1   **Putting RFMix and ADMIXTURE to the test in a complex admixed population**

2

3   Caitlin Uren[1,$], Eileen G. Hoal[1], Marlo Möller[1]

4

5   [1] DST-NRF Centre of Excellence for Biomedical Tuberculosis Research, South

6   African   Medical Research Council Centre for Tuberculosis Research, Division of

7   Molecular Biology and Human Genetics, Faculty of Medicine and Health Sciences,

8   Stellenbosch University, Cape Town, South Africa

9

10  $ Correspondence should be addressed to:

11  Dr. Caitlin Uren, Room 4036, 4[th] Floor Education Building, Francie van Zijl Drive,

12  DST-NRF Centre of Excellence for Biomedical Tuberculosis Research, South African

13  Medical Research Council Centre for Tuberculosis Research, Division of Molecular

14  Biology and Human Genetics, Faculty of Medicine and Health Sciences,

15  Stellenbosch University, Cape Town, 8000, South Africa

16  Phone: 021 938 9692

17  E-mail: caitlinu@sun.ac.za

18

19  **Running title:** Accuracy of global & local ancestry

20

21  **ORCID ID's**:

22  CU - 0000-0003-2358-0135

23  EGH - 0000-0002-6444-5688

24  MM - 0000-0002-0805-6741

25

26 **Abstract**

27

28 Global and local ancestry inference in admixed human populations can be performed

29 using computational tools implementing distinct algorithms, such as RFMix and

30 ADMIXTURE. The accuracy of these tools has been tested largely on populations

31 with relatively straightforward admixture histories but little is known about how well

32 they perform in more complex admixture scenarios. Using simulations, we show that

33 RFMix outperforms ADMIXTURE in determining global ancestry proportions in a

34 complex 5-way admixed population. In addition, RFMix correctly assigns local

35 ancestry with an accuracy of 89%. The increase in reported local ancestry inference

36 accuracy in this population (as compared to previous studies) can largely be

37 attributed to the recent availability of large-scale genotyping data for more

38 representative reference populations. The ability of RFMix to determine global and

39 local ancestry to a high degree of accuracy, allows for more reliable population

40 structure analysis, scans for natural selection, admixture mapping and case-control

41 association studies. This study highlights the utility of the extension of computational

42 tools to become more relevant to genetically structured populations, as seen with

43 RFMix. This is particularly noteworthy as modern-day societies are becoming

44 increasingly genetically complex and some genetic tools are therefore less

45 appropriate. We therefore suggest that RFMix be used for both global and local

46 ancestry estimation in complex admixture scenarios.

47

48

49 **Keywords:** South Africa; local ancestry inference; population genetics; RFMix;

50 ADMIXTURE

51 **Introduction**

52

53 Admixture, the exchange of genetic material between distinct populations, is a

54 hallmark of modern society - it can occur between closely or distantly related

55 populations (both genetically and geographically) (1000 Genomes Project

56 Consortium *et al.* 2012). This exchange of genetic material leads to population

57 structure; the pattern, timing and extent has been investigated in detail in a number

58 of populations (1000 Genomes Project Consortium *et al.* 2012; Gurdasani *et al.*

59 2015; Uren *et al.* 2016). Such studies on southern African populations are

60 particularly noteworthy as this area is postulated to be the geographical origin of

61 modern humans and therefore investigating population structure in modern southern

62 African populations may reveal more about the area's rich history (Henn *et al.* 2011).

63

64 Correctly and efficiently determining ancestral proportions in an admixed population

65 is possible by using computational and statistical algorithms that adapt to a variety of

66 demographic scenarios (Alexander *et al.* 2009; Maples *et al.* 2013; Brown and

67 Pasaniuc 2014). Furthermore, the ability to determine the ancestral origin of a

68 particular chromosomal region in an admixed individual has enabled the mapping of

69 the origins of genetic risk factors in complex disease i.e. admixture mapping

70 (Freedman *et al.* 2006; Cheng *et al.* 2009; Daya *et al.* 2014). The majority of the

71 computational and statistical tools used for global and local ancestry were however

72 tested on and tailored to 2- to 3-way admixed populations. The extension to a

73 complex 5-way admixed population and the evaluation of the resulting accuracy as

74 we present here, has rarely been done. (Daya *et al.* 2014; Uren *et al.* 2016).

