 2007. Genetic monitoring as a promising tool for conservation and management. *Trends Ecol Evol.* 22:25–33. <https://doi.org/10.1016/j.tree.2006.08.009>
- Slatkin M. 1987. Gene flow and the geographic structure of natural populations. *Science.* 236:787–792. <https://doi.org/10.1126/science.3576198>
- Stein C. 1956. Inadmissibility of the usual estimator for the mean of a multivariate distribution. *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability.* Vol. 1. Berkeley, CA: University of California Press. p. 197–206.
- Therkildsen NO, Hemmer-Hansen J, Hedeholm RB, Wisz MS, Pampoulie C, et al. 2013a. Spatiotemporal SNP analysis reveals pronounced biocomplexity at the northern range margin of Atlantic cod *Gadus morhua*. *Evol Appl.* 6:690–705. <https://doi.org/10.1111/eva.12055>
- Therkildsen NO, Hemmer-Hansen J, Hedeholm RB, Wisz MS, Pampoulie C, et al. 2013b. Data from: spatiotemporal SNP analysis reveals pronounced biocomplexity at the northern range margin of Atlantic cod *Gadus morhua*, v2, Dryad. Dataset. <https://doi.org/10.5061/dryad.rd250>
- Waples RS, Gaggiotti O. 2006. What is a population? An empirical evaluation of some genetic methods for identifying the number of gene pools and their degree of connectivity. *Mol Ecol.* 15: 1419–1439. <https://doi.org/10.1111/j.1365-294X.2006.02890.x>
- Weir BS, Cardon LR, Anderson AD, Nielsen DM, Hill WG. 2005. Measures of human population structure show heterogeneity among genomic regions. *Genome Res.* 15:1468–1476. <https://doi.org/10.1101/gr.4398405>
- Weir BS, Cockerham CC. 1984. Estimating F-statistics for the analysis of population structure. *Evolution.* 38:1358–1370.
- Weir BS, Goudet J. 2017. A unified characterization of population structure and relatedness. *Genetics.* 206:2085–2103. <https://doi.org/10.1534/genetics.116.198424>
- Weir BS, Hill WG. 2002. Estimating F-statistics. *Annu Rev Genet.* 36: 721–750. <https://doi.org/10.1146/annurev.genet.36.050802.093940>
- Whitlock MC, McCauley DE. 1999. Indirect measures of gene flow and migration: $F_{ST} \neq 1/(4Nm + 1)$. *Heredity.* 82: 117–125. <https://doi.org/10.1046/j.1365-2540.1999.00496.x>
- Wright S. 1931. Evolution in Mendelian populations. *Genetics.* 16: 97–158. PMID: 17246615
- Wright S. 1951. The genetical structure of populations. *Ann Eugen.* 15:323–354. <https://doi.org/10.1111/j.1469-1809.1949.tb02451.x>
- Wright S. 1965. The interpretation of population structure by F-statistics with special regard to systems of mating. *Evolution.* 19:395–420. <https://doi.org/10.2307/2406450>
- Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, et al. 2010. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science.* 329:75–78. <https://doi.org/10.1126/science.1190371>