

*Submitted to the Annals of Applied Statistics*

1.

1     **FAST INFERENCE OF INDIVIDUAL ADMIXTURE**  
2     **COEFFICIENTS USING GEOGRAPHIC DATA**

3     BY KEVIN CAYE\*, FLORA JAY†, OLIVIER MICHEL\*, AND OLIVIER  
4     FRANÇOIS\*

5     *Université Grenoble-Alpes\* and Université Paris Diderot†*

6             Accurately evaluating the distribution of genetic ancestry  
7             across geographic space is one of the main questions addressed  
8             by evolutionary biologists. This question has been commonly  
9             addressed through the application of Bayesian estimation pro-  
10            grams allowing their users to estimate individual admixture pro-  
11            portions and allele frequencies among putative ancestral pop-  
12            ulations. Following the explosion of high-throughput sequenc-  
13            ing technologies, several algorithms have been proposed to cope  
14            with computational burden generated by the massive data in  
15            those studies. In this context, incorporating geographic prox-  
16            imity in ancestry estimation algorithms is an open statistical  
17            and computational challenge. In this study, we introduce new  
18            algorithms that use geographic information to estimate ances-  
19            try proportions and ancestral genotype frequencies from pop-  
20            ulation genetic data. Our algorithms combine matrix factor-  
21            ization methods and spatial statistics to provide estimates of  
22            ancestry matrices based on least-squares approximation. We  
23            demonstrate the benefit of using spatial algorithms through ex-  
24            tensive computer simulations, and we provide an example of  
25            application of our new algorithms to a set of spatially refer-  
26            enced samples for the plant species *Arabidopsis thaliana*. With-  
27            out loss of statistical accuracy, the new algorithms exhibit run-  
28            times that are much shorter than those observed for previously  
29            developed spatial methods. Our algorithms are implemented  
30            in the R package, `tess3r`, which is available from [https://](https://github.com/BioShock38/TESS3_encho_sen)  
31            [github.com/BioShock38/TESS3\\_encho\\_sen](https://github.com/BioShock38/TESS3_encho_sen).

32     **1. Introduction.** High-throughput sequencing technologies have  
33     enabled studies of genetic ancestry for model and non-model species  
34     at an unprecedented pace. In this context, ancestry estimation algo-

---

*Keywords and phrases:* Ancestry Estimation Algorithms, Genotypic Data, Geo-  
graphic Data, Fast Algorithms  
imsart-aos ver. 2014/10/16 file: draft.tex date: October 7, 2016

35 rithms are important for demographic analysis, medical genetics, con-  
36 servation and landscape genetics (Pritchard, Stephens and Donnelly,  
37 2000; Tang et al., 2005; Schraiber and Akey, 2015; Segelbacher et al.,  
38 2010; François and Waits, 2016). With increasingly large data sets,  
39 Bayesian approaches to the inference of population structure, exem-  
40 plified by the computer program `structure` (Pritchard, Stephens and  
41 Donnelly, 2000), have been replaced by approximate algorithms that  
42 run several orders faster than the original version (Tang et al., 2005;  
43 Alexander and Lange, 2011; Frichot et al., 2014; Raj, Stephens and  
44 Pritchard, 2014). Considering  $K$  ancestral populations or genetic clus-  
45 ters, those algorithms estimate ancestry coefficients following two main  
46 directions: model-based and model-free approaches. In model-based ap-  
47 proaches, a likelihood function is defined for the matrix of ancestry  
48 coefficients, and estimation is performed by maximizing the logarithm  
49 of the likelihood function. For `structure` and derived models, model  
50 assumptions include linkage equilibrium and Hardy-Weinberg equilib-  
51 rium in ancestral populations. The first approximation to the original  
52 algorithm was based on an expectation-minimization algorithm (Tang  
53 et al., 2005), and more recent likelihood algorithms are implemented in  
54 the programs `admixture` and `faststructure` (Alexander and Lange,  
55 2011; Raj, Stephens and Pritchard, 2014). In model-free approaches,  
56 ancestry coefficients are estimated by using least-squares methods or  
57 factor analysis. Model-free methods make no assumptions about the bi-  
58 ological processes that have generated the data. To estimate ancestry

59 matrices, [Engelhardt and Stephens \(2010\)](#) proposed to use sparse fac-  
60 tor analysis, [Frichot et al. \(2014\)](#) used sparse non-negative matrix fac-  
61 torization algorithms, and [Popescu et al. \(2014\)](#) used kernel-principal  
62 component analysis. Least-squares methods accurately reproduce the  
63 results of likelihood approaches under the model assumptions of those  
64 methods ([Frichot et al., 2014](#); [Popescu et al., 2014](#)). In addition, model-  
65 free methods provide approaches that are valid when the assumptions  
66 of likelihood approaches are not met. Model-free methods are generally  
67 faster than model-based methods.

68 Among model-based approaches to ancestry estimation, an impor-  
69 tant class of methods have improved the Bayesian model of `structure`  
70 by incorporating geographic data through spatially informative prior  
71 distributions ([Chen et al., 2007](#); [Corander, Sirén and Arjas, 2008](#)). Un-  
72 der isolation-by-distance patterns ([Wright, 1943](#); [Malécot, 1948](#)), spa-  
73 tial algorithms provide more robust estimates of population structure  
74 than non-spatial algorithms which can lead to biased estimates of the  
75 number of clusters ([Durand et al., 2009](#)). Some Bayesian methods are  
76 based on Markov chain Monte Carlo algorithms which are computer-  
77 intensive ([François and Durand, 2010](#)). Recent efforts to improve the in-  
78 ference of ancestral relationships in a geographical context have mainly  
79 focused on the localization of recent ancestors ([Baran et al., 2013](#); [Lao](#)  
80 [et al., 2014](#); [Yang et al., 2014](#)). In these applications, spatial informa-  
81 tion is used in a predictive framework that assigns ancestors to putative  
82 geographic origins. While fast geographic estimation of individual an-

83 cetry proportions has been proposed previously (Caye et al., 2016),  
84 there is a growing need to develop individual ancestry estimation al-  
85 gorithms that reduce computational cost in a geographically explicit  
86 framework.

87 In this study, we present two new algorithms for the estimation of  
88 ancestry matrices based on geographic and genetic data. The new al-  
89 gorithms solve a least squares optimization problem as defined by Caye  
90 et al. (2016), based on Alternating Quadratic Programming (AQP) and  
91 Alternating Projected Least Squares (APLS). While AQP algorithms  
92 have a well-established theoretical background (Bertsekas, 1995), this  
93 is not the case of APLS algorithms. Using coalescent simulations, we  
94 provide evidence that the estimates computed by APLS algorithms  
95 are good approximations to the solutions of AQP algorithms. In ad-  
96 dition, we show that the performances of APLS algorithms scale with  
97 the dimensions of modern data sets. We discuss the application of our  
98 algorithms to data from European ecotypes of *Arabidopsis thaliana*,  
99 for which individual genomic and geographic data are available (Horton  
100 et al., 2012).

101 **2. New methods.** In this section we present two new algorithms  
102 for estimating individual admixture coefficients and ancestral genotype  
103 frequencies assuming  $K$  ancestral populations. In addition to geno-  
104 types, the new algorithms require individual geographic coordinates of  
105 sampled individuals.

106 *Q and G-matrices.* Consider a genotypic matrix,  $\mathbf{Y}$ , recording data  
107 for  $n$  individuals at  $L$  polymorphic loci for a  $p$ -ploid species (common  
108 values for  $p$  are  $p = 1, 2$ ). For autosomal SNPs in a diploid organism, the  
109 genotype at locus  $\ell$  is an integer number, 0, 1 or 2, corresponding to the  
110 number of reference alleles at this locus. In our algorithms, disjunctive  
111 forms are used to encode each genotypic value as the indicator of a  
112 heterozygote or a homozygote locus (Frichot et al. 2014). For a diploid  
113 organism each genotypic value, 0, 1, 2 is encoded as 100, 010 and 001.  
114 For  $p$ -ploid organisms, there are  $(p + 1)$  possible genotypic values at  
115 each locus, and each value corresponds to a unique disjunctive form.  
116 While our focus is on SNPs, the algorithms presented in this section  
117 extend to multi-allelic loci without loss of generality. Moreover, the  
118 method can be easily extended to genotype likelihoods by using the  
119 likelihood to encode each genotypic value (Korneliussen, Albrechtsen  
120 and Nielsen, 2014).