75

76    A South African population with unique genetic ancestry and 5-way admixuture (the

77    South African Coloured (SAC) population as termed in the South African census)

78    received ancestral contributions from Bantu-speaking African (~30%), KhoeSan

79    (~30%), European (~20%), East Asian (~10%) and South East Asian populations

80    (~10%) (de Wit *et al.* 2010; Chimusa *et al.* 2013; Uren *et al.* 2016). The admixture

81    began approximately 15 generations ago and followed a continuous migration model

82    (Uren *et al.* 2016). This number, mode and timing of admixture events is unique and

83    creates a highly complex population. Although we are able to describe the

84    demographic model for most populations, there are some gaps in knowledge. This

85    may include not knowing which populations or specific geographical locations are the

86    best proxies for the true ancestral populations. Therefore, any studies investigating

87    an association between genetics and disease risk needs to be able to correctly

88    account for population or even individual admixture proportions within the limits of

89    the availability of current genetic data.

90

91    The first step in a study design aimed at finding a link between ancestry and disease

92    (such as genome-wide association studies and admixture mapping) is to understand

93    the ancestral composition of the study population. Ancestral origins and contributions

94    to the 5-way admixed South African population have been estimated but there have

95    been very few studies that have investigated the accuracy of the results generated

96    by the computational algorithm used (de Wit *et al.* 2010; Chimusa *et al.* 2013;

97    Petersen *et al.* 2013; Daya *et al.* 2013; Uren *et al.* 2016). Here we have set out to

98    test the accuracy of global and local ancestry inference in one of the most complex

99    admixed populations world-wide, using newly available dense genotyping data. A

100    simulated 5-way admixed population is generated and global and local ancestry

101    estimates are compared to the true values to determine the accuracy of the

102    computational algorithm.

103

104    **Methods**

105

106    **Data merging and filtering**

107

108    KhoeSan genotype data from Martin and colleagues (Martin *et al.* 2017) was merged

109    with the PAGEII dataset (Wojcik *et al.* 2018). In order to increase the number of

110    European and South East Asian reference samples in the dataset, the data was

111    merged with Gujarati Indian and European genetic data from the 1000 Genomes

112    Project (1000 Genomes Project Consortium (2010) 2010).

113

114    Preliminary data filtering included a filter for minor allele frequency (0.003),

115    missingness per genotype (max 0.05) and missingness per individual (max 0.01). A

116    total of ~776k SNPs passed these filters and formed the initial merged dataset.

117    Further data filtering is described in the appropriate sections below. Data was

118    phased using SHAPEIT2 utilizing a recombination map averaged across European

119    and African populations (The International HapMap Consortium 2007; O'Connell *et*

120    *al.* 2014). A summary of the populations in the final dataset can be seen in Table 1.

121

122    **Simulations**

123

124    The computational workflow is summarized in Figure 1. A subset of reference

125    individuals from the final merged dataset described in Table 1 was used to generate

126  a simulated dataset using admix-simu (Williams 2016). A demographic model

127  consisting of the ancestry proportions  described above and a continuous migration

128  model starting at 15 generations ago (Uren et al. (2016)), was used to generate a

129  simulated 5-way admixed population (Uren *et al.* 2016). This simulation results in a

130  heterogenous population, reminiscent of a real-world SAC population. The simulation

131  does not take post-admixture selection into account since it is highly unlikely that 350

132  years would result in distinct selection signals, rather, the inherent selection signals

133  in the source populations will be transferred in a random manner to the simulated

134  admixed population. Genotype as well as local ancestry calls were generated for this

135  simulated dataset from real reference haplotypes, thus capturing the complexity of

136  this heterogenous 5-way admixed South African population.

137

138

139  **Software choices**

140

141  Although there are a number of software programs that are able to estimate global

142  ancestry, ADMIXTURE is the most utilized. Reasons for this include the ability to

143  include related individuals in one run and to generate accurate admixture proportions

144  using relatively low-density SNP-array data (Alexander *et al.* 2009). The other widely

145  used global ancestry algorithm, STRUCTURE has been shown to overestimate

146  admixture proportions in complex populations (Cheng *et al.* 2017).

