

Moment Estimators of Relatedness From Low-Depth Whole-Genome Sequencing Data

Anthony Francis Herzig (✉ anthony.herzig@inserm.fr)

Inserm, Univ Brest, EFS, UMR 1078, GGB

Marina Ciullo

Institute of Genetics and Biophysics A. Buzzati-Traverso - CNR and IRCCS Neuromed

Anne-Louise Leutenegger

Inserm, Université de Paris, UMR 1141

Hervé Perdry

CESP Inserm U1018, Université Paris-Saclay, Villejuif

FranceGenRef Consortium

LABEX GENMED, Centre National de Recherche en Génomique Humaine

Research Article

Keywords: Kinship, fraternity coefficient, low-depth, sequencing data, genotype likelihoods, moment estimators

Posted Date: January 25th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1109592/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Moment estimators of relatedness from low-depth whole-**
2 **genome sequencing data**

3 Herzig AF¹, Ciullo M^{2,3}, FranceGenRef Consortium⁴, Leutenegger A-L⁵, Perdry H⁶

4

5 1. Inserm, Univ Brest, EFS, UMR 1078, GGB, F-29200 Brest, France

6 2. Institute of Genetics and Biophysics A. Buzzati-Traverso - CNR, Naples, Italy

7 3. IRCCS Neuromed, Pozzilli, Isernia, Italy

8 4. LABEX GENMED, Centre National de Recherche en Génomique Humaine, Evry,
9 Paris

10 5. Inserm, Université de Paris, UMR 1141, NeuroDiderot, F-75019 Paris, France

11 6. CESP Inserm U1018, Université Paris-Saclay, Villejuif, France

12

13 Corresponding Author:

14 Anthony Francis Herzig

15 anthony.herzig@inserm.fr

16 +33298017361

17 Inserm UMR 1078, 22 Avenue Camille Desmoulins, 29238 Brest, France

18 **Abstract**

19 **Background.** Estimating relatedness is an important step for many genetic study designs. A
20 variety of methods for estimating coefficients of pairwise relatedness from genotype data
21 have been proposed. Both the kinship coefficient φ and the fraternity coefficient ψ for all pairs
22 of individuals are of interest. However, when dealing with low-depth sequencing or
23 imputation data, individual level genotypes cannot be confidently called. To ignore such
24 uncertainty is known to result in biased estimates. Accordingly, methods have recently been
25 developed to estimate kinship from uncertain genotypes.

26 **Results.** We present new method-of-moment estimators of both the coefficients φ and ψ
27 calculated directly from genotype likelihoods. We have simulated low-depth genetic data for
28 a sample of individuals with extensive relatedness by using the complex pedigree of the known
29 genetic isolates of Cilento in South Italy. Through this simulation, we explore the behaviour of
30 our estimators, demonstrate their properties, and show advantages over alternative methods.
31 A demonstration of our method is given for a sample of 150 French individuals with down-
32 sampled sequencing data.

33 **Conclusions.** We find that our method can provide accurate relatedness estimates whilst
34 holding advantages over existing methods in terms of robustness, independence from
35 external software, and required computation time. The method presented in this paper is
36 referred to LowKi (**Low**-depth **K**inship) and has been made available in an R package
37 (<https://github.com/genostats/LowKi>).

38 **Keywords:** Kinship, fraternity coefficient, low-depth, sequencing data, genotype likelihoods,
39 moment estimators.

40

41 **Background**

42 Accurate estimates of genetic relatedness between individual organisms are essential for a
43 wide range of study designs and analyses strategies currently at play in plant, animal, or
44 human genetics. These coefficients that describe the similarity and extent of shared origin
45 between genomes may currently be estimated in a large variety of ways and a multitude of
46 methods have been proposed. One's data characteristics and envisaged analyses will dictate
47 the most appropriate method to be used. For overviews of the current options for relatedness
48 estimation, and its utility, we point the reader to (1–4) and references therein.

49 In recent years the cost of whole-genome sequencing (WGS) has continued to tumble.
50 Accordingly, more and more study designs have emerged that require large sample sizes to
51 power their analyses. The depth of sequencing carried out over a large sample will have a
52 significant effect on a researcher's budget. Whilst the accuracy of genotyping is highly
53 dependent on the depth (5), there are often more advantages to being able to sequence a
54 large number of individuals but at a low depth. Recent high profile association studies using
55 this approach include (6) and (7). Indeed, low-depth sequencing data was used in many of
56 cohorts participating in the Haplotype Reference Consortium panel (8). Furthermore, shallow
57 sequencing is often unavoidable in the expanding field of ancient DNA, where the possibilities
58 of sequencing DNA from remains of long deceased organisms (9) are being widely explored.
59 Whilst technological advances allow for greater and greater accuracy in this field, in some
60 circumstances, sequencing to a high depth may simply not be feasible due to the paucity of
61 available genetic material. Another area where genetic material of high quality might be
62 difficult to ascertain is in the study of wild animal populations where DNA is collected from
63 more challenging sources such as hair, feathers, egg membranes or similar (10).

64 In adaptation to this recent trend of low-depth sequencing studies, a number of methods have
65 been proposed to estimate relatedness coefficients from such datasets. The specificity of
66 these methods is that they work upon genotype likelihoods or posterior probabilities, thus
67 incorporating the uncertainty of genotype calls. These include Hidden Markov Model (HMM)
68 based methods, maximum likelihood expectation based methods, and method-of-moment
69 estimates. The former two approaches can be computationally heavy while moment-based-
70 estimators present a quick and simple alternative. However, the loss of information entailed
71 by analysing genotype likelihoods as a proxy for true genotypes will lead to biased estimates
72 of relatedness which methods using moment-based-estimators need to account for.

73 Moment-based methods have so far only been developed for estimation of the kinship
74 coefficient and the one software that performs this estimation requires an intermediate step
75 from an existing HMM method. We propose here LowKi, a method to directly estimate genetic
76 relatedness matrices from genotype likelihoods in a single step, which has now been
77 incorporated into the genetic data management and analysis R-package 'Gaston' (11). LowKi
78 calculates moment estimates of relatedness in the form of a genetic relatedness matrices
79 (GRMs) with suitable adjustments for the genotype uncertainty that is present with low-depth
80 WGS data. For a pair of individuals i and i' , LowKi provides estimates of $\varphi^{ii'}$, the kinship
81 coefficient of individuals i and i' . This is the probability that a pair of randomly drawn alleles
82 from individual i and i' at the same locus will be in a state of Identity-by-Descent (IBD). LowKi
83 also provides a moment-estimate of $\psi^{ii'}$ the fraternity coefficient which is the proportion of
84 the genome for which individuals i and i' share two pairs of alleles at the same locus (IBD=2).
85 First, 'naïve' moment estimators were defined by an approximation of the construction of the
86 classical moment estimators used for genotype data; but based on individual genotypes

87 likelihoods. These estimators, which are referred to as ‘unadjusted estimates’ in this study,
88 are biased. Simulation studies show that as the average read depth decreases, the bias
89 increased. It is indeed intuitive that additional uncertainty or ‘fuzziness’ in the genotype
90 likelihoods gives a stronger downward bias in a moment-estimate. This makes sense when
91 considering that the additional fuzziness represents an increasing lack of information about
92 the genotypes as random error contributions to the genotype likelihoods (occurring
93 independently between individuals) become more and more prevalent. We found that it
94 sufficed to fit regression models between point-wise moment estimators and a summary
95 statistic of the genotype likelihood fuzziness to obtain accurate estimates (denominated in the
96 text as ‘adjusted estimates’).

