Edinburgh Research Explorer

# Accurate determination of breed origin of alleles in a simulated smallholder crossbred dairy cattle population

OPEN ACCESS

# Accurate determination of breed origin of alleles in a simulated smallholder crossbred dairy cattle population

Berihu Welderufael[1,2,3], Isidore Houaga[1,2*], Chris R Gaynor[1], Gregor Gorjanc[1], John M Hickey[1]

[1]The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Easter Bush, Midlothian, Scotland, UK

[2]Centre for Tropical Livestock Genetics and Health (CTLGH), Easter Bush Campus, EH25 9RG, UK

[3]Department of Animal, Rangeland and Wildlife Sciences, College of Dryland Agriculture and Natural Resources, Mekelle University, P.O. Box: 231, Mekelle, Ethiopia


[*]Corresponding author


Email addresses:


BW: berihu.gebremedhin2@mu.edu.et

IH: Isidore.Houaga@roslin.ed.ac.uk

RCG: chris.gaynor@roslin.ed.ac.uk

GG: gregor.gorjanc@roslin.ed.ac.uk

JMH: john.hickey@roslin.ed.ac.uk

# Abstract

**Background**

Accurate assignment of breed origin of alleles at a heterozygote locus may help to introduce a resilient or adaptive haplotype in crossbreeding. In this study, we developed and tested a method to assign breed of origin for individual alleles in crossbred dairy cattle. After generations of mating within and between local breeds as well as the importation of exotic bulls, five rounds of selected crossbred cows were simulated to mimic a dairy breeding programme in the low- and middle-income countries (LMICs). In each round of selection, the alleles of those crossbred animals were phased and assigned to their breed of origin (being either local or exotic).

**Results**

Across all core lengths and modes of phasing (with offset or no), the average percentage of alleles correctly assigned a breed origin was 95.76%, with only 1.39% incorrectly assigned and 2.85% missing or unassigned. On consensus, the average percentage of alleles correctly assigned a breed origin was 93.21%, with only 0.46% incorrectly assigned and 6.33% missing or unassigned. This high proportion of alleles correctly assigned a breed origin resulted in a high core-based mean accuracy of 0.99 and a very high consensus-based mean accuracy of 1.00. The algorithm's assignment yield and accuracy were affected by the choice of threshold levels for the best match of assignments. The threshold level had the opposite effect on assignment yield and assignment accuracy. A less stringent threshold generated higher assignment yields and lower assignment accuracy.

**Conclusions**

We developed an algorithm that accurately assigns a breed origin to alleles of crossbred animals designed to represent breeding programmes in the LMICs. The

50    developed algorithm is straightforward in its application and does not require prior

51    knowledge of pedigree, which makes it more relevant and applicable in LMICs

52    breeding programmes.

53

# Background

Dairy cattle production in low- and middle-income countries (LMICs) is characterised by low-input and low-output production systems. To increase the milk productivity of dairy cattle, crossbreeding between the high-producing breeds of developed countries and the low-producing, but resilient breeds of LMICs has been practised for decades. Crossbreeding, either via the importation of semen from elite bulls or the use of imported bulls, has substantially increased milk production and farmers' income [1]. However, this genetic gain has not always been observed, and overreliance on import without judicious use of best alleles is not expected to deliver the best possible genetic gains.

In many LMICs, including those in Eastern Africa, efforts are being undertaken to establish sustainable breeding programmes for long-term genetic gains with a focus on smallholder farmers [2]. In partnership with government and non-government organizations, projects like the African Dairy Genetic Gains (ADGG, https://africadgg.wordpress.com) have been able to import and provide improved dairy genetics to smallholder farmers in the Eastern Africa. However, because of the differences in environmental factors and breeding infrastructure, the importation and provision of improved genetics have not yet been sustainable and successful [2]. Instead, such crossbreeding practices have led to haphazardly admixed cattle populations with no or poor pedigree records [2].