121 Our algorithms provide statistical estimates for the matrix  $\mathbf{Q} \in$   
122  $\mathbb{R}^{K \times n}$  which contains the admixture coefficients,  $\mathbf{Q}_{i,k}$ , for each sam-  
123 pled individual,  $i$ , and each ancestral population,  $k$ . The algorithms  
124 also provide estimates for the matrix  $\mathbf{G} \in \mathbb{R}^{(p+1)L \times K}$ , for which the

125 entries,  $\mathbf{G}_{(p+1)\ell+j,k}$ , correspond to the frequency of genotype  $j$  at locus  
126  $\ell$  in population  $k$ . Obviously, the  $\mathbf{Q}$  and  $\mathbf{G}$ -matrices must satisfy the  
127 following set of probabilistic constraints

$$\mathbf{Q}, \mathbf{G} \geq 0, \quad \sum_{k=1}^K \mathbf{Q}_{i,k} = 1, \quad \sum_{j=0}^p \mathbf{G}_{(p+1)\ell+j,k} = 1, \quad j = 0, 1, \dots, p,$$

128 for all  $i, k$  and  $\ell$ . Using disjunctive forms and the law of total probabil-  
129 ity, estimates of  $\mathbf{Q}$  and  $\mathbf{G}$  can be obtained by factorizing the genotypic  
130 matrix as follows  $\mathbf{Y} = \mathbf{Q} \mathbf{G}^T$  (Frichot et al., 2014). Thus the inference  
131 problem can be solved by using constrained nonnegative matrix factor-  
132 ization methods (Lee and Seung, 1999; Cichocki et al., 2009). In the  
133 sequel, we shall use the notations  $\Delta_Q$  and  $\Delta_G$  to represent the sets of  
134 probabilistic constraints put on the  $\mathbf{Q}$  and  $\mathbf{G}$  matrices respectively.

135 *Geographic weighting.* Geography is introduced in the matrix factor-  
136 ization problem by using weights for each pair of sampled individuals.  
137 The weights impose regularity constraints on ancestry estimates over  
138 geographic space. The definition of geographic weights is based on the  
139 spatial coordinates of the sampling sites,  $(x_i)$ . Samples close to each  
140 other are given more weight than samples that are far apart. The com-  
141 putation of the weights starts with building a complete graph from the  
142 sampling sites. Then the weight matrix is defined as follows

$$w_{ij} = \exp(-\text{dist}(x_i, x_j)^2 / \sigma^2),$$

143 where  $\text{dist}(x_i, x_j)$  denotes the geodesic distance between sites  $x_i$  and

144  $x_j$ , and  $\sigma$  is a range parameter. Values for the range parameter can be  
145 investigated by using spatial variograms (Cressie, 1993). To evaluate  
146 variograms, we extend the univariate variogram to genotypic data as  
147 follows

$$(2.1) \quad \gamma(h) = \frac{1}{2|N(h)|} \sum_{i,j \in N(h)} \frac{1}{L} \sum_{l=1}^{(p+1)L} |Y_{i,l} - Y_{j,l}|,$$

148 where  $N(h)$  is defined as the set of individuals separated by geographic  
149 distance  $h$ . In applications, computing and visualizing the  $\gamma$  function  
150 provides useful information on the level of spatial autocorrelation be-  
151 tween individuals in the data.

152 Next, we introduce the *Laplacian matrix* associated with the geo-  
153 graphic weight matrix,  $\mathbf{W}$ . The Laplacian matrix is defined as  $\mathbf{\Lambda} =$   
154  $\mathbf{D} - \mathbf{W}$  where  $\mathbf{D}$  is a diagonal matrix with entries  $\mathbf{D}_{i,i} = \sum_{j=1}^n \mathbf{W}_{i,j}$ ,  
155 for  $i = 1, \dots, n$  (Belkin and Niyogi, 2003). Elementary matrix algebra  
156 shows that (Cai et al., 2011)

$$\text{Tr}(\mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q}) = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \|\mathbf{Q}_{i,\cdot} - \mathbf{Q}_{j,\cdot}\|^2.$$

157 In our approach, assuming that geographically close individuals are  
158 more likely to share ancestry than individuals at distant sites is thus  
159 equivalent to minimizing the quadratic form  $\mathcal{C}(\mathbf{Q}) = \text{Tr}(\mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q})$  while  
160 estimating the matrix  $\mathbf{Q}$ .

161 *Least-squares optimization problems.* Estimating the matrices  $\mathbf{Q}$  and  
162  $\mathbf{G}$  from the observed genotypic matrix  $\mathbf{Y}$  is performed through solving

163 an optimization problem defined as follows (Caye et al., 2016)

$$(2.2) \quad \begin{aligned} \min_{\mathbf{Q}, \mathbf{G}} \quad & \text{LS}(\mathbf{Q}, \mathbf{G}) = \|\mathbf{Y} - \mathbf{Q}\mathbf{G}^T\|_{\text{F}}^2 + \alpha' \frac{(p+1)L}{K\lambda_{\max}} \mathcal{C}(\mathbf{Q}), \\ \text{s.t.} \quad & \mathbf{Q} \in \Delta_Q, \\ & \mathbf{G} \in \Delta_G. \end{aligned}$$

164 The notation  $\|\mathbf{M}\|_{\text{F}}$  denotes the Frobenius norm of a matrix,  $\mathbf{M}$ . The  
165 regularization term is normalized by  $(p+1)L/K\lambda_{\max}$ , where  $\lambda_{\max}$  is the  
166 largest eigenvalue of the Laplacian matrix. With this normalization,  
167 both terms of the optimization problem (2.2) are given the same order  
168 of magnitude. The regularization parameter  $\alpha'$  controls the regularity  
169 of ancestry estimates over geographic space. Large values of  $\alpha'$  imply  
170 that ancestry coefficients have similar values for nearby individuals,  
171 whereas small values ignore spatial autocorrelation in observed allele  
172 frequencies. In the rest of the article, we will use  $\alpha' = 1$  and  $\alpha = (p +$   
173  $1)L/K\lambda_{\max}$ . Using the least-squares approach, the number of ancestral  
174 populations,  $K$ , can be chosen after the evaluation of a cross-validation  
175 criterion for each  $K$  (Alexander and Lange, 2011; Frichot et al., 2014;  
176 Frichot and François, 2015).

177 *The Alternating Quadratic Programming (AQP) method.* Because the  
178 polyhedrons  $\Delta_Q$  and  $\Delta_G$  are convex sets and the LS function is convex  
179 with respect to each variable  $\mathbf{Q}$  or  $\mathbf{G}$  when the other one is fixed,  
180 the problem (2.2) is amenable to the application of block coordinate  
181 descent (Bertsekas, 1995). The APQ algorithm starts from initial values  
182 for the  $G$  and  $Q$ -matrices, and alternates two steps. The first step



183 computes the matrix  $\mathbf{G}$  while  $\mathbf{Q}$  is kept fixed, and the second step  
 184 permutes the roles of  $\mathbf{G}$  and  $\mathbf{Q}$ . Let us assume that  $\mathbf{Q}$  is fixed and  
 185 write  $\mathbf{G}$  in a vectorial form,  $g = \text{vec}(\mathbf{G}) \in \mathbb{R}^{K(p+1)L}$ . The first step  
 186 of the algorithm actually solves the following quadratic programming  
 187 subproblem. Find

$$(2.3) \quad g^* = \arg \min_{g \in \Delta_G} (-2v_Q^T g + g^T \mathbf{D}_Q g),$$

188 where  $\mathbf{D}_Q = \mathbf{I}_{(p+1)L} \otimes \mathbf{Q}^T \mathbf{Q}$  and  $v_Q = \text{vec}(\mathbf{Q}^T \mathbf{Y})$ . Here,  $\otimes$  denotes the  
 189 Kronecker product and  $\mathbf{I}_d$  is the identity matrix with  $d$  dimensions.  
 190 Note that the block structure of the matrix  $\mathbf{D}_Q$  allows us to decom-  
 191 pose the subproblem (2.3) into  $L$  independent quadratic programming  
 192 problems with  $K(p+1)$  variables. Now, consider that  $\mathbf{G}$  is the value  
 193 obtained after the first step of the algorithm, and write  $\mathbf{Q}$  in a vec-  
 194 torial form,  $q = \text{vec}(\mathbf{Q}) \in \mathbb{R}^{nK}$ . The second step solves the following  
 195 quadratic programming subproblem. Find

$$(2.4) \quad q^* = \arg \min_{q \in \Delta_Q} (-2v_G^T q + q^T \mathbf{D}_G q),$$

196 where  $\mathbf{D}_G = \mathbf{I}_n \otimes \mathbf{G}^T \mathbf{G} + \alpha \mathbf{\Lambda} \otimes \mathbf{I}_K$  and  $v_G = \text{vec}(\mathbf{G}^T \mathbf{Y}^T)$ . Unlike sub-  
 197 problem (2.3), subproblem (2.4) can not be decomposed into smaller  
 198 problems. Thus, the computation of the second step of the AQP al-  
 199 gorithm implies to solve a quadratic programming problem with  $nK$   
 200 variables which can be problematic for large samples ( $n$  is the sample

201 size). The AQP algorithm is described in details in Appendix [A.1](#). For  
202 AQP, we have the following convergence result.