97 To assess our approach, we have analysed both simulated and real data. Firstly, we used a
98 simulation dataset which constitutes 1,444 simulated whole-genome sequences derived from
99 the complex pedigree structure of the genetic isolates of Cilento (12–14). This simulation
100 dataset was first produced to assess phasing and imputation methods (15) before being used
101 as a tool to explore heritability estimation (16). Here we overlay a second layer of data
102 simulation to convert our simulated sequencing data into low-depth sequencing data. To
103 complement our simulation analysis, we also apply our models to a real dataset of 150
104 individuals from the FranceGenRef WGS panel (Labex GENMED <http://www.genmed.fr/>).
105 These individuals have been sequenced to a depth of 30-40× so we down-sampled individual
106 bam files to create a dataset representative of WGS data at a depth of 2.5×. Our aim was to
107 show that we can recover relatedness matrices similar to GRMs calculated on high quality
108 genotypes from low-depth data in an expedient manner. We compared our approach to two
109 existing methods: SEEKIN (v1.01) (17) and NGSRelateV2 (v2) (18,19). We show that our

110 estimates are accurate and are less time consuming to compute than those provided by
111 alternative software.

112 **Results**

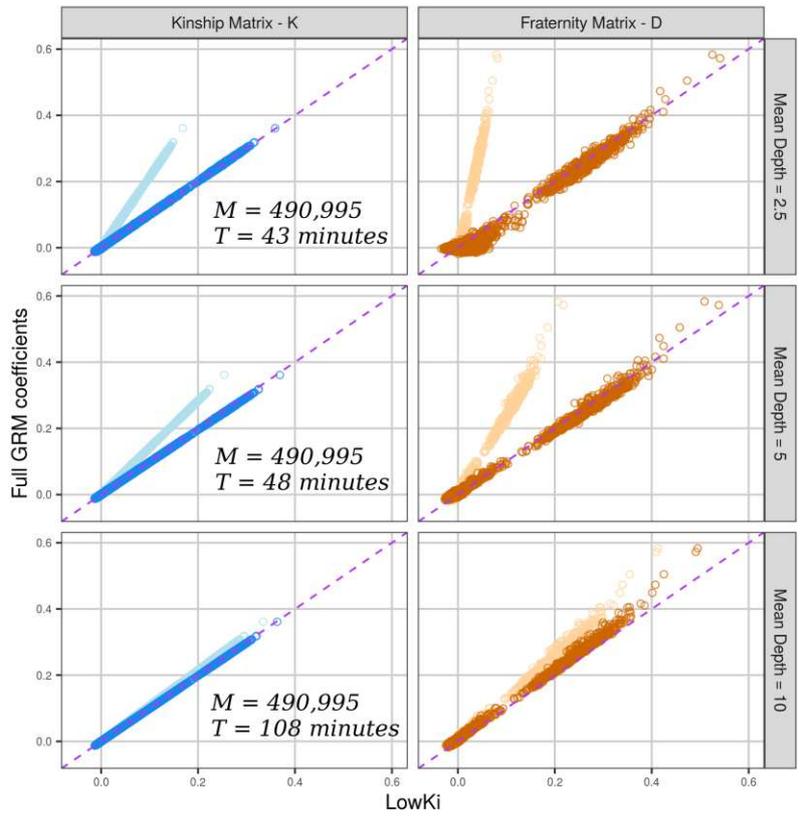
113 **Relatedness estimation from genotype likelihoods in CilentoSim**

114 Our primary simulation dataset (here denoted as ‘CilentoSim’ and described in the Methods)
115 comprises 1,444 individuals and 490,995 genetic variants across the 22 autosomal
116 chromosomes (see Methods). We established that this variant set was appropriate for the
117 calculation of a GRM as this set captured the known pedigree structure of Cilento. This is seen
118 by comparing the kinship and fraternity GRMs (calculated from the simulated genotypes) to
119 the true IBD sharing matrices calculated based on records of all haplotype mosaics created in
120 the simulation (see Methods) (Supplementary Figure 1). For Kinship, the GRM gives a very
121 precise estimate of the exact simulated IBD-sharing fractions. For fraternity, the GRM
122 estimates are highly correlated with the simulated IBD-sharing despite a lower precision
123 compared to kinship.

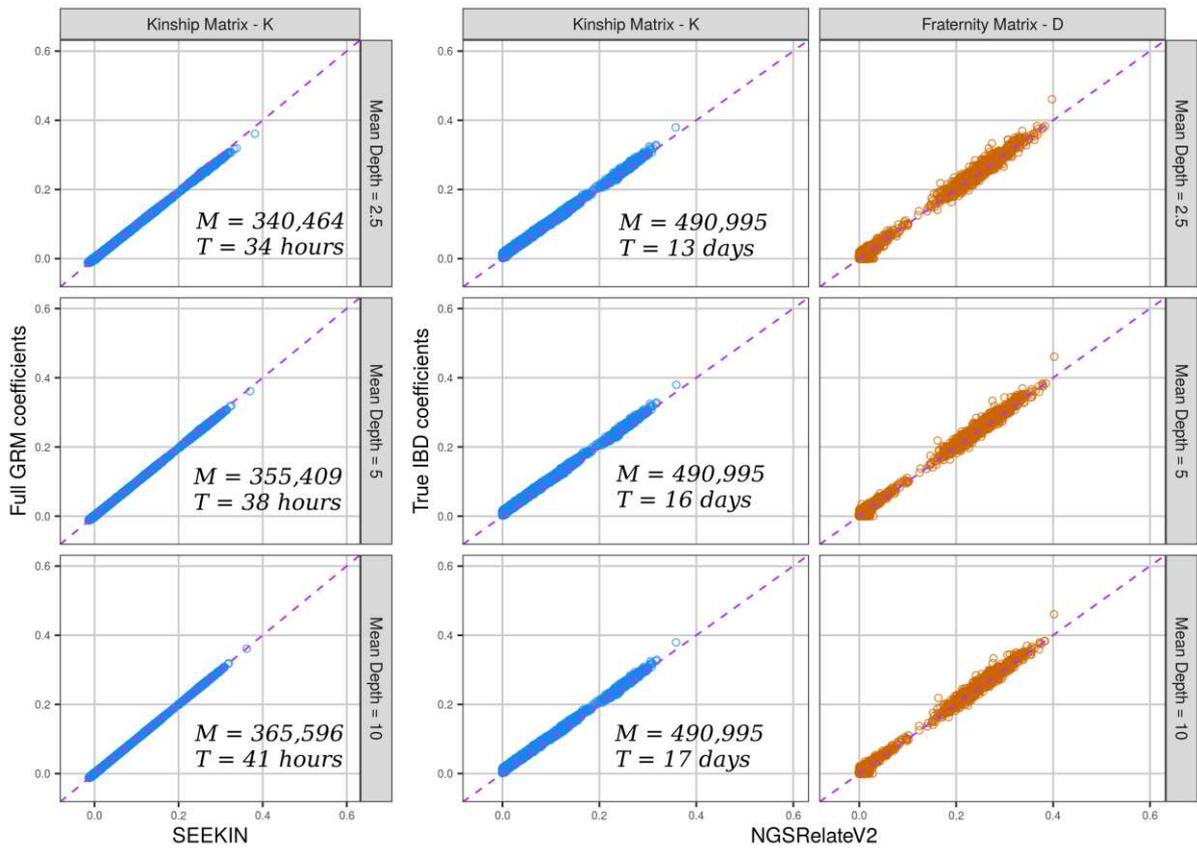
124 We artificially reduced the depth of our simulated sequencing data by drawing random alleles
125 from each simulated individual genotype to a specified depth and then replacing in our
126 simulation the true genotype with three genotype likelihoods. The method used here is based
127 on the simulation proposed by Kim et al. (20), uses a simplified version of the genotype
128 likelihood model of GATK (21–23), and is described fully in the Methods. We used this
129 additional layer of simulation to give new datasets with average read depths of 2.5×, 5×, and
130 10×.