147

148  RFMix was chosen as the local ancestry inference algorithm of choice as it allows for

149  parameter optimization given the number of ancestral populations, has the inherent

150  ability to calculate local and global ancestry simultaneously, allows for array-based

151   input data as well as whole genome sequencing data, and has a proven track record

152   with admixed populations (Maples *et al.* 2013; Padhukasahasram 2014).

153   Furthermore, during a preliminary study by Daya and colleagues, RFMix was shown

154   to be the most accurate tool for local ancestry estimation in this 5-way admixed

155   South African population (Daya *et al.* 2014).

156

157   **GAI accuracy**

158

159   Reference individuals not included in the dataset used for the simulation, were

160   allocated to the dataset used for global and local ancestry inference. Global ancestry

161   proportions were determined by ADMIXTURE (Alexander *et al.* 2009) and RFMix

162   (Maples *et al.* 2013).

163

164   The ADMIXTURE analysis was performed in an unsupervised manner after filtering

165   the dataset for linkage disequilibrium as per recommendations in the ADMIXTURE

166   manual (50kb window size, step size of 10kb and $R^2$ threshold of 0.1). Relatedness

167   in the reference dataset was assessed using king (Manichaikul *et al.* 2010) and all

168   second degree relatives were removed prior to admixture analysis.

169

170   RFMix was run using default parameters, a time since admixture of 15 generations

171   (in line with the simulation) as well as 3 expectation-maximization (EM) iterations

172   (further EM iterations were not shown to increase accuracy (Maples *et al.* 2013)).

173   The correlation of the two methods by means of the Root Mean Squared Error was

174   performed in R.

175

176 **LAI accuracy**

177

178 Local ancestry calls were generated by RFMix using the same parameters as

179 described in the previous section.

180

181 The ability to correctly assign local ancestry was calculated in two ways. The first

182 determined the global accuracy i.e. how often the computational tool assigned the

183 correct ancestry (as per the simulations) and the second looked at this accuracy per

184 ancestral population (Atkinson 2018). These accuracy estimators were then

185 averaged over all individuals in the simulated 5-way admixed dataset.

186

187

188 **Data Availability**

189

190 No new genetic data was generated for this study however all reference data

191 supporting the findings of this study are available via the original publication.

192

193

194 **Results and discussion**

195

196 The aim of this study was to determine the accuracy of global and local ancestry

197 inference. In order to do this, a highly complex 5-way admixed population was

198 simulated. The local and global ancestry estimates were then compared to the true

199 simulated data.

200

201

202 **Global Ancestry Inference Accuracy**

203 The genetic diversity inherent in an admixed South African population was simulated

204 using 5 reference populations (see Methods). The average ancestry proportions

205 across these individuals were in line with what is seen in the real-world (Table 2)

206 (Uren *et al.* 2016). The simulations provided the basis with which the global ancestry

207 proportions as calculated by ADMIXTURE (Alexander *et al.* 2009) and RFMix

208 (Maples *et al.* 2013) could be compared.

209

210 Unsupervised admixture analysis of the simulated dataset by ADMIXTURE and

211 RFMix confirmed that the simulated 5-way admixed population is highly

212 heterogenous. Average ancestral proportions for both computational tools are given

213 in Table 2. The comparisons across the 5 ancestries for each simulated individual

214 are also depicted in Figure 2. Root Mean Squared Errors (RMSE) (RFMix vs

215 Simulation and ADMIXTURE vs simulation) were calculated for each ancestry. As

216 per the RMSE's, RFMix outperforms ADMIXTURE in correctly estimating admixture

217 proportions in the 5-way admixed population, with the exception of East Asian

218 ancestry where the accuracy is equal. ADMIXTURE over-estimates the Bantu-

219 speaking African contribution and under-estimates the KhoeSan ancestral

220 proportions. ADMIXTURE also overestimates European ancestry and

221 underestimates South East Asian ancestry. This is most likely due to inherent

222 European ancestry present in South East Asian populations and similarly, Bantu-

223 speaking ancestry in the KhoeSan reference population. It is likely that if more

224 homogenous reference populations were chosen, this trend would be negated but,

225    as previously mentioned, most modern day populations are admixed and therefore

226    computational tools should be able to account for this within the algorithms.