203 **THEOREM 2.1.** *The AQP algorithm converges to a critical point of*  
204 *problem (2.2).*

205 **PROOF.** The quadratic convex functions defined in subproblems (2.3)  
206 and (2.4) have finite lower bounds. The convex sets  $\Delta_Q$  and  $\Delta_G$  are not  
207 empty sets, and they are compact sets. Thus the sequence generated  
208 by the AQP algorithm is well-defined, and has limit points. According  
209 to Corollary 2 of [Grippo and Sciandrone \(2000\)](#), we conclude that the  
210 AQP algorithm converges to a critical point of problem (2.2).

211 *Alternating Projected Least-Squares (APLS).* In this paragraph, we  
212 introduce an APLS estimation algorithm which approximates the so-  
213 lution of problem (2.2), and reduces the complexity of the AQP al-  
214 gorithm. The APLS algorithm starts from initial values of the  $G$  and  
215  $Q$ -matrices, and alternates two steps. The matrix  $\mathbf{G}$  is computed while  
216  $\mathbf{Q}$  is kept fixed, and *vice versa*. Assume that the matrix  $\mathbf{Q}$  is known.  
217 The first step of the APLS algorithm solves the following optimization  
218 problem. Find

$$(2.5) \quad \mathbf{G}^* = \arg \min \|\mathbf{Y} - \mathbf{Q}\mathbf{G}^T\|_{\mathbb{F}}^2.$$

219 This operation can be done by considering  $(p + 1)L$  (the number of  
220 columns of  $\mathbf{Y}$ ) independent optimization problems running in parallel.

221 The operation is followed by a projection of  $\mathbf{G}^*$  on the polyedron of  
222 constraints,  $\Delta_G$ . For the second step, assume that  $\mathbf{G}$  is set to the value  
223 obtained after the first step is completed. We compute the eigenvec-  
224 tors,  $\mathbf{U}$ , of the Laplacian matrix, and we define the diagonal matrix  $\mathbf{\Delta}$   
225 formed by the eigenvalues of  $\mathbf{\Lambda}$  (The eigenvalues of  $\mathbf{\Lambda}$  are non-negative  
226 real numbers). According to the spectral theorem, we have

$$\mathbf{\Lambda} = \mathbf{U}^T \mathbf{\Delta} \mathbf{U}.$$

227 After this operation, we project the data matrix  $\mathbf{Y}$  on the basis of  
228 eigenvectors as follows

$$\text{proj}(\mathbf{Y}) = \mathbf{U} \mathbf{Y},$$

229 and, for each individual, we solve the following optimization problem

$$(2.6) \quad q_i^* = \arg \min \| \text{proj}(\mathbf{Y})_i - \mathbf{G}^T q \|^2 + \alpha \lambda_i \| q \|^2,$$

230 where  $\text{proj}(\mathbf{Y})_i$  is the  $i$ th row of the projected data matrix,  $\text{proj}(\mathbf{Y})$ ,  
231 and  $\lambda_i$  is the  $i$ th eigenvalue of  $\mathbf{\Lambda}$ . The solutions,  $q_i$ , are then concate-  
232 nated into a matrix,  $\text{conc}(q)$ , and  $\mathbf{Q}$  is defined as the projection of the  
233 matrix  $\mathbf{U}^T \text{conc}(q)$  on the polyedron  $\Delta_Q$ . The complexity of step (2.6)  
234 grows linearly with  $n$ , the number of individuals. While the theoretic-  
235 al convergence properties of AQP algorithms are lost for APLS algo-  
236 rithms, the APLS algorithms are expected to be good approximations

237 of AQP algorithms. The APLS algorithm is described in details in Ap-  
238 pendix [A.2](#).

239 *Comparison with tess3.* The algorithm implemented in a previous  
240 version of `tess3` also provides approximation of of solution of [\(2.2\)](#).  
241 The `tess3` algorithm first computes a Cholesky decomposition of the  
242 Laplacian matrix. Then, by a change of variables, the least-squares  
243 problem is transformed into a sparse nonnegative matrix factorization  
244 problem ([Caye et al., 2016](#)). Solving the sparse non-negative matrix fac-  
245 torization problem relies on the application of existing methods ([Kim  
246 and Park, 2011](#); [Frichot et al., 2014](#)). The methods implemented in  
247 `tess3` have an algorithmic complexity that increases linearly with the  
248 number of loci and the number of clusters. They lead to estimates that  
249 accurately reproduce those of the Monte Carlo algorithms implemented  
250 in the Bayesian method `tess 2.3` ([Caye et al., 2016](#)). Like for the AQP  
251 method, the `tess3` previous algorithms have an algorithmic complexity  
252 that increases quadratically with the sample size.

253 *Ancestral population differentiation statistics and local adaptation scans.*  
254 Assuming  $K$  ancestral populations, the  $Q$  and  $G$ -matrices obtained  
255 from the AQP and from the APLS algorithms were used to compute  
256 single-locus estimates of a population differentiation statistic similar to  
257  $F_{ST}$  ([Martins et al., 2016](#)), as follows

$$F_{ST}^Q = 1 - \sum_{k=1}^K q_k \frac{f_k(1 - f_k)}{f(1 - f)},$$

258 where  $q_k$  is the average of ancestry coefficients over sampled individuals,  
259  $q_k = \sum_{i=1}^n q_{ik}/n$ , for the cluster  $k$ ,  $f_k$  is the ancestral allele frequency in  
260 population  $k$  at the locus of interest, and  $f = \sum_{k=1}^K q_k f_k$  (Martins et al.  
261 2016). The locus-specific statistics were used to perform statistical tests  
262 of neutrality at each locus, by comparing the observed values to their  
263 expectations from the genome-wide background. The test was based  
264 on the squared  $z$ -score statistic,  $z^2 = (n - K)F_{ST}^Q/(1 - F_{ST}^Q)$ , for which  
265 a chi-squared distribution with  $K - 1$  degrees of freedom was assumed  
266 under the null-hypothesis (Martins et al., 2016). The calibration of  
267 the null-hypothesis was achieved by using genomic control to adjust  
268 the test statistic for background levels of population structure (Devlin  
269 and Roeder, 1999; François et al., 2016). After recalibration of the null-  
270 hypothesis, the control of the false discovery rate was achieved by using  
271 the Benjamini-Hochberg algorithm (Benjamini and Hochberg, 1995).

272 *R package.* We implemented the AQP and APLS algorithms in the R  
273 package `tess3r`, available from Github and submitted to the Compre-  
274 hensive R Archive Network (R Core Team, 2016).

### 275 3. Simulated and real data sets.

276 *Coalescent simulations.* We used the computer program `ms` to per-  
277 form coalescent simulations of neutral and outlier SNPs under spatial  
278 models of admixture (Hudson, 2002). Two ancestral populations were  
279 created from the simulation of Wright’s two-island models. The sim-  
280 ulated data sets contained admixed genotypes for  $n$  individuals for

281 which the admixture proportions varied continuously along a longitu-  
282 dinal gradient (Durand et al., 2009; François and Durand, 2010). In  
283 those scenarios, individuals at each extreme of the geographic range  
284 were representative of their population of origin, while individuals at  
285 the center of the range shared intermediate levels of ancestry in the two  
286 ancestral populations (Caye et al., 2016). For those simulations, the  $Q$   
287 matrix,  $\mathbf{Q}_0$ , was entirely described by the location of the sampled in-  
288 dividuals.

289 Neutrally evolving ancestral chromosomal segments were generated  
290 by simulating DNA sequences with an effective population size  $N_0 =$   
291  $10^6$  for each ancestral population. The mutation rate per bp and gener-  
292 ation was set to  $\mu = 0.25 \times 10^{-7}$ , the recombination rate per generation  
293 was set to  $r = 0.25 \times 10^{-8}$ , and the parameter  $m$  was set to obtained  
294 neutral levels of  $F_{ST}$  ranging between values of 0.005 and 0.10. The  
295 number of base pairs for each DNA sequence was varied between 10k  
296 to 300k to obtain numbers of polymorphic locus ranging between 1k  
297 and 200k after filtering out SNPs with minor allele frequency lower than  
298 5%. To create SNPs with values in the tail of the empirical distribution  
299 of  $F_{ST}$ , additional ancestral chromosomal segments were generated by  
300 simulating DNA sequences with a migration rate  $m_s$  lower than  $m$ .  
301 The simulations reproduced the reduced levels of diversity and the in-  
302 creased levels of differentiation expected under hard selective sweeps  
303 occurring at one particular chromosomal segment in ancestral popula-  
304 tions (Martins et al., 2016). For each simulation, the sample size was

305 varied in the range  $n = 50-700$ .