131 We applied our method alongside SEEKIN and NGSRelateV2. In Figure 1a, the off-diagonal
132 elements of our estimated GRM matrices are compared to the the ‘Full GRM’ estimates from

133 complete simulated genotypes. As a moment estimator, SEEKIN is also compared to the 'Full
134 GRM' but NGSRelateV2 (a maximum likelihood estimator) is compared to the simulated IBD
135 sharing probabilities (which are similar but not identical to the Full GRM estimates, see
136 Supplementary Figure 1) to give a fairest assessment of the three methods. LowKi is able to
137 recover the structure of the full kinship and fraternity matrices ('Adjusted Estimates'; dark
138 blue and dark brown). Included in Figure 1a are also the 'Unadjusted Estimates' (light blue and
139 light brown) from our model that are downwardly biased by a multiplicative factor. This
140 demonstrates the efficiency of the adjustment procedure which has been sketched in the
141 Introduction and is described in detail in the Methods.



Estimator ● Kinship Matrix - K - Unadjusted ● Fraternity Matrix - D - Unadjusted
● Kinship Matrix - K - Adjusted ● Fraternity Matrix - D - Adjusted



Estimator ● Kinship Matrix - K

Estimator ● Kinship Matrix - K ● Fraternity Matrix - D

143 **Figure 1a**
144 *LowKi estimates for kinship and fraternity for CilentoSim. Off-diagonal elements of the*
145 *estimated kinship and fraternity matrices against the true simulated IBD sharing coefficient in*
146 *CilentoSim at three different simulated mean read depths (2.5×, 5×, and 10×). Lighter colours*
147 *represent the unadjusted estimated from our method and the darker colours give the final*
148 *recalibrated estimates. The number of variants (M) and the time (T) required for the calculation*
149 *of the two matrices are overlaid on the figure.*

150 **Figure 1b-c**
151 *Corresponding estimates from SEEKIN (kinship only) and NGSRelateV2.*
152

153 **Comparison to existing software on simulated data**

154 In Figure 1b and 1c, relatedness estimates given by two alternative algorithms, SEEKIN and
155 NGSRelateV2, respectively, along with running times for these programs. SEEKIN only
156 produces an estimate for the kinship matrix and indeed uses in part a similar moment-
157 estimator to our method presented here. SEEKIN gave very accurate kinship estimates. The
158 key specificities of SEEKIN involve an intermediate step of the imputation software BEAGLE
159 (v4.1) (24), the leveraging of an external reference panel (here the 1000 Genomes Project
160 phase 3 haplotype reference panel (25) was used) and a re-weighting based on the imputation
161 quality of variants in the summation that forms each GRM entry. As the initial step of BEAGLE
162 cannot be avoided, we include the runtime of BEAGLE into the run time of SEEKIN. For low-
163 depth data, running BEAGLE is very time consuming. We followed the recommendations for
164 using BEAGLE as described by the authors of SEEKIN. The use of BEAGLE will change the data
165 in two particular ways: firstly, the uncertainty present in the initial data will be largely removed
166 as BEAGLE will effectively take the prior information given to it in the form of genotype
167 likelihoods and add more precision based on similarities between pairs of individuals in the
168 sample or between individuals in the sample and the external panel of reference haplotypes
169 using the same haplotype clustering HMM machinery as is applied in BEAGLE's haplotype
170 phasing and genotype imputation methods. Secondly, running BEAGLE is likely to require the

171 removal of some variants. For example, on our simulated dataset with genotype likelihoods
172 created using a mean depth of 2.5×, for our initial dataset of 490,995, BEAGLE returns
173 information for only 340,464 variants. In Supplementary Figure 2, the difference in precision
174 is displayed between the initial genotype likelihoods supplied to BEAGLE (for a random
175 selection of 25,000 variants) and the posterior genotype probabilities. This demonstrates the
176 importance of the use of BEAGLE to the SEEKIN method.

177 A different approach is proposed by NGSRelateV2 which directly estimates relatedness
178 parameters through maximum likelihood estimation. Indeed, for the fraternity matrix, the
179 true IBD coefficients were estimated more precisely with NGSRelateV2 than with a GRM on
180 simulated genotypes (a comparison can be made between Figure 1 and Supplementary Figure
181 1). This software also produces additional information as it gives estimates for all nine
182 condensed identity-by-descent states. NGSRelateV2 gave very accurate estimates for both
183 kinship and fraternity in the CilentoSimu analysis though did require extensive amounts of
184 running time on default settings. NGSRelateV2 is multithreaded and uses four threads as a
185 default; we did not alter this default setting. Using more than the four default threads would
186 give a big increase in speed but this software will remain computationally expensive for large
187 sample sizes.

188 **Testing with real data**

189 We also applied our method, SEEKIN, and NGSRelateV2 to a set of real genotypes. 150
190 individuals with WGS data were made available to us from the FranceGenRef panel; of which
191 all individuals are not closely-related except for two pairs of siblings. We down-sampled this
192 dataset from 30-40× to 2.5× in order to create realistic low-depth WGS data. The estimates
193 of the Kinship and Fraternity matrix entries for these two sibling pairs are given in Table 1.

194 Moment-estimators on down-sampled data can were compared against moment-estimators
 195 on the original 30-40× data with one exception: the estimates from NGSRelateV2 on down-
 196 sampled data were compared to the estimates of NGSRelateV2 applied the original 30-40×
 197 data. This was because with a sample size of 150, the fraternity coefficients in particular may
 198 be imprecisely estimated even on the original data and so it would not necessarily be fair to
 199 benchmark NGSRelateV2 against GRM estimates.

Table 1		WGS data at 30-40×		Down-sampled WGS data at ~ 2.5×				
		Gaston Full GRM	NGS-RelateV2	LowKi (Unadjusted)*	LowKi	BEAGLE+ LowKi	BEAGLE+ SEEKIN	NGS-RelateV2‡
Sibling Pair 1	$\hat{\phi}$	0.258	0.268	0.122	0.262	0.254	0.268	0.256
	$\hat{\psi}$	0.294	0.379	0.059	0.192	0.237	†	0.465
Sibling Pair 2	$\hat{\phi}$	0.216	0.228	0.100	0.215	0.214	0.226	0.216
	$\hat{\psi}$	0.196	0.270	0.035	0.113	0.150	†	0.377
M (number of variants)				1,009,181	1,009,181	949,075	917,715	1,051,789
Time				4 minutes	4 minutes	15 hours **	15 hours **	6 hours
$\hat{\phi}$: estimate of the two siblings' unobserved kinship coefficient ϕ $\hat{\psi}$: estimate of the two siblings' unobserved fraternity coefficient ψ * We also give the unadjusted values from LowKi, similar to in Figure 1a ** The time required for these estimators is almost entirely due to BEAGLE, both LowKi and SEEKIN require only a few minutes † SEEKIN does not provide estimates of ψ ‡ NGSRelateV2 results on down-sampled data should be compared to results from the same software but applied to the original data that had not been down-sampled								

200

201 All methods were able to distinguish the two pairs as being closely related compared to all
 202 other pairs (Supplementary Figure 3). All methods produced accurate estimates for the kinship
 203 coefficient but the fraternity coefficient proved difficult for all methods to estimate
 204 accurately. One intuitive solution is to use the same approach as SEEKIN and first performed

205 imputation with BEAGLE; this allowed us to significantly improve our estimates of the
206 fraternity coefficients between the two sibling pairs (Supplementary Figure 3, panel (c)). This
207 allowed us to have improved estimates for the fraternity matrix. In every case, LowKi
208 underestimates both fraternity coefficients; however, NGSRelateV2 gave overestimations.
209 Even without knowing the true fraternity coefficients, an overestimation is apparent as
210 NGSRelateV2 returned fraternity coefficients outside of the expected range of likely
211 coefficients between siblings in an outbred population (approximately between 0.1 and 0.4)
212 (26). LowKi's estimates remain within this range. When running NGSRelateV2 on data without
213 down-sampling, we found significantly lower fraternity coefficients for the two sibling pairs
214 (Table 1, Supplementary Figure 3).