227

228    In addition, we hypothesize that the discrepancy in admixture proportions between

229    RFMix and ADMIXTURE can also be attributed to the increase in prior information

230    given to RFMix in order to determine admixture proportions i.e. phase and

231    recombination rate.

232

233

234    **Local Ancestry Inference Accuracy**

235    Beyond global ancestry proportions, the simulation of a 5-way admixed population

236    resulted in known local ancestry calls, to which calls by a computational tool can be

237    compared. The ancestral origin of each parental chromosomal region was

238    determined using RFMix. RFMix has been shown to outperform other computational

239    tools in the estimation of local ancestry in complex admixture scenarios (Daya *et al.*

240    2014). The local ancestry calls by RFMix were compared to the "true" simulated

241    ancestral origin of each region (Figure 3). The overall local ancestry inference

242    accuracy across all individuals and ancestries is ~89%; 88% accurate in calling

243    Bantu-speaking African ancestry, 87% calling KhoeSan ancestry, 95% calling

244    European ancestry, 86% calling East Asian ancestry and 85% calling South East

245    Asian ancestry. The statistical significance of RFMix's ability to call a specific

246    ancestry over another was assessed. RFMix is able to call European ancestry more

247    precisely than any of the African or Asian ancestries and it was able to call KhoeSan

248    ancestry more accurately than East Asian ancestry (Figure 3). This is consistent with

249    what was previously found and confirms that these algorithms are not tailored to

250    African populations (Maples *et al.* 2013).

251

252    The local ancestry inference accuracy estimates presented here are substantially

253    higher than previously obtained for this 5-way admixed South African population

254    (Daya *et al.* 2014). This increased accuracy can be attributed largely to the recent

255    availability of large-scale genotyping array data from the KhoeSan population which

256    is used as a reference for this admixed population.  This, in addition to the overall

257    higher SNP density, increased the accuracy from ~70% (as previously reported

258    (Daya *et al.* 2014)) to ~89%. As new datasets become available and the overlap

259    between datasets increases, we envisage this accuracy increasing even further.

260

261

262    **Conclusion**

263

264    In conclusion, the findings presented here detail the accuracy of global and local

265    ancestry inference of one of the most complex populations worldwide, which puts

266    ADMIXTURE and RFMix to the ultimate test. Due to the accuracy and versatility of

267    RFMix in determining global and local ancestry in a single program, it should be the

268    algorithm of choice to characterize more complex admixture scenarios. The inclusion

269    of accurate admixture proportions as a covariate in association studies is vital, and it

270    is our opinion that researchers studying complex admixed populations should use

271    RFMix for this purpose.

272

273     Furthermore, we demonstrate that computational tools *are* able to decipher the

274     complex African genetic history with a high degree of accuracy, but there is still

275     some room for improvement regarding the tailoring of computational tools to handle

276     diverse, admixed reference and target populations under study.

277

278     As populations become increasingly mobile, the probability of admixture is greater

279     and the extension of these and future computational tools to more genetically

280     complex populations across the world is vital and, as we have demonstrated, is

281     possible. The conclusions of this study will be increasingly relevant and

282     generalizable.

283

284

285     **Acknowledgements**

286

290

291

292     **Author Contributions**

293

294     CU designed the study, wrote the first draft of the manuscript and performed the

295     computational analyses. MM and EH helped to develop the research  and edited the

296     manuscript.

297

298

**Funding**

299

300

301 This research was funded (partially or fully) by the South African government through

302 the South African Medical Research Council and the National Research Foundation.

303 CU was supported by a fellowship from the Claude Leon Foundation.

304

305 **Conflict of Interest**

306

307 The authors declare that they have no conflict of interest.

308

309

310 **Ethical Approval and Informed Consent**

311

312 All procedures performed in studies involving human participants were in accordance

313 with the ethical standards of the institutional and/or national research committee and

314 with the 1964 Helsinki declration and its later amendments or comparable ethical

315 standards. Informed consent was obtained from all individual participants included in

316 the study. This study was approved by the Stellenbosch University Health Research

317 Ethics Committee (Reference #N11/07/210).

318

319

320 **References**

321

322  1000 Genomes Project Consortium (2010), 2010 A map of human genome variation

323      from population-scale sequencing. Nature 467: 1061–1073.

324  1000 Genomes Project Consortium, G. R. Abecasis, A. Auton, L. D. Brooks, M. A.