306 We compared the AQP and APLS algorithm estimates with those ob-  
307 tained with the `tess3` algorithm. Each program was run 5 times. Using  
308  $K = 2$  ancestral populations, we computed the root mean squared error  
309 (RMSE) between the estimated and known values of the  $Q$ -matrix, and  
310 between the estimated and known values of the  $G$ -matrix. To evaluate  
311 the benefit of spatial algorithms, we compared the statistical errors of  
312 APLS algorithms to the errors obtained with `snmf` method that re-  
313 produces the outputs of the `structure` program accurately (Frichot  
314 et al., 2014; Frichot and François, 2015). To quantify the performances  
315 of neutrality tests as a function of ancestral and observed levels of  $F_{ST}$ ,  
316 we used the area under the precision-recall curve (AUC) for several  
317 values of the selection rate. Subsamples from a real data set were used  
318 to perform a runtime analysis of the AQP and APLS algorithms (*A.*  
319 *thaliana* data, see below). Runtimes were evaluated by using a single  
320 computer processor unit Intel Xeon 2.0 GHz.

321 *Application to European ecotypes of Arabidopsis thaliana.* We used  
322 the APLS algorithm to survey spatial population genetic structure and  
323 to investigate the molecular basis of adaptation by considering SNP  
324 data from 1,095 European ecotypes of the plant species *A. thaliana*  
325 (214k SNPs, Horton et al. (2012)). The cross-validation criterion was  
326 used to evaluate the number of clusters in the sample, and a statis-  
327 tical analysis was performed to evaluate the range of the variogram  
328 from the data. We used R functions of the `tess3r` package to display

329 interpolated admixture coefficients on a geographic map of Europe (R  
330 Core team 2016). A gene ontology enrichment analysis using the soft-  
331 ware AMIGO (Carbon et al., 2009) was performed in order to evaluate  
332 which molecular functions and biological processes might be involved  
333 in local adaptation in Europe.



334 **4. Results.**

335 *Statistical errors.* We used coalescent simulations of neutral polymor-  
336 phisms under spatial models of admixture to compare the statistical  
337 errors of the AQP and APLS estimates with those of the `tess3` al-  
338 gorithm. The ground truth for the  $Q$ -matrix ( $\mathbf{Q}_0$ ) was computed from  
339 the mathematical model for admixture proportions used to generate the  
340 data. For the  $G$ -matrix, the ground truth matrix ( $\mathbf{G}_0$ ) was computed  
341 from the empirical genotype frequencies in the two population samples  
342 before an admixture event. The root mean squared errors (RMSE) for  
343 the  $\mathbf{Q}$  and  $\mathbf{G}$  estimates decreased as the sample size and the number of  
344 loci increased (Figure 1). For all algorithms, the statistical errors were  
345 generally small when the number of loci was greater than 10k SNPs.  
346 Those results provided evidence that the three algorithms produced  
347 equivalent estimates of the matrices  $\mathbf{Q}_0$  and  $\mathbf{G}_0$ . The results also pro-  
348 vided a formal check that the APLS and `tess3` algorithms converged  
349 to the same estimates as those obtained after the application of the  
350 AQP algorithm, which is guaranteed to converge mathematically.

351 *The benefit of including spatial information in algorithms.* Using neu-  
352 tral coalescent simulations of spatial admixture, we compared the sta-  
353 tistical estimates obtained from a spatial algorithm (APLS) and a non-  
354 spatial algorithm (sNMF, Frichot et al. 2014). For various levels of an-  
355 cestral population differentiation, estimates obtained from the spatial  
356 algorithm were more accurate than for those obtained using non-spatial  
357 approaches (Figure 2). For the larger samples, much finer population

358 structure was detected with the spatial method than with the non-  
359 spatial algorithm (Figure 2).

360 In simulations of outlier loci, we used the area under the precision-  
361 recall curve (AUC) for quantifying the performances of tests based on  
362 the estimates of ancestry matrices,  $\mathbf{Q}$  and  $\mathbf{G}$ . In addition, we computed  
363 AUCs for  $F_{ST}$ -based neutrality tests using truly ancestral genotypes. As  
364 they represented the maximum reachable values, AUCs based on truly  
365 ancestral genotypes were always higher than those obtained for tests  
366 based on reconstructed matrices. For all values of the relative selection  
367 intensity, AUCs were higher for spatial methods than for non-spatial  
368 methods (Figure 3, the relative selection intensity is the ratio of migra-  
369 tion rates at neutral and adaptive loci). For high selection intensities,  
370 the performances of tests based on estimates of ancestry matrices were  
371 close to the optimal values reached by tests based on true ancestral  
372 frequencies. These results provided evidence that including spatial in-  
373 formation in ancestry estimation algorithms improves the detection of  
374 signatures of hard selective sweeps having occurred in unknown ances-  
375 tral populations.

376 *Runtime and convergence analyses.* We subsampled a large SNP data  
377 set for *A. thaliana* ecotypes to compare the convergence properties and  
378 runtimes of the `tess3`, AQP, and APLS algorithms. In those exper-  
379 iments, we used  $K = 6$  ancestral populations, and replicated 5 runs  
380 for each simulation. For  $n = 100 - 600$  individuals ( $L = 50k$  SNPs),  
381 the APLS algorithm required more iterations (25 iterations) than the

382 AQP algorithm (20 iterations) to converge to its solution (Figure 4).  
383 This was less than for `tess3` (30 iterations). For  $L = 10 - 200\text{k}$  SNPs  
384 ( $n = 150$  individuals), similar results were observed. For 50k SNPs, the  
385 runtimes were significantly lower for the APLS algorithm than for the  
386 `tess3` and AQP algorithms. For  $L = 50\text{k}$  SNPs and  $n = 600$  individ-  
387 uals, it took on average 0.956 min for the APLS and 100 min for the  
388 AQP algorithm to compute ancestry estimates. For `tess3`, the runtime  
389 was on average 66.3 min. For  $L = 100\text{k}$  SNPs and  $n = 150$  individuals,  
390 it took on average 0.628 min (8.97 min) for the APLS (AQP) algo-  
391 rithm to compute ancestry estimates. For `tess3`, the runtime was on  
392 average 1.27 min. For those values of  $n$  and  $L$ , the APLS algorithm im-  
393 plementation ran about 2 to 100 times faster than the other algorithm  
394 implementations.

395 *Application to European ecotypes of Arabidopsis thaliana.* We used  
396 the APLS algorithm to survey spatial population genetic structure and  
397 perform a genome scan for adaptive alleles in European ecotypes of  
398 the plant species *A. thaliana*. The cross validation criterion decreased  
399 rapidly from  $K = 1$  to  $K = 3$  clusters, indicating that there were three  
400 main ancestral groups in Europe, corresponding to geographic regions  
401 in Western Europe, Eastern and Central Europe and Northern Scan-  
402 dinavia. For  $K$  greater than four, the values of the cross validation  
403 criterion decreased in a slower way, indicating that subtle substruc-  
404 ture resulting from complex historical isolation-by-distance processes  
405 could also be detected (Figure 5). The spatial analysis provided an ap-

406 proximate range of  $\sigma = 150\text{km}$  for the spatial variogram (Figure 5).  
407 Figure 6 displays the  $Q$ -matrix estimate interpolated on a geographic  
408 map of Europe for  $K = 6$  ancestral groups. The estimated admixture  
409 coefficients provided clear evidence for the clustering of the ecotypes in  
410 spatially homogeneous genetic groups.

411 *Targets of selection in *A. thaliana* genomes.* Tests based on the  $F_{ST}^Q$   
412 statistic were applied to the 241k SNP data set to reveal new targets  
413 of natural selection in the *A. thaliana* genome. *A. thaliana* occurs in a  
414 broad variety of habitats, and local adaptation to the environment is  
415 acknowledged to be important in shaping its genetic diversity through  
416 space (Hancock et al., 2011; Fournier-Level et al., 2011). The APLS  
417 algorithm was run on the 1,095 European lines of *A. thaliana* with  
418  $K = 6$  ancestral populations and  $\sigma = 1.5$  for the range parameter.  
419 After controlling the FDR at the level 1%, the program produced a  
420 list of 12,701 candidate SNPs, including linked loci and representing  
421 3% of the total number of loci. The top 100 candidates included SNPs  
422 in the flowering-related genes SHORT VEGETATIVE PHASE (SVP),  
423 COP1-interacting protein 4.1 (CIP4.1) and FRIGIDA (FRI) ( $p$ -values  
424  $< 10^{-300}$ ). These genes were detected by previous scans for selection  
425 on this dataset (Horton et al., 2012). We performed a gene ontology  
426 enrichment analysis using AmiGO in order to evaluate which biological  
427 functions might be involved in local adaptation in Europe. We found  
428 a significant over-representation of genes involved in cellular processes  
429 (fold enrichment of 1.06,  $p$ -value equal to 0.0215 after Bonferonni cor-

FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 21

430 rection).