For a sustainable breed improvement through genetic intervention and for the appropriate design of breeding programmes, accurate breed identification, on both the level of the individual and of the individual genetic variant, is important. In livestock

4

79    populations with little or no pedigree records, the use of genomic data could be

80    transformational in determining breed composition and establishment of breeding

81    programmes [2]. Estimates of breed composition and the breed origin of alleles from

82    genomic data is superior to estimates from pedigree data due to invariably missing or

83    inaccurate records and deviations from expected compositions due to Mendelian

84    sampling [3,4]. Especially in populations with complex ancestries like the dairy cattle

85    in the Eastern Africa, genomic data and knowledge of breed composition is essential

86    to evaluate the performance and adaptability of the crossbreds [4], and to predict the

87    effectiveness of any foreign germplasm in the production systems.

88

89    Selection, genetic discovery and management decisions can be aided by determining

90    the breed origin of alleles, particularly for genetic variants that only occur in one of

91    the constituent populations of crossbred animals [5]. Unlike determining the average

92    breed composition of an individual, determining the breed origin of an individual's

93    haplotypes and associated alleles can allow breed-specific genetic evaluations to be

94    conducted, which can increase the accuracy of genetic selection, particularly when the

95    linkage disequilibrium pattern is different in the two breeds [6]. Thus, recent studies

96    in admixed cattle populations have shown that the Breed Origin of Allele (BOA)

97    method has increased the accuracy of genomic prediction [7,8].

98

99    Using only genomic data and no pedigree data, Vandenplas et al. [5] developed an

100   approach that traces haplotypes of crossbred animals and assigns each allele of the

101   haplotypes to their breed of origin. To develop the algorithm that assigns alleles of

102   crossbreds a breed origin, they simulated a three-way pig-crossbreeding programme

103   with five generations of random selection. They evaluated the accuracy of the

104    algorithm and reported that more than 90% of alleles of crossbred animals were

105    correctly assigned a breed origin. Thus, for up to 10% of all alleles of crossbred

106    animals, they could not assign a breed origin. However, accurate determination of the

107    breed origin of alleles of crossbred populations is very important to estimate breed-

108    specific effects of alleles when performing genomic evaluations [9]. If we could

109    accurately assign breed origin for alleles at heterozygote loci of crossbred animals, we

110    may be able to detect which haplotypes should be promoted to genetically improve

111    dairy cows in the LMICs.

112

113    In the current study, we developed an algorithm to assign a breed of origin for alleles

114    in crossbred dairy cattle and tested it on a simulated smallholder dairy cattle

115    population dataset. To resolve the breed origin of alleles, the algorithm aligns the

116    haplotypes of crossbred dairy cows to the haplotypes of likely constituent breeds, i.e.,

117    imported (exotic) and/or local breeds and assigns the breed of origin based on the best

118    match. We then evaluated the algorithm's accuracy using a simulated crossbreeding

119    programme designed to mimic the ADGG smallholder genotype data. The average

120    percentage of alleles correctly assigned a breed origin was 95.76%, resulting in a high

121    core-based mean accuracy of 0.99 and a very high consensus-based mean accuracy of

122    1.00. The developed algorithm does not require prior pedigree knowledge and is,

123    hence, straightforward to apply in LMIC breeding programmes.