325      DePristo *et al.*, 2012 An integrated map of genetic variation from 1,092 human

326      genomes. Nature 491: 56–65.

327  Alexander, D. H., J. Novembre, and K. Lange, 2009 Fast model-based estimation of

328      ancestry in unrelated individuals. Genome Res. 19: 1655–1664.

329  Atkinson, E., 2018 *Calculations of accuracy comparing Williams lab simulations to*

330      *RFmix runs: eatkinson/LAIaccuracy.*

331  Brown, R., and B. Pasaniuc, 2014 Enhanced methods for local ancestry assignment

332      in sequenced admixed individuals. PLoS Comput. Biol. 10: e1003555.

333  Cheng, C. Y., W. H. Kao, N. Patterson, A. Tandon, C. A. Haiman *et al.*, 2009

334      Admixture mapping of 15,280 African Americans identifies obesity

335      susceptibility loci on chromosomes 5 and X. PLoS.Genet 5: e1000490.

336  Cheng, J. Y., T. Mailund, and R. Nielsen, 2017 Fast admixture analysis and

337      population tree estimation for SNP and NGS data. Bioinformatics 33: 2148–

338      2155.

339  Chimusa, E. R., M. Daya, M. Möller, R. Ramesar, B. M. Henn *et al.*, 2013

340      Determining ancestry proportions in complex admixture scenarios in South

341      Africa using a novel proxy ancestry selection method. PLoS ONE 8: e73971.

342  Daya, M., L. van der Merwe, U. Galal, M. Möller, M. Salie *et al.*, 2013 A Panel of

343      Ancestry Informative Markers for the Complex 5-way Admixed South African

344      Coloured Population. PLoS ONE 8: e82224.

345    Daya, M., L. van der Merwe, C. R. Gignoux, P. D. van Helden, M. Möller *et al.*, 2014

346        Using multi-way admixture mapping to elucidate TB susceptibility in the South

347        African Coloured population. BMC Genomics 15: 1021.

348    Freedman, M. L., C. A. Haiman, N. Patterson, G. J. McDonald, A. Tandon *et al.*,

349        2006 Admixture mapping identifies 8q24 as a prostate cancer risk locus in

350        African-American men. Proc.Natl.Acad.Sci.U.S.A 103: 14068–14073.

351    Gurdasani, D., T. Carstensen, F. Tekola-Ayele, L. Pagani, I. Tachmazidou *et al.*,

352        2015 The African Genome Variation Project shapes medical genetics in

353        Africa. Nature 517: 327–332.

354    Henn, B. M., C. R. Gignoux, M. Jobin, J. M. Granka, J. M. Macpherson *et al.*, 2011

355        Hunter-gatherer genomic diversity suggests a southern African origin for

356        modern humans. Proc. Natl. Acad. Sci. U.S.A. 108: 5154–5162.

357    Manichaikul, A., J. C. Mychaleckyj, S. S. Rich, K. Daly, M. Sale *et al.*, 2010 Robust

358        relationship inference in genome-wide association studies. Bioinformatics 26:

359        2867–2873.

360    Maples, B. K., S. Gravel, E. E. Kenny, and C. D. Bustamante, 2013 RFMix: a

361        discriminative modeling approach for rapid and robust local-ancestry

362        inference. Am. J. Hum. Genet. 93: 278–288.

363    Martin, A. R., M. Lin, J. M. Granka, J. W. Myrick, X. Liu *et al.*, 2017 An Unexpectedly

364        Complex Architecture for Skin Pigmentation in Africans. Cell 171: 1340-

365        1353.e14.

366    O'Connell, J., D. Gurdasani, O. Delaneau, N. Pirastu, S. Ulivi *et al.*, 2014 A General

367        Approach for Haplotype Phasing across the Full Spectrum of Relatedness.

368        PLOS Genet 10: e1004234.

369    Padhukasahasram, B., 2014 Inferring ancestry from population genomic data and its

370        applications. Front. Genet. 5:.

371    Petersen, D. C., O. Libiger, E. A. Tindall, R.-A. Hardie, L. I. Hannick *et al.*, 2013

372        Complex Patterns of Genomic Admixture within Southern Africa. PLOS Genet

373        9: e1003309.