431 **5. Discussion.** Including geographic information on sample loca-  
432 tions in the inference of ancestral relationships among organisms is a  
433 major objective of population genetic studies (Malécot, 1948; Cavalli,  
434 Menozzi and Piazza, 1994; Epperson, 2003). Assuming that geographi-  
435 cally close individuals are more likely to share ancestry than individuals  
436 at distant sites, we introduced two new algorithms for estimating ances-  
437 try proportions using geographic information. Based on least-squares  
438 problems, the new algorithms combine matrix factorization approaches  
439 and spatial statistics to provide accurate estimates of individual ances-  
440 try coefficients and ancestral genotype frequencies. The two methods  
441 share many similarities, but they differ in the approximations they  
442 make in order to decrease algorithmic complexity. More specifically,  
443 the AQP algorithm was based on quadratic programming, whereas the  
444 APLS algorithm was based on the spectral decomposition of the Lapla-  
445 cian matrix. The algorithmic complexity of APLS algorithm grows lin-  
446 early with the number of individuals in the sample while the method  
447 has the same statistical accuracy as more complex algorithms.

448 To measure the benefit of using spatial algorithms, we compared the  
449 statistical errors observed for spatial algorithms with those observed  
450 for non-spatial algorithms. The errors of spatial methods were lower  
451 than those observed with non-spatial methods, and spatial algorithms  
452 allowed the detection of more subtle population structure. In addition,  
453 we implemented neutrality tests based on the spatial estimates of the  $Q$   
454 and  $G$ -matrices (Martins et al., 2016), and we observed that those tests

455 had higher power to reject neutrality than those based on non-spatial  
456 approaches. Thus spatial information helped improving the detection  
457 of signatures of selective sweeps having occurred in ancestral popu-  
458 lations prior to admixture events. We applied the neutrality tests to  
459 perform a genome scan for selection in European ecotypes of the plant  
460 species *A. thaliana*. The genome scan confirmed the evidence for selec-  
461 tion at flowering-related genes *CIP4.1*, *FRI* and *DOG1* differentiating  
462 Fennoscandia from North-West Europe (Horton et al., 2012).

463 Estimation of ancestry coefficients using fast algorithms that extend  
464 non-spatial approaches – such as `structure` – has been intensively  
465 discussed during the last years (Wollstein and Lao, 2015). In these im-  
466 provements, spatial approaches have received less attention than non-  
467 spatial approaches. In this study, we have proposed a conceptual frame-  
468 work for developing fast spatial ancestry estimation methods, and a  
469 suite of computer programs implements this framework in the R pro-  
470 gram `tess3r`. Our package provides an integrated pipeline for esti-  
471 mating and visualizing population genetic structure, and for scanning  
472 genomes for signature of local adaptation. The algorithmic complexity  
473 of our algorithms allow their users to analyze samples including hun-  
474 dreds to thousands of individuals. For example, analyzing more than  
475 one thousand *A. thaliana* genotypes, each including more than 210k  
476 SNPs, took less than a few minutes using a single CPU. In addition,  
477 the algorithms have multithreaded versions that run on parallel com-  
478 puters by using multiple CPUs. The multithreaded algorithm, which is

479 available from the R program, allows using our programs in large-scale  
480 genomic sequencing projects.

## APPENDIX A: ALGORITHMS

481 ALGORITHM A.1. AQP algorithm pseudo code. To solve optimiza-  
482 tion problem (2.2).

**Input:** the data matrix  $\mathbf{Y} \in \{0, 1\}^{n \times (p+1)L}$ , the Laplacian matrix  $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$ , the number of ancestral populations  $K$ , the regularization coefficient  $\alpha$ , the maximum number of iteration *itMax*

**Output:** the admixture coefficients matrix  $\mathbf{Q} \in \mathbb{R}^{n \times K}$ , the ancestral genotype frequencies matrix  $\mathbf{G} \in \mathbb{R}^{K \times (p+1)L}$

Initialize  $\mathbf{Q}$  at random ;

**for**  $it = 1..itMax$  **do**

    // G optimization step

**for**  $l = 1..L$  **do**

483      $Y^l \leftarrow \mathbf{Y}_{.,(p+1)l..(p+1)l+d}$  ;  
     $\mathbf{D}_Q \leftarrow \mathbf{I}_{p+1} \otimes \mathbf{Q}^T \mathbf{Q}$  ;  
     $v_Q \leftarrow \text{Vec}(\mathbf{Q}^T Y^l)$  ;  
     $g^* \in \arg \min_{g \in \Delta_G} -2v_Q^T g + g^T \mathbf{D}_Q g$  ;  
     $\text{Vec}(\mathbf{G}_{(p+1)l..(p+1)l+d,.}) \leftarrow g^*$  ;

**end**

    // Q optimization step

$\mathbf{D}_G \leftarrow Id_n \otimes \mathbf{G}^T \mathbf{G} + \alpha \mathbf{\Lambda} \otimes \mathbf{I}_K$  ;

$v_G \leftarrow \text{Vec}(\mathbf{G}^T \mathbf{Y}^T)$  ;

$\text{Vec}(\mathbf{Q}^T) \in \arg \min_{q \in \Delta_Q} -2v_G^T q + q^T \mathbf{D}_G q$  ;

**end**

484 ALGORITHM A.2. APLS algorithm pseudo code. To solve the op-



485 timization problem (2.2).

**Input:** the data matrix  $\mathbf{Y} \in \{0, 1\}^{n \times (d+1)L}$ , the eigen values and vectors matrices  $\mathbf{U}$  and  $\mathbf{\Delta}$  such that  $\mathbf{\Lambda} = \mathbf{U}^T \mathbf{\Delta} \mathbf{U}$ , the number of ancestral populations  $K$ , the regularization coefficient  $\alpha$ , the maximum number of iteration *itMax*

**Output:** the admixture coefficients matrix  $\mathbf{Q} \in \mathbb{R}^{n \times K}$ , the ancestral genotype frequencies matrix  $\mathbf{G} \in \mathbb{R}^{K \times (d+1)L}$

Initialize  $\mathbf{Q}$  at random ;

$proj(\mathbf{Y}) \leftarrow \mathbf{R} \mathbf{Y}$  ;

**for**  $it = 1..itMax$  **do**

    // G optimization step

**for**  $j = 1..(p+1)L$  **do**

$g^* \in \arg \min_{g \in \mathbb{R}^K} \|\mathbf{Y}_{:,j} - \mathbf{Q}g\|^2$ ;

$\mathbf{G}_{j,\cdot} \leftarrow g^*$ ;

**end**

    Project  $\mathbf{G}$  such that  $\mathbf{G} \in \Delta_G$  ;

    // Q optimization step

**for**  $i = 1..n$  **do**

$g_i^* \in \arg \min_{q \in \mathbb{R}^K} \|\text{proj}(\mathbf{Y})_{i,\cdot} - \mathbf{G}^T q\|^2 + \alpha \mathbf{\Delta}_{i,i} \|q\|^2$ ;

$proj(\mathbf{Q})_{i,\cdot} \leftarrow g_i^*$ ;

**end**

$\mathbf{Q} \leftarrow \mathbf{U}^T \text{proj}(\mathbf{Q})$ ;

    Project  $\mathbf{Q}$  such that  $\mathbf{Q} \in \Delta_Q$  ;

**end**

## ACKNOWLEDGEMENTS

487 This work has been partially supported by the LabEx PERSYVAL-  
488 Lab (ANR-11-LABX-0025-01) funded by the French program Investisse-  
489 ment d’Avenir. Olivier François acknowledges support from Grenoble  
490 INP and from the Agence Nationale de la Recherche, project AFRICROP  
491 ANR-13-BSV7-0017.