215 **Discussion**

216 It is intuitive that with data of the huge breadth of the whole human genome, even when the
217 quality of sequencing data is extremely low, relatedness between individuals should still be
218 captured. Existing methods have either involved maximum likelihood estimation or moment-
219 estimators of relatedness coefficients. The former estimators carry a high computation
220 burden and require a modelling of the mechanism that links true genotypes to genotype
221 likelihoods. The latter moment estimators have a lower computational burden, but will
222 however suffer from bias. This can either be dealt with using an intermediate imputation
223 algorithm to improve the data as is the case of SEEKIN that requires BEAGLE; or by attempting
224 to explicitly account for the bias as in the method we have developed here.

225 By estimating orthogonal components for additive and non-additive genotypic effects, we
226 constructed a moment estimator for the kinship the fraternity coefficient from low-depth
227 data. Such a moment-estimator has never been provided for fraternity by an alternative

228 software. Estimation of fraternity is important for classifying relatives and also for exploring
229 the effects of non-additive genetic effects (16). Our moment-estimators for fraternity were
230 sufficient to distinguish pairs of siblings in the analysis of FranceGenRef where even
231 NGSRelateV2 was unable to give perfect estimates of fraternity from data down-sampled to
232 2.5×. For the fraternity matrix our re-adjusted estimators were not as accurate as for kinship,
233 though it was clear that fraternity coefficients are harder to estimate not always perfectly
234 estimated with moment-estimators even when using genotype data. The NGSRelateV2 model
235 expects genotype likelihoods to accurately reflect the probabilities of the unobserved value of
236 the hidden genotype. This may explain why it gave very good estimates in the simulation
237 based on Cilento (where data was created with a naïve version of the GATK models and was
238 generated from simulated genotypes) but performed less well in the more complex real data
239 scenario.

240 To correct for bias in the estimates of LowKi, we introduced an innovative simulation
241 extrapolation type approach that could account for unseen error models and the inherent
242 reduction in variance typical to low-depth genetic data. These adjustment methods proposed
243 should be robust to different sources of bias arising from genotype uncertainty coming from
244 different types of bioinformatics pipeline. This could give more flexibility than likelihood based
245 methods such as lcMLkin (27), NGSRelateV2, as well as similar methods proposed for
246 estimating inbreeding coefficients (28). Note that we did not test lcMLkin here following its
247 assessment in the publication presenting SEEKIN (17). The adjustment technique developed
248 for LowKi could be harnessed in other areas of research involving shallow sequencing data.

249 We showed that the alternative to such bias-correction, using an external imputation
250 algorithm, could also lead to lengthy run times and a reliance on the accuracy of the external

251 algorithm. Notwithstanding such observations, the methodology of SEEKIN that uses the
252 intermediate step of BEAGLE is clearly often highly effective and can also be used in
253 conjunction with LowKi. Indeed, BEAGLE probably worked exceptionally well in the case of our
254 CilentoSim dataset due to the many pairs of very closely related individuals in the sample.
255 However, in the circumstance where only a small sample size is present or when an
256 appropriate reference panel cannot be ascertained (both might be the case in studies of
257 ancient DNA or small isolated populations for examples), it is beneficial to have a method for
258 proceeding directly to relatedness estimates from genotype likelihood data.

259 **Conclusion**

260 LowKi was effective at computing very accurate GRMs in a large sample of individuals with a
261 full spectrum of IBD-sharing between pairs of related individuals in a detailed simulation
262 study. We complemented this analysis by assessing an example of real low-depth genetic data
263 from FranceGenRef, where our re-adjusted relatedness coefficient estimates were able to
264 quickly and accurately identify the pairs of siblings in the sample. By analysing real data, we
265 have illustrated that our estimators perform well outside of the idealised setting of a
266 simulations. Real data will harbour phenomena such as allele-balance bias (29) or region-
267 specific sequencing error rates (30) so it was important to verify our estimators on an example
268 of true sequencing data.

269 When comparing to existing methods, LowKi does not require the use of intermediate
270 software such as BEAGLE and requires the least computation time. The innovative adjustment
271 method applied in LowKi gives flexibility to the method to account for different possible
272 sources of bias. The LowKi methods proposed here for estimating relatedness have been made
273 available at <https://github.com/genostats/LowKi> and work in conjunction with the existing R-

274 package Gaston. This represents a fast and accurate standalone option for computing kinship
275 and fraternity coefficients from low-depth sequencing data.

276 **Methods**

277 Throughout, the index $i \in 1, \dots, N$ will denote individuals (with two different individuals
278 denoted as i and i') and $j \in 1, \dots, M$ will indicate bi-allelic genetic variants. Individual level
279 genotype data are denoted as G^{ij} which takes values in $\{0,1,2\}$ for the three possible
280 genotypes AA, Aa , and aa , respectively.

281 **Simulation of low-depth data**

282 Existing simulated WGS data for 1,444 individuals based on the pedigree of the Cilento isolates
283 was our starting point (16). These simulated individuals were constructed as mosaics of
284 haplotype chunks sourced from the UK10K imputation panel. The formation of mosaic
285 haplotypes from the UK10K imputation reference panel (31) has been described in two
286 previous studies (15,16). Through gene-dropping, the individuals share chunks in accordance
287 with the known pedigree of Cilento by means of gene-dropping (32) on to the pedigree. By
288 recording the source of each chunk (within the UK10K), we have knowledge of the exact IBD-
289 sharing probabilities in the simulated population. For this study, we added an additional layer
290 of simulation to translate simulated genotypes into simulated genotype-likelihoods typical of
291 low-depth WGS data. We also only retained 490,995 variants by first selecting those with a
292 minor allele frequency above 5% and then by performing pruning on linkage disequilibrium
293 with Gaston.

294 In the Cilento cohort, there are 19 individuals with whole-genome sequencing data. These
295 individuals were sequenced to an average depth of 50-60 \times . From this dataset, we took a list
296 of per-variant mean read depths and scaled each entry so that the global mean read depth

307 would either be 10, 5, or 2.5. These lists became the lists of mean depths for each variant for
 308 our simulation. For each individual level genotype G^{ij} and for an assigned average read depth
 309 d_j for the position, we draw three sets of reads to represent the number of reads carrying the
 300 reference allele A , the minor allele a , and error reads that carry a base that matches neither
 301 A or a . The size of these three groups are denoted as R_A , R_a , and R_ε . We model the occurrence
 302 of reads with Poisson distributions and thus draw R_A as Poisson with parameter $\rho_A^{ij} d_j$, R_a as
 303 Poisson with parameter $\rho_a^{ij} d_j$, and R_ε as Poisson with parameter $\rho_\varepsilon^{ij} d_j$. These parameters
 304 have values depending on the true genotypes G^{ij} and the error rate ε_j at the position as shown
 305 in Table 2.