374    The International HapMap Consortium, 2007 A second generation human haplotype

375        map of over 3.1 million SNPs. Nature 449: 851–861.

376    Uren, C., M. Kim, A. R. Martin, D. Bobo, C. R. Gignoux *et al.*, 2016 Fine-Scale

377        Human Population Structure in Southern Africa Reflects Ecogeographic

378        Boundaries. Genetics 204: 303–314.

379    Williams, A., 2016 *admix-simu: admix-simu: program to simulate admixture between*

380        *multiple populations*.

381    de Wit, E., W. Delport, C. E. Rugamika, A. Meintjes, M. Moller *et al.*, 2010 Genome-

382        wide analysis of the structure of the South African Coloured Population in the

383        Western Cape. Hum.Genet. 128: 145–153.

384    Wojcik, G., M. Graff, K. K. Nishimura, R. Tao, J. Haessler *et al.*, 2018 The PAGE

385        Study: How Genetic Diversity Improves Our Understanding of the Architecture

386        of Complex Traits. bioRxiv 188094.

387

388    **Figure Legends:**

389    **Figure 1:** Computational workflow

390

391    The number (n) of individuals included in each dataset, over all ancestral

392    populations. For details, please see the methods section.

393

394    **Figure 2:** Comparison between observed global ancestry proportions and "true"

395    proportions showing RFMix performs more accurately than ADMIXTURE in ancestry

396    determination.

397

398    Admixture proportions calculated by ADMIXTURE are in black and RFMix in red.

399    Root Mean Square Errors for every comparison are shown.

400

401

402    **Figure 3:**  Boxplot showing the accuracy with which RFMix assigns an ancestral

403    origin to a genetic region, stratified by reference population.

404

405    The median (bold horizontal line) and the upper and lower quartiles are shown. Data

406    faliing outside this range are plotted as outliers. The differences in accuracies across

407    ancestries were assessed using a Wilcoxon non-parametric test. All statistically

408    significant p values ($< 0.01$) are shown.

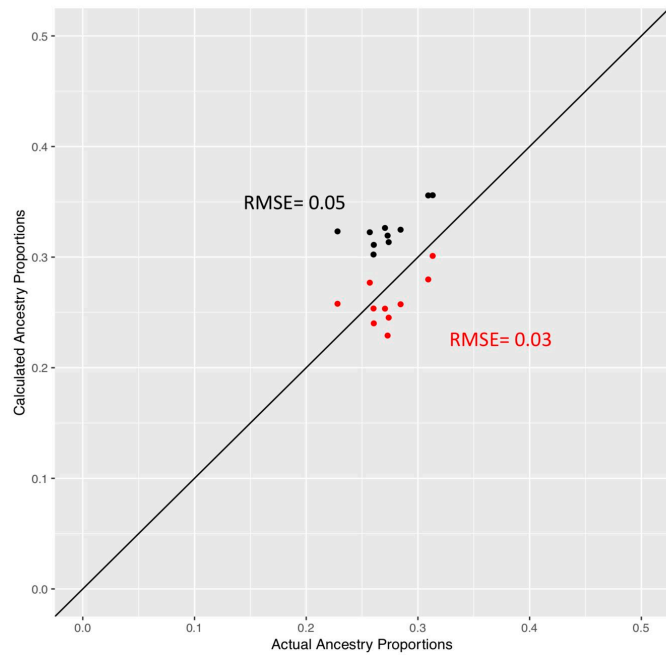**Table 1: Population characteristics of the final merged dataset**

| Population | Number of individuals included |
|---|---|
| KhoeSan (Nama and ≠Khomani San) | 284 |
| European (British) | 79 |
| African (Yoruba and Luhya) | 35 |
| East Asian (Han) | 50 |
| South East Asian (Gujarati) | 103 |

**Table 2: Average admixture proportions**

| | Previously Reported (Uren *et al.* 2016) (%) | Simulation (%) | ADMIXTURE (%) | RFMix (%) |
|---|---|---|---|---|
| Bantu-speaking African | 32 | 27 | 33 | 26 |
| KhoeSan | 30 | 33 | 25 | 33 |
| European | 19 | 22 | 26 | 23 |
| East Asian | 7 | 6 | 7 | 6 |
| South East Asian | 12 | 12 | 9 | 12 |