## REFERENCES

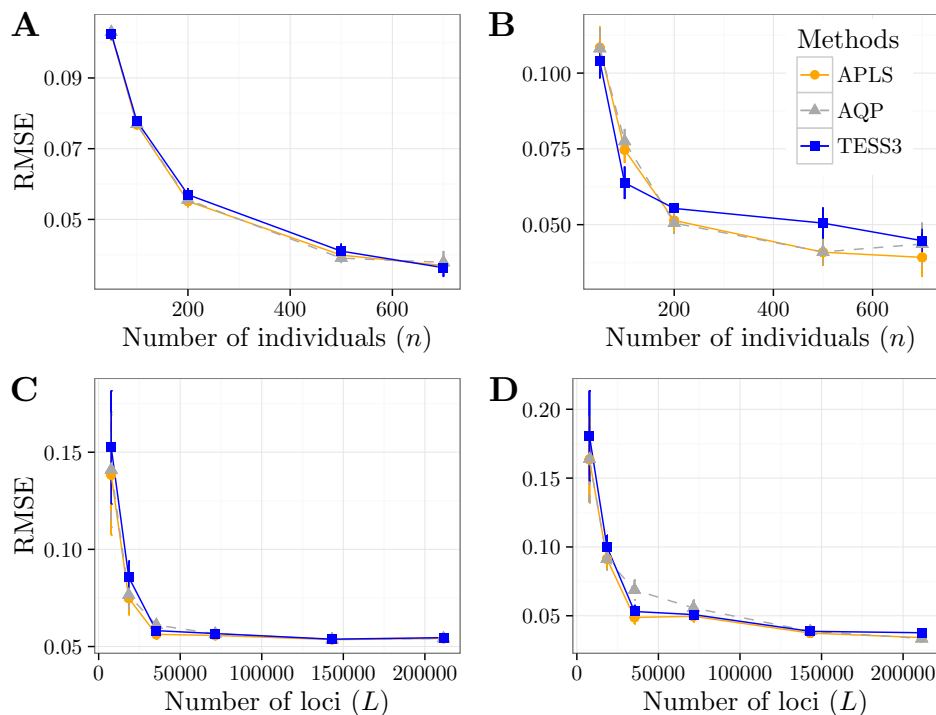
- 492 ALEXANDER, D. H. and LANGE, K. (2011). Enhancements to the ADMIXTURE  
493 algorithm for individual ancestry estimation. *BMC bioinformatics* **12** 246.
- 494 BARAN, Y., QUINTELA, I., CARRACEDO, Á., PASANIUC, B. and HALPERIN, E.  
495 (2013). Enhanced localization of genetic samples through linkage-disequilibrium  
496 correction. *American Journal of Human Genetics* **92** 882–894.
- 497 BELKIN, M. and NIYOGI, P. (2003). Laplacian Eigenmaps for Dimensionality Re-  
498 duction and Data Representation. *Neural Computation* **15** 1373–1396.
- 499 BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the False Discovery Rate:  
500 A Practical and Powerful Approach to Multiple Testing. *Source Journal of the*  
501 *Royal Statistical Society. Series B (Methodological)* **57** 289–300.
- 502 BERTSEKAS, D. P. (1995). *Nonlinear programming*. Athena Scientific.
- 503 CAI, D., HE, X., HAN, J. and HUANG, T. S. (2011). Graph Regularized Non-  
504 negative Matrix Factorization for Data Representation. *IEEE Transactions on*  
505 *Pattern Analysis and Machine Intelligence* **33** 1548–1560.
- 506 CARBON, S., IRELAND, A., MUNGALL, C. J., SHU, S., MARSHALL, B., LEWIS, S.,  
507 AMIGO HUB and WEB PRESENCE WORKING GROUP (2009). AmiGO: online  
508 access to ontology and annotation data. *Bioinformatics (Oxford, England)* **25**  
509 288–9.
- 510 CAVALLI, L. L., MENOZZI, P. and PIAZZA, A. (1994). *The History and Geography*  
511 *of Human Genes*. Princeton University Press.
- 512 CAYE, K., DEIST, T. M., MARTINS, H., MICHEL, O. and FRANÇOIS, O. (2016).  
513 TESS3: Fast inference of spatial population structure and genome scans for se-  
514 lection. *Molecular Ecology Resources* **16** 540–548.
- 515 CHEN, C., DURAND, E., FORBES, F. and FRANÇOIS, O. (2007). Bayesian cluster-  
516 ing algorithms ascertaining spatial population structure: A new computer pro-  
517 gram and a comparison study. *Molecular Ecology Notes* **7** 747–756.
- 518 CICHOCKI, A., ZDUNEK, R., PHAN, A. H. and AMARI, S. I. (2009). *Nonnegative*  
519 *Matrix and Tensor Factorizations: Applications to Exploratory Multi-Way Data*  
520 *Analysis and Blind Source Separation*. John Wiley & Sons, Ltd, Chichester, UK.
- 521 CORANDER, J., SIRÉN, J. and ARJAS, E. (2008). Bayesian spatial modeling of  
522 genetic population structure. *Computational Statistics* **23** 111–129.

FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 27

- 523 CRESSIE, N. A. C. (1993). *Statistics for Spatial Data*. Wiley Series in Probability  
524 and Statistics. John Wiley & Sons, Inc., Hoboken, NJ, USA.
- 525 DEVLIN, B. and ROEDER, K. (1999). Genomic control for association studies. *Bio-*  
526 *metrics* **55** 997–1004.
- 527 DURAND, E., JAY, F., GAGGIOTTI, O. E. and FRANÇOIS, O. (2009). Spatial infer-  
528 ence of admixture proportions and secondary contact zones. *Molecular Biology*  
529 *and Evolution* **26** 1963–1973.
- 530 ENGELHARDT, B. E. and STEPHENS, M. (2010). Analysis of population structure:  
531 A unifying framework and novel methods based on sparse factor analysis. *PLoS*  
532 *Genetics* **6** e1001117.
- 533 EPPERSON, B. K. (2003). *Geographical genetics*. Princeton University Press.
- 534 FOURNIER-LEVEL, A., KORTE, A., COOPER, M. D., NORDBORG, M.,  
535 SCHMITT, J. and WILCZEK, A. M. (2011). A map of local adaptation in *Arabidopsis thaliana*. *Science (New York, N.Y.)* **334** 86–89.
- 536 FRANÇOIS, O. and DURAND, E. (2010). Spatially explicit Bayesian clustering mod-  
537 els in population genetics. *Molecular Ecology Resources* **10** 773–784.
- 538 FRANÇOIS, O. and WAITS, L. P. (2016). Clustering and Assignment Methods in  
539 Landscape Genetics. *Landscape Genetics: Concepts, Methods, Applications* 114–  
540 128.
- 541 FRANÇOIS, O., MARTINS, H., CAYE, K. and SCHOVILLE, S. D. (2016). Controlling  
542 false discoveries in genome scans for selection. *Molecular Ecology* **25** 454–469.
- 543 FRICHOT, E. and FRANÇOIS, O. (2015). LEA: An R package for landscape and  
544 ecological association studies. *Methods in Ecology and Evolution* **6** 925–929.
- 545 FRICHOT, E., MATHIEU, F., TROUILLON, T., BOUCHARD, G. and FRANÇOIS, O.  
546 (2014). Fast and efficient estimation of individual ancestry coefficients. *Genetics*  
547 **196** 973–983.
- 548 GRIPPO, L. and SCIANDRONE, M. (2000). On the convergence of the block nonlin-  
549 ear Gauss-Seidel method under convex constraints. *Operations Research Letters*  
550 **26** 127–136.
- 551 HANCOCK, A. M., BRACHI, B., FAURE, N., HORTON, M. W., JARY-  
552 MOWYCZ, L. B., SPERONE, F. G., TOOMAJIAN, C., ROUX, F. and BERGEL-  
553 SON, J. (2011). Adaptation to climate across the *Arabidopsis thaliana* genome.  
554 *Science (New York, N.Y.)* **334** 83–86.
- 555 HORTON, M. W., HANCOCK, A. M., HUANG, Y. S., TOOMAJIAN, C.,  
556 ATWELL, S., AUTON, A., MULIYATI, N. W., PLATT, A., SPERONE, F. G.,  
557 VILHJÁLMSSON, B. J., NORDBORG, M., BOREVITZ, J. O. and BERGELSON, J.  
558 (2012). Genome-wide patterns of genetic variation in worldwide *Arabidopsis*  
559 *thaliana* accessions from the RegMap panel. *Nature genetics* **44** 212–216.
- 560 HUDSON, R. R. (2002). Generating samples under a WrightFisher neutral model  
561 of genetic variation. *Bioinformatics* **18** 337–338.
- 562 KIM, J. and PARK, H. (2011). Fast Nonnegative Matrix Factorization: an Active-  
563 Set-Like Method and Comparisons. *SIAM Journal on Scientific Computing* **33**  
564 3261–3281.
- 565 KORNELIUSSEN, T. S., ALBRECHTSEN, A. and NIELSEN, R. (2014). ANGSD: Anal-  
566 ysis of Next Generation Sequencing Data. *BMC bioinformatics* **15** 356.