124

# Methods

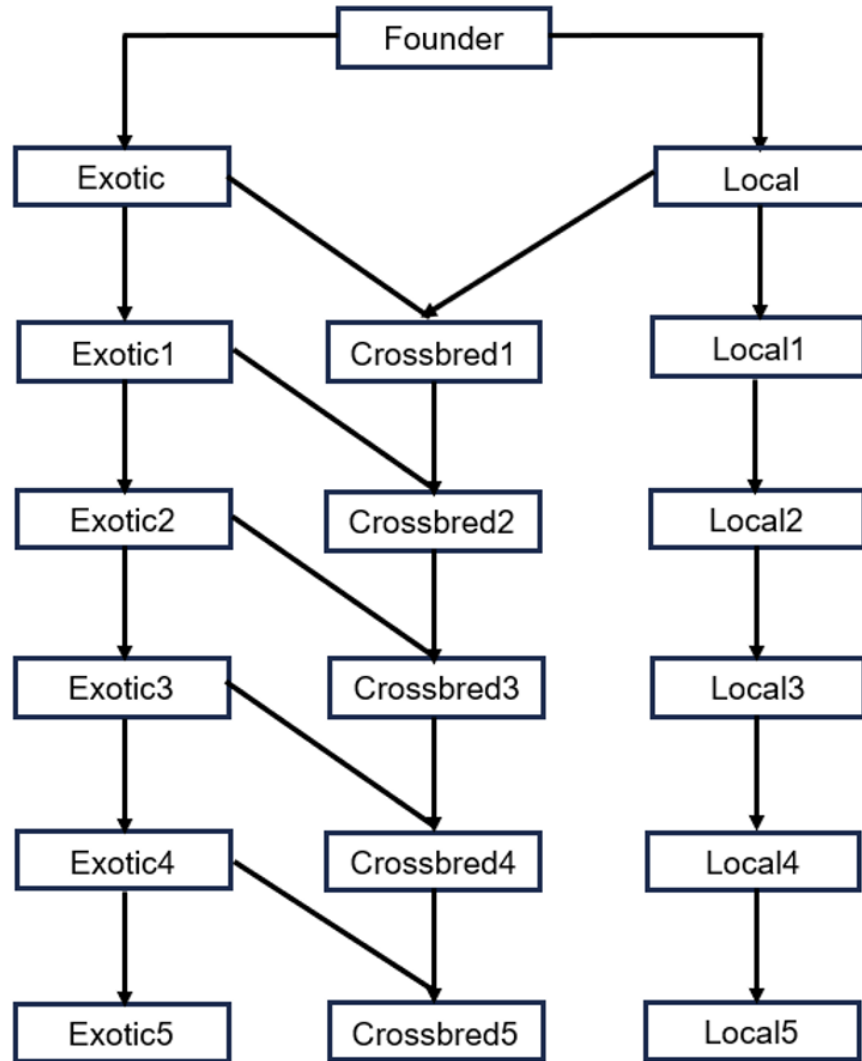The design of the breeding programme and development of the allele assignment algorithm involved two steps.

1. We designed a breeding programme and simulated genotype data on which we tested the algorithm's performance. The simulated genotype data had an ancient cattle founder that is assumed to have split into African (local) and European (exotic) cattle populations. After generations of mating within and between local breeds and the importation of exotic bulls, crossbred dairy cows were created to mimic the dairy cows kept by smallholders in the LMICs.

2. We developed an allele assignment algorithm that traces haplotypes and assigns a breed origin for each allele of the crossbred cows. The haplotypes are phased and defined for different core lengths to improve the accuracy and efficiency of the allele assignment algorithm.

The following subsections describe the details for simulating and phasing genotypes and developing the allele assignment algorithm.

## Simulation of genotype data

Genotype and haplotype data for an ancient cattle breed were simulated using the AlphaSimR package [10], designed for stochastic simulations of breeding programmes. A total of 2500 individual animals with a genome structure of 1000 SNPs in one autosomal chromosome were simulated. The ancient cattle breed split into two, each representing an exotic breed and an indigenous breed. The indigenous breed further split into four more closely related local founder populations. Variation in the demographic history of these founder populations were accounted for in the simulated biallelic haplotypes of the breeds using the Markovian Coalescent Simulator (MaCS) software [11] embedded in the AlphaSimR package [10]– [See

150     Additional file 1, Script S1] for details. As described in the AlphaSimR, the

151     genotypes and haplotypes of the descendants, i.e., the crossbred animals, were then

152     derived from these haplotypes using simulated mating between the exotic and local

153     breeds. After within and between breed random mating of indigenous animals for 10