Table 2	$G^{ij} = 0$ Genotype AA	$G^{ij} = 1$ Genotype Aa	$G^{ij} = 2$ Genotype aa
ρ_A^{ij}	$1 - \varepsilon_j$	$\frac{1}{2}(1 - \varepsilon_j) + \frac{1}{6}\varepsilon_j$	$\frac{1}{3}\varepsilon_j$
ρ_a^{ij}	$\frac{1}{3}\varepsilon_j$	$\frac{1}{2}(1 - \varepsilon_j) + \frac{1}{6}\varepsilon_j$	$1 - \varepsilon_j$
ρ_ε^{ij}	$\frac{2}{3}\varepsilon_j$	$\frac{2}{3}\varepsilon_j$	$\frac{2}{3}\varepsilon_j$

306
 307 The values of ε_j were drawn randomly as 10^{-u_j} with u_j drawn uniformly between 2 and 3. In
 308 any case where $R_A = R_a = 0$, we set the all three genotype likelihoods to missing. In order to
 309 compute genotype likelihoods, we apply a flat prior and binomial likelihoods as used in the
 310 simplest interpretation of the GATK calling algorithm. This leads to the likelihood of the
 311 observed reads occurring given the true genotypes as proportional to $(\rho_A^{ij})^{R_A} \times (\rho_a^{ij})^{R_a} \times (\rho_\varepsilon^{ij})^{R_\varepsilon}$.

312 **Moment estimators of relatedness from low-depth**

313 In order to define our new moment-estimators for relatedness matrices, we give first a brief
 314 introduction and explanation of notations and theory. Here, the concepts of additive and non-
 315 additive components are being borrowed from the literature of quantitative genetics and in

316 particular the polygenic models first proposed by RA Fisher (33) where the genetic effects of
 317 each variant can be split into two orthogonal components. The first being the additive
 318 contribution, describing the effect that increases linearly with the number of minor alleles in
 319 the genotype, and the second being the non-additive contribution which describe the
 320 deviations away from the additive model caused by interactions between the two alleles and
 321 a single locus as is observed for example in recessive or dominant models.

322 Genotype-based GRM estimates for kinship and fraternity matrices (denoted as K and D ,
 323 respectively) can be defined as follows:

$$K_{ii'} = \frac{1}{M} \sum_{j=1}^M X_A^{ij} \times X_A^{i'j} \quad \& \quad D_{ii'} = \frac{1}{M} \sum_{j=1}^M X_D^{ij} \times X_D^{i'j}$$

324 Where X_A^{ij} and X_D^{ij} are the classical additive and non-additive components of the individual
 325 level genotypes G^{ij} which are defined as follows:

$$326 \quad X_A^{ij} = \alpha_0^j 1_{\{G^{ij}=0\}} + \alpha_1^j 1_{\{G^{ij}=1\}} + \alpha_2^j 1_{\{G^{ij}=2\}}$$

$$327 \quad X_D^{ij} = \delta_0^j 1_{\{G^{ij}=0\}} + \delta_1^j 1_{\{G^{ij}=1\}} + \delta_2^j 1_{\{G^{ij}=2\}}$$

328 where

$$329 \quad \alpha_k^j = \frac{k - 2q_j}{\sqrt{2p_jq_j}}, (k = 0, 1, 2) \text{ and } \delta_0^j = \frac{q_j}{p_j}, \delta_1^j = -1, \delta_2^j = \frac{p_j}{q_j},$$

330 q_j being the minor allele frequency of variant j and $p_j = 1 - q_j$. Alternative notations are
 331 presented in (34) and (35) but give the same moment-estimators. The values of $(\alpha_0^j, \alpha_1^j, \alpha_2^j)$
 332 are obtained through standardisation of G^j , interpreted as a random variable (the SNP index
 333 j is fixed, the sample is constituted of the values G^{ij} for $i = 1, \dots, N$): its expected value is $2q_j$
 334 and its standard deviation is $\sqrt{2p_jq_j}$ (assuming Hardy-Weinberg proportions). The resulting
 335 random variable X_A^j has expected value 0 and variance 1. The values of $(\delta_0^j, \delta_1^j, \delta_2^j)$ can then

336 be determined by imposing three constraints on the resulting variable X_D^j : $E(X_D^j) = 0$,
337 $\text{var}(X_D^j) = 1$, and $E(X_A^j X_D^j) = 0$ (the two variables are independent – or ‘orthogonal’ – in
338 the sample).

339 It is well established that under the circumstances of correct Hardy-Weinberg proportions in
340 the population and of having in-hand the correct value of the minor allele frequencies, $K_{ii'}$
341 will be an unbiased estimator of $2\varphi^{ii'}$ and $D_{ii'}$ will be an unbiased estimator of $\psi^{ii'}$. Such
342 genetic relatedness matrices were first introduced in (36,37) for kinship and in (38) for
343 heritability and have been repurposed for many other uses.

344 Such moment estimators necessitate allele frequency information. For low-depth sequencing
345 data, it is possible to estimate allele frequencies directly from genotype probabilities. This is
346 however problematic as the additional uncertainty in the data will characteristically lead to
347 increased estimates of allele frequencies as well as potential perturbations to Hardy-Weinberg
348 proportions. This can be observed in Supplementary Figure 4 where we compared observed
349 minor alleles frequencies and heterozygosity statistics from the original simulated genotypes
350 of CilentoSim against those estimated from genotype likelihoods at a depth of $2.5\times$. The
351 perturbation to allele frequencies is difficult to avoid, but the issue of potential Hardy-
352 Weinberg deviations may be circumvented by defining our additive and non-additive
353 components on estimated genotype frequencies (rather than allele) in order to correctly
354 achieve orthogonality. The derivations that we give here are equivalent to those found in
355 Vitezica et al. (35).

356 Across the sample, we estimate genotype probabilities by averaging across all genotype
357 probabilities in the sample. First, individual genotype likelihood data (typically available on a
358 log-scale) in the form GL_{AA}^{ij} , GL_{Aa}^{ij} , and GL_{aa}^{ij} are rescaled to genotype probabilities P_{AA}^{ij} ,

359 P_{Aa}^{ij} , and P_{aa}^{ij} . Then we estimate genotype frequencies in the sample as: $\bar{P}_{AA}^j = \frac{1}{N} \sum_{i=1}^N P_{AA}^{ij}$,

360 $\bar{P}_{Aa}^j = \frac{1}{N} \sum_{i=1}^N P_{Aa}^{ij}$, and $\bar{P}_{aa}^j = \frac{1}{N} \sum_{i=1}^N P_{aa}^{ij}$.