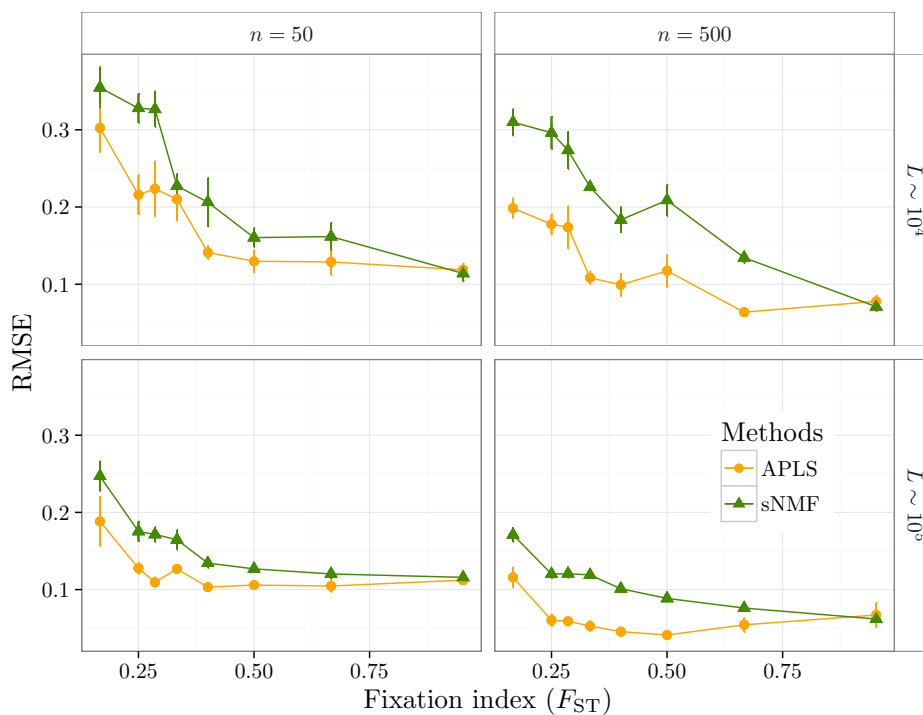
- 568 LAO, O., LIU, F., WOLLSTEIN, A. and KAYSER, M. (2014). GAGA: A New Algo-  
569 rithm for Genomic Inference of Geographic Ancestry Reveals Fine Level Popu-  
570 lation Substructure in Europeans. *PLoS Computational Biology* **10** e1003480.
- 571 LEE, D. D. and SEUNG, H. S. (1999). Learning the parts of objects by non-negative  
572 matrix factorization. *Nature* **401** 788–791.
- 573 MALÉCOT, G. (1948). *Les mathématiques de l'hérédité*. Masson et Cie, Paris.
- 574 MARTINS, H., CAYE, K., LUU, K., BLUM, M. G. B. and FRAN, O. (2016). Ident-  
575 ifying outlier loci in admixed and in continuous populations using ancestral  
576 population differentiation statistics. *Molecular ecology*.
- 577 POPESCU, A. A., HARPER, A. L., TRICK, M., BANCROFT, I. and HUBER, K. T.  
578 (2014). A novel and fast approach for population structure inference using  
579 Kernel-PCA and optimization. *Genetics* **198** 1421–1431.
- 580 PRITCHARD, J. K., STEPHENS, M. and DONNELLY, P. (2000). Inference of popu-  
581 lation structure using multilocus genotype data. *Genetics* **155** 945–959.
- 582 RAJ, A., STEPHENS, M. and PRITCHARD, J. K. (2014). FastSTRUCTURE: Vari-  
583 ational inference of population structure in large SNP data sets. *Genetics* **197**  
584 573–589.
- 585 SCHRAIBER, J. G. and AKEY, J. M. (2015). Methods and models for unravelling  
586 human evolutionary history. *Nature Reviews Genetics* **16** 727–740.
- 587 SEGELBACHER, G., CUSHMAN, S. A., EPPERSON, B. K., FORTIN, M. J., FRAN-  
588 COIS, O., HARDY, O. J., HOLDEREGGER, R., TABERLET, P., WAITS, L. P. and  
589 MANEL, S. (2010). Applications of landscape genetics in conservation biology:  
590 Concepts and challenges. *Conservation Genetics* **11** 375–385.
- 591 TANG, H., PENG, J., WANG, P. and RISCH, N. J. (2005). Estimation of individual  
592 admixture: analytical and study design considerations. *Genetic Epidemiology* **28**  
593 289–301.
- 594 WOLLSTEIN, A. and LAO, O. (2015). Detecting individual ancestry in the human  
595 genome. *Investigative genetics* **6** 7.
- 596 WRIGHT, S. (1943). Isolation by Distance. *Genetics* **28** 114–138.
- 597 YANG, W.-Y., PLATT, A., CHIANG, C. W.-K., ESKIN, E., NOVEMBRE, J. and  
598 PASANIUC, B. (2014). Spatial localization of recent ancestors for admixed indi-  
599 viduals. *G3 (Bethesda, Md.)* **4** 2505–2518.

FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 29



600

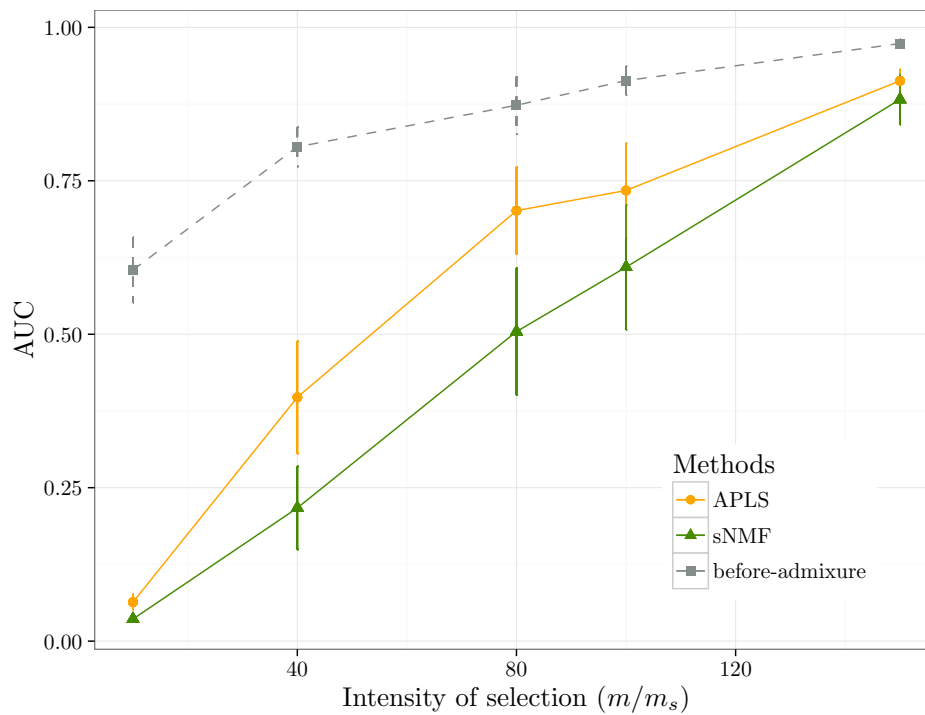
601 **Figure 1. Root Mean Squared Errors (RMSEs) for the  $Q$  and  $G$**   
602 **matrix estimates.** Simulations of spatially admixed populations. A-  
603 B) Statistical errors for APLS, AQP and `tess3` estimates as a function  
604 of the sample size,  $n$  ( $L \sim 10^4$ ). C-D) Statistical errors for APLS, AQP  
605 and `tess3` estimates as a function of the number of loci,  $L$  ( $n = 200$ ).



606

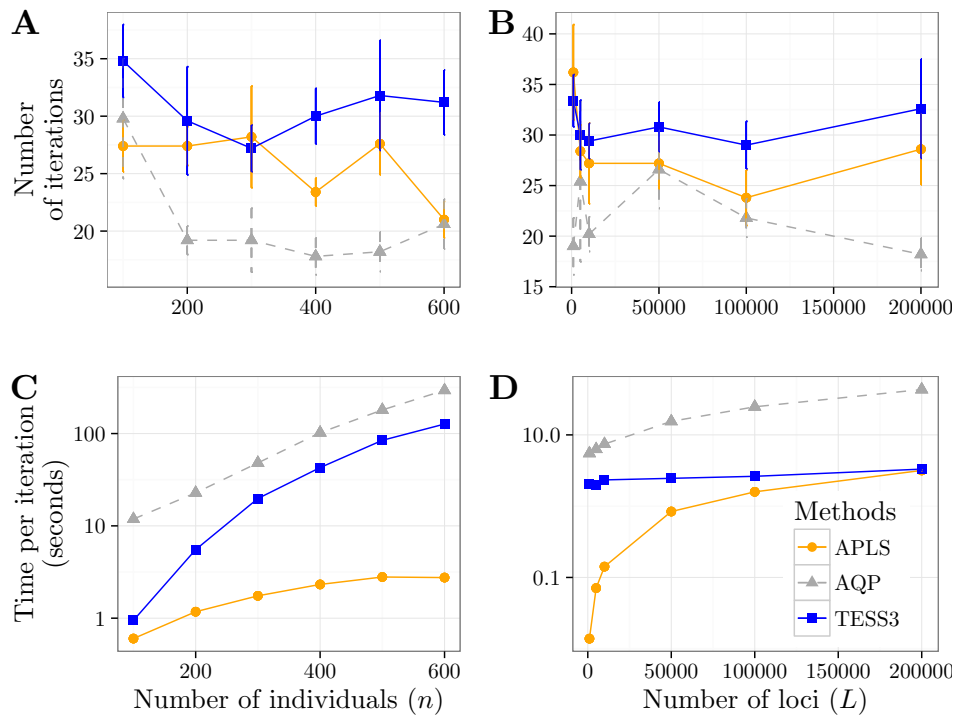
607 **Figure 2. Root Mean Squared Errors (RMSEs) for the  $Q$  esti-**  
608 **mates.** Simulations of spatially admixed populations for several values  
609 of fixation index ( $F_{ST}$ ) between ancestral populations. Ancestral popu-  
610 lations are simulated with Wright's two-island models and the fixation  
611 index is defined as  $1/(1 + 4N_0m)$  where  $m$  is the migration rate and  
612  $N_0$  the effective population size. The statistical errors for sNMF and  
613 APLS are represented as a function of  $F_{ST}$ .

FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 31



614

615 **Figure 3. Area under the precision-recall curve (AUC).** Neu-  
616 trality tests applied to simulations of spatially admixed populations.  
617 AUCs for tests based on  $F_{ST}$  with the true ancestral populations, spa-  
618 tial ancestry estimates computed with APLS algorithms, non-spatial  
619 (**structure-like**) ancestry estimates computed with the **snmf** algo-  
620 rithm. The relative intensity of selection in ancestral populations, de-  
621 fined as the ratio  $m/m_s$ , was varied in the range 1 – 160.

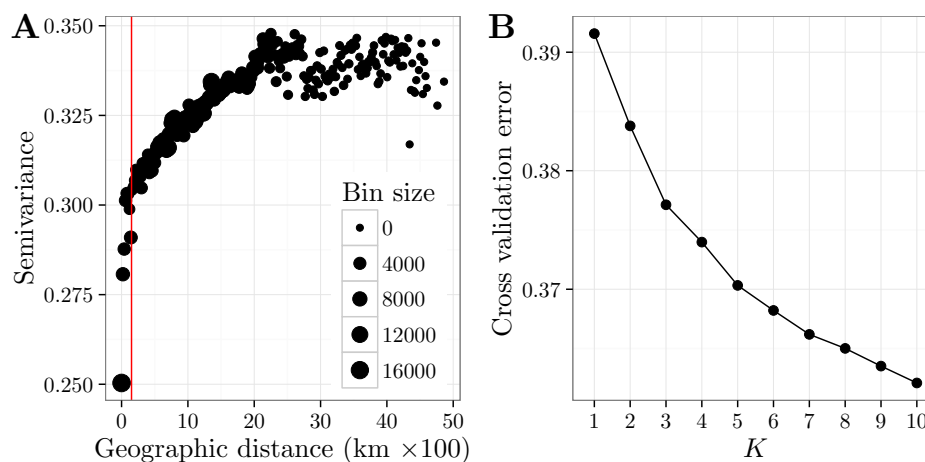


622

623 **Figure 4. Number of iterations and runtimes for the AQP,**  
624 **APLS and tess3 algorithm implementations.** A-B) Total num-  
625 ber of iterations before an algorithm reached a steady solution. C-D)  
626 Runtime for a single iteration (seconds). The number of SNPs was kept  
627 fixed to  $L = 50k$  in A and C. The number of individuals was kept fixed  
628 to  $n = 150$  in B and D.



FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 33



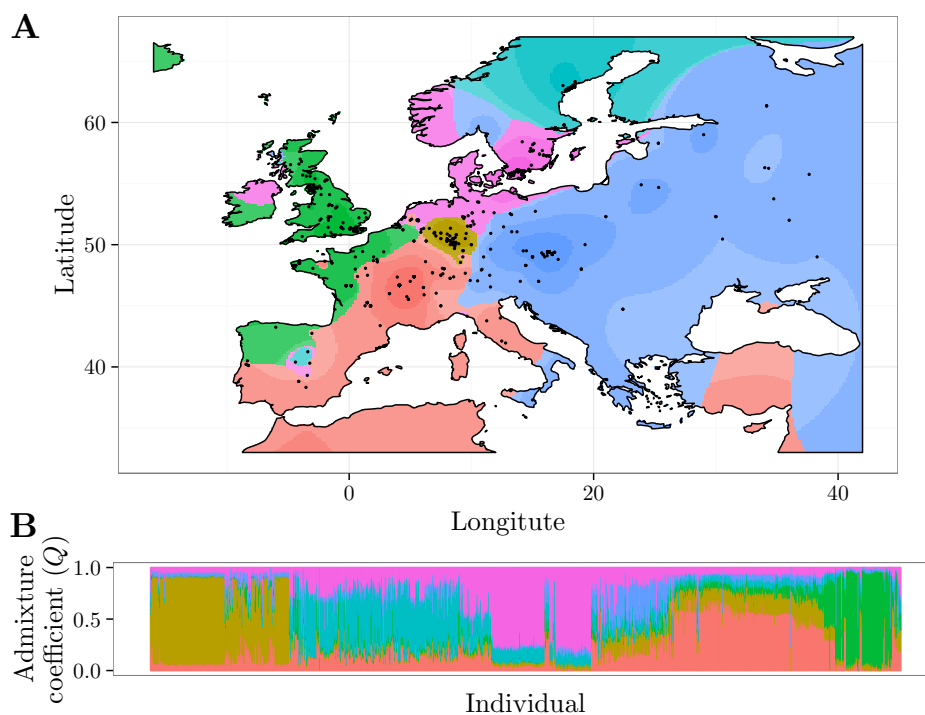
629

630 **Figure 5. Range  $\sigma$  and choice of  $K$  for the APLS algorithm.**

631 A) Empirical variogram for the *A. thaliana* data. The red vertical line

632 shows the range value  $\sigma = 1.5$ . B) Cross validation error as function of

633 the number of ancestral populations,  $K$ .

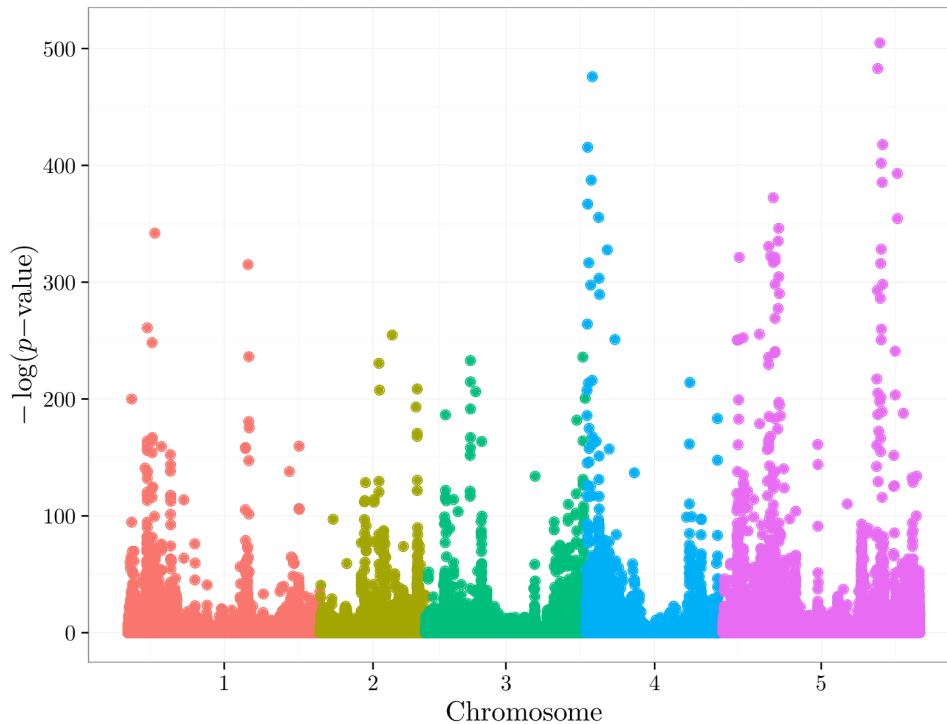


634

635 **Figure 6.** *A. thaliana* ancestry coefficients. Ancestry coefficient es-  
636 timates computed by the APLS algorithm with  $K = 6$  ancestral pop-  
637 ulations and  $\sigma = 1.5$  for the range parameter. A) Geographic map of  
638 ancestry coefficients. B) Barplot of ancestry coefficients.

FAST INFERENCE OF INDIVIDUAL ADMIXTURE COEFFICIENTS 35

639



640

641 **Figure 7. Local adaptation in European lines of *A. thaliana*.**

642 Manhattan plot of  $-\log(p\text{-value})$ .  $p$ -value were computed from popula-

643 tion structure estimated by the APLS algorithm with  $K = 6$  ancestral

644 populations and  $\sigma = 1.5$  for the range parameter.

645 KEVIN CAYE AND OLIVIER FRANÇOIS  
UNIVERSITÉ GRENOBLE-ALPES  
CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE  
TIMC-IMAG UMR 5525  
GRENOBLE, 38042, FRANCE  
E-MAIL: [kevin.caye@imag.fr](mailto:kevin.caye@imag.fr)  
[olivier.francois@imag.fr](mailto:olivier.francois@imag.fr)

FLORA JAY  
UNIVERSIT PARIS DIDEROT  
CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE  
ECO-ANTHROPOLOGIE ET ETHNOBIOLOGIE UMR 7206  
PARIS, 75013, FRANCE  
E-MAIL: [flora.jay@lri.fr](mailto:flora.jay@lri.fr)

646 OLIVIER MICHEL  
UNIVERSITÉ GRENOBLE-ALPES  
CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE  
GIPSA-LAB UMR 5216  
GRENOBLE, 38042, FRANCE  
E-MAIL: [olivier.michel@gipsa-lab.grenoble-inp.fr](mailto:olivier.michel@gipsa-lab.grenoble-inp.fr)