361 The additive and dominant component are defined as

362
$$\tilde{X}_A^{ij} = \tilde{\alpha}_0^j P_{AA}^{ij} + \tilde{\alpha}_1^j P_{Aa}^{ij} + \tilde{\alpha}_2^j P_{aa}^{ij}$$

363
$$\tilde{X}_D^{ij} = \tilde{\delta}_0^j P_{AA}^{ij} + \tilde{\delta}_1^j P_{Aa}^{ij} + \tilde{\delta}_2^j P_{aa}^{ij}$$

364 As previously, the values of the triplet $(\tilde{\alpha}_0^j, \tilde{\alpha}_1^j, \tilde{\alpha}_2^j)$ are obtained by standardizing the vector

365 with $(0, 1, 2)$ using the observed mean and variance of the expected minor allele count (or

366 genotype dosage) \tilde{G}^j which is constituted of the values \tilde{G}^{ij} for $i = 1, \dots, N$, where $\tilde{G}^{ij} = P_{AA}^{ij} +$

367 $2P_{Aa}^{ij}$. The values of $(\tilde{\delta}_0^j, \tilde{\delta}_1^j, \tilde{\delta}_2^j)$ are derived from the constraints $E(\tilde{X}_D^j) = 0$, $\text{var}(\tilde{X}_D^j) = 1$,

368 and $E(\tilde{X}_A^j \tilde{X}_D^j) = 0$, where, as before, expected values are computed across the sample (j is

369 fixed and i goes from 1 to N). We obtain

370
$$(\tilde{\alpha}_0^j, \tilde{\alpha}_1^j, \tilde{\alpha}_2^j) = \left(\bar{P}_{Aa}^j + 4\bar{P}_{aa}^j \bar{P}_{AA}^j - \bar{P}_{Aa}^{j2} \right)^{-\frac{1}{2}} \times (-\bar{P}_{Aa}^j - 2\bar{P}_{aa}^j, 1 - \bar{P}_{Aa}^j - 2\bar{P}_{aa}^j, 2 - \bar{P}_{Aa}^j - 2\bar{P}_{aa}^j)$$

371 and

372
$$(\tilde{\delta}_0^j, \tilde{\delta}_1^j, \tilde{\delta}_2^j) = \left(\bar{P}_{Aa}^j + 4 \frac{\bar{P}_{aa}^j \bar{P}_{AA}^j}{\bar{P}_{Aa}^j} + \bar{P}_{AA}^j \right)^{-\frac{1}{2}} \times \left(\sqrt{\frac{\bar{P}_{Aa}^j}{\bar{P}_{AA}^j}}, -2 \sqrt{\frac{\bar{P}_{aa}^j \bar{P}_{AA}^j}{\bar{P}_{Aa}^{j2}}}, \sqrt{\frac{\bar{P}_{AA}^j}{\bar{P}_{Aa}^j}} \right).$$

373 Finally, the GRM matrices using genotype likelihoods are computed as

374
$$\tilde{K}_{ii'} = \frac{1}{M} \sum_{j=1}^M \tilde{X}_A^{ij} \times \tilde{X}_A^{i'j} \text{ and } \tilde{D}_{ii'} = \frac{1}{M} \sum_{j=1}^M \tilde{X}_D^{ij} \times \tilde{X}_D^{i'j}.$$

375 **Correcting the bias**

376 Our initial simulation results indicated a clear relationship between the average depth and the
 377 biases in the estimates of both off-diagonal and diagonal elements in the GRMs. Indeed, the
 378 bias observed appeared similar to the bias that occurs when hard-called genotypes (setting
 379 the genotypes to the most probable genotype) are used for estimating GRMs as reported by
 380 Dou et al. (17). For a given average read depth, our simulation results suggest that $E[\tilde{K}_{ii'}] =$
 381 $2\beta_1\varphi_{ii'}$ and $E[\tilde{D}_{ii'}] = \beta_2\psi_{ii'}$ for some unknown constants β_1 and β_2 .
 382 Each off-diagonal element of matrices \tilde{K} and \tilde{D} are themselves averages over many point
 383 estimated from each genetic variant. These point-wise come from genetic variants with
 384 differing read depths and hence we should expect some variants to be giving greater or lesser
 385 biased point-wise estimations. When the depth is low, the three genotype probabilities tend
 386 to become less certain, we move further away from a tuple of probabilities such as (1,0,0)
 387 (which represents a certain genotype of AA) towards a tuple such as $(\frac{1}{2}, \frac{1}{2}, 0)$ or even $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$
 388 where there is no certainty as to what the true genotype may be. This uncertainty or
 389 ‘fuzziness’ of the data can be summarised by the variance of the genotype (here thought of as
 390 a random variable taking values in {0,1,2} occurring at probabilities P_{AA}^{ij} , P_{Aa}^{ij} , and P_{aa}^{ij} ,
 391 respectively. We denote this measure as $v^{ij} := P_{Aa}^{ij}(1 - P_{Aa}^{ij}) + 4P_{aa}^{ij}(1 - P_{aa}^{ij}) - 4P_{Aa}^{ij}P_{aa}^{ij}$. To
 392 demonstrate this, we simulated repeatedly low-depth data for a single variant shared (with
 393 IBD status at random) by two siblings; varying values of the average depth and minor allele
 394 frequency for the variant. Pairs of siblings are expected to share at least one haplotype IBD
 395 for 50% of their genome and to share both haplotypes IBD for 25%. By varying the depth, we
 396 could see the change in the expected bias (Supplementary Figure 5) suggesting clearly that
 397 additional uncertainty or ‘fuzziness’ in the genotype likelihoods gives a stronger downward
 398 bias in a GRM moment-estimate. In Supplementary Figure 5, the average point-wise estimates

399 of kinship are plotted against the varying values of $v^{ii'j} := \frac{1}{2}(v^{ij} + v^{i'j})$. Different mean
400 values of $v^{ii'j}$ came from simulating read data with depths varying between 2× and 25×.
401 Here, we observed broadly linear relationships, with the slope changing slightly depending on
402 the minor allele frequency of the variant. We can also see that as $v^{ii'j}$ tends to zero, the
403 multiplicative bias in our estimate tends to one; and thus the estimator becomes unbiased.
404 This suggests that if we can have a model for this relationship between bias and the fuzziness
405 of each variant, it should be possible to gain an estimation of the unbiased value of the
406 relatedness coefficients between i and i' . Hence we use an idea similar to simulation
407 extrapolation (39) though rather than artificially adding more noise to our data, we simple
408 take advantage to the different levels of noise between different SNPs and extrapolate what
409 our relatedness estimators would be with zero noise.

410 Our point wise estimates for the two matrices are written as $\tilde{K}_{ii'}^j$ and $\tilde{D}_{ii'}^j$, and we use linear
411 regression to perform what was found to be the most appropriate medialisation which was
412 the following:

$$413 \quad E[K_{ii'}^j] = (z_1 + z_2\varphi_{ii'}) (1 + z_3(v^{ij} + v^{i'j}) + z_4v^{ij}v^{i'j})$$

$$414 \quad E[D_{ii'}^j] = (u_1 + u_2\psi_{ii'}) (1 + u_3(v^{ij} + v^{i'j}) + u_4v^{ij}v^{i'j})$$

415 Here, $\varphi_{ii'}$ and $\psi_{ii'}$ are the kinship and fraternity coefficients between i and i' . The model
416 represents the intuition that when the fuzziness (v^{ij} and $v^{i'j}$) is null, the pointwise estimators
417 should have expected values of the ‘true’ pointwise estimator from full WGS data; though we
418 allow for the expectation to be linear in the ‘true’ estimator by introducing quantities z_1 and
419 z_2 for kinship and u_1 and u_2 for fraternity. Indeed, all quantities z_{1-4} and u_{1-4} are nuisance

420 parameters that allow a flexible modelling of potential biases that could be created by
421 studying low-depth data.

422 Using this model, regressing values of $\tilde{K}_{ii'}^j$ or $\tilde{K}_{ii'}^j$ across values of j against corresponding
423 values of v^{ij} and $v^{i'j}$ leads to estimates of $\varphi_{ii'}$ and $\psi_{ii'}$ from the intercepts of the linear
424 regression models. Our adjustment procedure circumvents the nuisance parameters by firstly
425 performing the aforementioned regression on the diagonal elements of the matrices \tilde{K} and \tilde{D}
426 ($i = i'$) with the knowledge that $2\varphi_{ii}$ and ψ_{ii} should be equal to 1. Then in a second step, we
427 regress the mean unadjusted estimates ($\tilde{K}_{ii'}^j$ or $\tilde{D}_{ii'}^j$) against the intercepts from the
428 aforementioned linear regression models that compared $\tilde{K}_{ii'}^j$ or $\tilde{K}_{ii'}^j$ with v^{ij} and $v^{i'j}$. These
429 steps combined provide new regression-based estimates of ϕ and φ .

430 This adjustment procedure carries a reasonable computational burden, so we simply apply it
431 to a subset of pairs which are chosen to represent a good range of relatedness estimates ($\hat{\varphi}$
432 or $\hat{\psi}$) among the unadjusted estimates in the sample calculated by LowKi. The adjustment
433 proceeds by comparing the outcome of the unadjusted estimates and the regression-based
434 estimates (described here) in order to calculate the appropriate multiplicative biases β_1 and
435 β_2 in order to finally adjust the initial unadjusted estimates of LowKi.

436 **Testing existing software**

437 To run SEEKIN (v1.01), we first applied BEAGLE (v4.1). BEAGLE was given reference haplotypes
438 from the 1000 Genomes project (Phase 3) and was run in windows of 750 variants with buffers
439 of 250 variants. We found that BEAGLE required very long runtimes, hence we set the
440 parameter 'modelscale' equal to 3 which the authors of BEAGLE suggested in the software's
441 manual as an appropriate setting to increase both speed and accuracy when applying BEAGLE

442 to genotype likelihood data. Otherwise, both NGSRelateV2 (v2) and SEEKIN were run with the
443 default recommended parameters.

444 **Testing on real data**

445 In order to test our method on a real dataset, we were given access to 150 individuals from
446 FranceGenRef and down-sampled their individual bam files to an average of 2.5× coverage.
447 The FranceGenRef panel comprises 856 individuals from the population of France and
448 combines individuals from the GAZEL cohort (www.gazel.inserm.fr/en), from the PREGO
449 cohort (www.vacarme-project.org), and 50 blood donors from the Finistere region. The down-
450 sampling was achieved by simply counting the number of reads in the original bam-files, and
451 randomly sampling the appropriate proportion of these reads given that full bam files
452 correspond to average read depth of 35×. This set of 150 individuals contains two sibling pairs
453 who have an expected kinship of 0.25 and expected fraternity coefficient of 0.25. All other
454 pairs are expected to have kinship and fraternity coefficients very close to zero. There may be
455 residual population structure in the sample as individuals of FranceGenRef come from
456 different regions in France; a country with substantial fine-scale population structure (40).
457 Down-sampling and calling were performed with samtools (v0.1.19) (41), Sambamba (v0.7.1)
458 (42), and GATK HaplotypeCaller (v3.7) which provides the genotype likelihoods that we supplied
459 to LowKi as well as NGSRelateV2, and BEAGLE followed by SEEKIN.

460 We observed that LowKi's estimators were improved if variants with a very small observed
461 expected minor allele frequency were removed from the calculation and such a filter has been
462 added. Specifically, the quantity $\bar{P}_{Aa}^j + 2\bar{P}_{aa}^j$ should be in the range 0.05 to 1.95. In the example
463 of the 150 individuals of FranceGenRef, 1,009,181 variants out of a possible 1,051,789 were
464 used in the calculation.

465 **Declarations**

466 **Ethics approval and consent to participate**

467 Not applicable.

468 **Consent for publication**

469 Not applicable.

470 **Availability of data and materials**

471 LowKi is freely available at <https://github.com/genostats/LowKi> and implemented in R.
472 Instructions for download and implementation as well as example datasets are also provided
473 at this location. The package contains a small simulated dataset allowing to test the method.

474 Contact for applications for access to genetic data from the Cilento isolates: Marina Ciullo
475 (marina.ciullo@igb.cnr.it).

476 Contacts for applications for access to simulation datasets based on the pedigree structure of
477 the Cilento Isolates: Marina Ciullo (marina.ciullo@igb.cnr.it), Anne-Louise Leutenegger (anne-louise.leutenegger@inserm.fr) and Anthony F. Herzig (anthony.herzig@inserm.fr).

479 Data from the FranceGenRef panel will be submitted to the French Centralized Data Center of
480 the France Medicine Genomic Plan that is under construction. Enquiries for the use of this
481 data can be addressed to GENMED LABEX (<http://www.genmed.fr/index.php/en/contact>).

482 **Competing interests**

483 The authors declare that they have no competing interests.

484 **Funding**

485 This work was supported by LABEX GENMED funded as part of “Investissement d’avenir”
486 program managed by Agence Nationale pour la Recherche (grant number ANR-10-LABX-
487 0013), and by the French regional council of Pays-de-le-Loire (VaCaRMe project). This work
488 was also supported by the POPGEN project as part of the Plan Médecine Génomique 2025
489 (FMG2025/POPGEN) and by Inserm cross cutting project GOLD.

490 **Authors’ contributions**

491 AFH and HP designed the method and wrote the R package. AFH was the main writer of the
492 manuscript. MC gave access to the Cilento data. A-LL participated significantly to the
493 manuscript elaboration. All authors read and approved the final manuscript.

494 **Acknowledgments**

495 We would like to thank kindly all of the participants of the Vallo di Diano project and the
496 Cilento cohort. We acknowledge the Center of Biological Resources (CHU Nantes, Hotel Dieu,
497 CRB, Nantes, F-44093, France), the Dijon CRB, the CEPH and the Genomics and Bioinformatics
498 Core Facility of Nantes (GenoBiRD, Biogenouest). We thank the FranceGenRef Consortium for
499 giving us the opportunity to work with the FranceGenRef sequencing data for this project. We
500 are also very grateful to Emmanuelle Génin for her valuable advice and insights regarding this
501 project.

502

503 References

- 504 1. Speed D, Balding DJ. Relatedness in the post-genomic era: is it still useful? *Nature Reviews*
505 *Genetics*. 2015;16(1):33–44.
- 506 2. Thompson EA. Identity by Descent: Variation in Meiosis, Across Genomes, and in Populations.
507 *Genetics*. 2013;194(2):301.
- 508 3. Weir BS, Anderson AD, Hepler AB. Genetic relatedness analysis: modern data and new
509 challenges. *Nat Rev Genet*. 2006;7(10):771–80.
- 510 4. Goudet J, Kay T, Weir BS. How to estimate kinship. *Molecular Ecology*. 2018;27(20):4121–35.
- 511 5. Sims D, Sudbery I, Iltott NE, Heger A, Ponting CP. Sequencing depth and coverage: key
512 considerations in genomic analyses. *Nat Rev Genet*. 2014;15(2):121–32.
- 513 6. Gilly A, Ritchie GR, Southam L, Farmaki A-E, Tsafantakis E, Dedoussis G, et al. Very low-depth
514 sequencing in a founder population identifies a cardioprotective APOC3 signal missed by
515 genome-wide imputation. *Hum Mol Genet*. 2016;25(11):2360–5.
- 516 7. Converge Consortium, Cai N, Bigdeli TB, Kretzschmar W, Li Y, Liang J, et al. Sparse whole-
517 genome sequencing identifies two loci for major depressive disorder. *Nature*.
518 2015;523(7562):588–91.
- 519 8. the Haplotype Reference Consortium, McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood
520 AR, et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nature Genetics*.
521 2016;48(10):1279–83.
- 522 9. Hofreiter M, Paijmans JLA, Goodchild H, Speller CF, Barlow A, Fortes GG, et al. The future of
523 ancient DNA: Technical advances and conceptual shifts. *Bioessays*. 2015;37(3):284–93.
- 524 10. Städele V, Vigilant L. Strategies for determining kinship in wild populations using genetic data.
525 *Ecology and Evolution*. 2016;6(17):6107–20.
- 526 11. Perdry H, Dandine-Rolland C, Bandyopadhyay D, Kettner L. Gaston: Genetic data handling (QC,
527 GRM, LD, PCA) & linear mixed models. CRAN. 2018; [https://cran.r-](https://cran.r-project.org/web/packages/gaston/index.html)
528 [project.org/web/packages/gaston/index.html](https://cran.r-project.org/web/packages/gaston/index.html).
- 529 12. Colonna V, Nutile T, Astore M, Guardiola O, Antoniol G, Ciullo M, et al. Campora: a young
530 genetic isolate in South Italy. *Hum Hered*. 2007;64(2):123–35.
- 531 13. Colonna V, Nutile T, Ferrucci RR, Fardella G, Aversano M, Barbujani G, et al. Comparing
532 population structure as inferred from genealogical versus genetic information. *Eur J Hum*
533 *Genet*. 2009;17(12):1635–41.
- 534 14. Nutile T, Ruggiero D, Herzig AF, Tirozzi A, Nappo S, Sorice R, et al. Whole-Exome Sequencing in
535 the Isolated Populations of Cilento from South Italy. *Scientific Reports*. 2019;9(1).
- 536 15. Herzig AF, Nutile T, Babron M-C, Ciullo M, Bellenguez C, Leutenegger A-L. Strategies for phasing
537 and imputation in a population isolate. *Genetic Epidemiology*. 2018;42(2).

- 538 16. Herzig AF, Nutile T, Ruggiero D, Ciullo M, Perdry H, Leutenegger A-L. Detecting the dominance
539 component of heritability in isolated and outbred human populations. *Scientific Reports*.
540 2018;8(1).
- 541 17. Dou J, Sun B, Sim X, Hughes JD, Reilly DF, Tai ES, et al. Estimation of kinship coefficient in
542 structured and admixed populations using sparse sequencing data. *PLOS Genetics*.
543 2017;13(9):e1007021.
- 544 18. Hanghøj K, Moltke I, Andersen PA, Manica A, Korneliussen TS. Fast and accurate relatedness
545 estimation from high-throughput sequencing data in the presence of inbreeding. *Gigascience*.
546 2019;8(5).
- 547 19. Korneliussen TS, Moltke I. NgsRelate: a software tool for estimating pairwise relatedness from
548 next-generation sequencing data. *Bioinformatics*. 2015;31(24):4009–11.
- 549 20. Kim SY, Lohmueller KE, Albrechtsen A, Li Y, Korneliussen T, Tian G, et al. Estimation of allele
550 frequency and association mapping using next-generation sequencing data. *BMC*
551 *Bioinformatics*. 2011;12:231.
- 552 21. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for
553 variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet*.
554 2011;43(5):491–8.
- 555 22. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The Genome
556 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.
557 *Genome Res*. 2010;20(9):1297–303.
- 558 23. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, et al. From
559 FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline.
560 *Curr Protoc Bioinformatics*. 2013;43:11.10.1-33.
- 561 24. Browning BL, Browning SR. Genotype Imputation with Millions of Reference Samples. *Am J*
562 *Hum Genet*. 2016;98(1):116–26.
- 563 25. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*.
564 2015;526(7571):68–74.
- 565 26. Visscher PM, Medland SE, Ferreira MAR, Morley KI, Zhu G, Cornes BK, et al. Assumption-Free
566 Estimation of Heritability from Genome-Wide Identity-by-Descent Sharing between Full
567 Siblings. *PLOS Genetics*. 2006 Mar 24;2(3):e41.
- 568 27. Lipatov M, Sanjeev K, Patro R, Veeramah KR. Maximum Likelihood Estimation of Biological
569 Relatedness from Low Coverage Sequencing Data. *bioRxiv*. 2015;023374.
- 570 28. Vieira FG, Fumagalli M, Albrechtsen A, Nielsen R. Estimating inbreeding coefficients from NGS
571 data: Impact on genotype calling and allele frequency estimation. *Genome Res*.
572 2013;23(11):1852–61.
- 573 29. Muyas F, Bosio M, Puig A, Susak H, Domènech L, Escaramis G, et al. Allele balance bias identifies
574 systematic genotyping errors and false disease associations. *Hum Mutat*. 2018/11/23 ed.
575 2019;40(1):115–26.

- 576 30. Li H. Toward better understanding of artifacts in variant calling from high-coverage samples.
577 *Bioinformatics*. 2014;30(20):2843–51.
- 578 31. The UK10K Consortium, Walter K, Min JL, Huang J, Crooks L, Memari Y, et al. The UK10K project
579 identifies rare variants in health and disease. *Nature*. 2015;526:82.
- 580 32. Wijsman EM, Rothstein JH, Thompson EA. Multipoint linkage analysis with many multiallelic or
581 dense diallelic markers: Markov chain-Monte Carlo provides practical approaches for genome
582 scans on general pedigrees. *Am J Hum Genet*. 2006;79(5):846–58.
- 583 33. Fisher RA. XV.—The Correlation between Relatives on the Supposition of Mendelian
584 Inheritance. *Earth and Environmental Science Transactions of The Royal Society of Edinburgh*.
585 1919;52(2):399–433.
- 586 34. Zhu Z, Bakshi A, Vinkhuyzen AAE, Hemani G, Lee SH, Nolte IM, et al. Dominance genetic
587 variation contributes little to the missing heritability for human complex traits. *Am J Hum*
588 *Genet*. 2015;96(3):377–85.
- 589 35. Vitezica ZG, Legarra A, Toro MA, Varona L. Orthogonal Estimates of Variances for Additive,
590 Dominance, and Epistatic Effects in Populations. *Genetics*. 2017;206(3):1297–307.
- 591 36. VanRaden PM. Genomic measures of relationship and inbreeding. *Interbull Annual Meeting*
592 *Proceedings*. 2007;(37):33–33.
- 593 37. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci*.
594 2008;91(11):4414–23.
- 595 38. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common SNPs explain
596 a large proportion of the heritability for human height. *Nature Genetics*. 2010;42(7):565–9.
- 597 39. Cook JR, Stefanski LA. Simulation-Extrapolation Estimation in Parametric Measurement Error
598 Models. *Journal of the American Statistical Association*. 1994;89(428):1314–28.
- 599 40. Saint Pierre A, Giemza J, Alves I, Karakachoff M, Gaudin M, Amouyel P, et al. The genetic history
600 of France. *European Journal of Human Genetics*. 2020;28(7):853–65.
- 601 41. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
602 format and SAMtools. *Bioinformatics*. 2009;25(16):2078–9.
- 603 42. Tarasov A, Vilella AJ, Cuppen E, Nijman IJ, Prins P. Sambamba: fast processing of NGS alignment
604 formats. *Bioinformatics*. 2015;31(12):2032–4.

605

606

607 **Figure legends**

608 **Figure 1a**

609 *LowKi estimates for kinship and fraternity for CilentoSim. Off-diagonal elements of the*
610 *estimated kinship and fraternity matrices against the true simulated IBD sharing coefficient in*
611 *CilentoSim at three different simulated mean read depths (2.5×, 5×, and 10×). Lighter colours*
612 *represent the unadjusted estimated from our method and the darker colours give the final*
613 *recalibrated estimates. The number of variants (M) and the time (T) required for the calculation*
614 *of the two matrices are overlaid on the figure.*

615 **Figure 1b-c**

616 *Corresponding estimates from SEEKIN (kinship only) and NGSRelateV2.*

617

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [LowDepthAFH20211123sup.docx](#)