154     generations, the 1000 best females were selected on genetic merit of a single

155     hypothetical trait with a small amount of dominance (mean dominance degree of 0.1

156     and variance of dominance degree of 0.1) and heritability of $h^2=0.3$. The 1000

157     selected local cows were then mated to 25 imported Holstein bulls to produce the first

158     crossbred animals (crossbred1). The local cows were allowed to calve twice

159     producing a total of 2000 offsprings with the assumption of 1000 female and 1000

160     male calves. The breeding programme continued by using all the 1000 female calves

161     (crossbred1) as replacement heifers and mating these to 25 newly imported Holstein

162     bulls to produce the next crossbred cows (crossbred2), while both exotic and local

163     populations were kept as purebred and source of purebred animals. This importation

164     of exotic bulls and mating to the crossbred cows was repeated for up to five rounds of

165     selections, hereafter referred as generations (Fig. 1). Simulated genotype and

166     haplotype data were generated in 10 replicates.

167

**Figure 1 Schematic representation of the simulated breeding programme.** A
founder population on the top of the figure is split into exotic and local breeds.

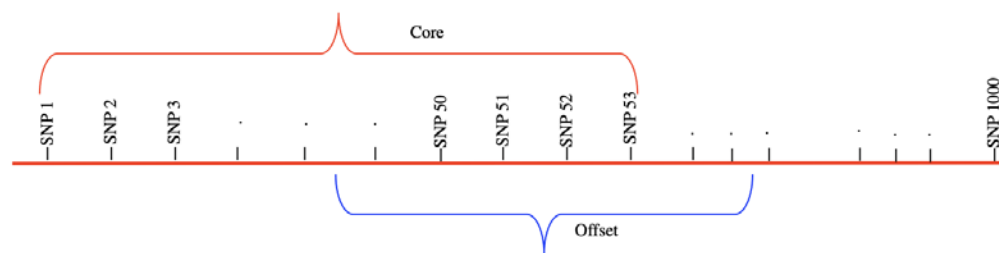## Genetic structure of the simulated SNP genotype data

To assess the genetic similarity between the founders and developed crossbred
animals, we performed principal component analysis (PCA) of SNP genotypes on the
simulated data. The PCA was performed using the prcomp command of the R
statistical software [12].

175 **Phasing of simulated genotype data**

176 True simulated genotype and haplotype data enabled us to calculate the phasing yield

177 and allele assignment accuracy. From the genotype data, haplotypes were

178 reconstructed and compared with the simulated haplotypes. The reconstruction of

179 possible haplotypes from the genotype data via phasing was performed using the

180 software AlphaPhase [13]. Different core and tail lengths govern the length of desired

181 haplotype segments used to phase the alleles in the genotype data. As illustrated in

182 Fig. 2, a core is a string of consecutive SNP loci used to phase a given genome region

183 [13].

184 Phasing of the simulated genotype data was performed using a wide range of core and

185 tail lengths. Preliminarily analyses suggested that core lengths of 100 to 280 SNPs

186 would yield optimum allele assignments. Therefore, for the final analyses, we defined

187 10 different core lengths centred around 280 SNPs (Table 1) and phasing was

188 performed for each core length both in the offset and no-offset modes of the

189 AlphaPhase [13]. We moved 50% of the core length forward to define Offset. In total,

190 there were 2000 scenarios: 10 (replicates) x 10 (core lengths) x 10 (thresholds) x 2

191 (offset or no offset modes).



192

193 **Figure 2 Illustration of a core and offset.** Phasing was performed in two modes:

194 either using the whole length of a core or by moving it forward 50% of the core length

195 (offset) to define the begging of a given core.

196

197 **Table 1 Core lengths (in terms of numbers of SNPs) used to phase the genotype**

198 **data**

| Core | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Core length (SNPs) | 100 | 120 | 140 | 160 | 180 | 200 | 220 | 240 | 260 | 280 |

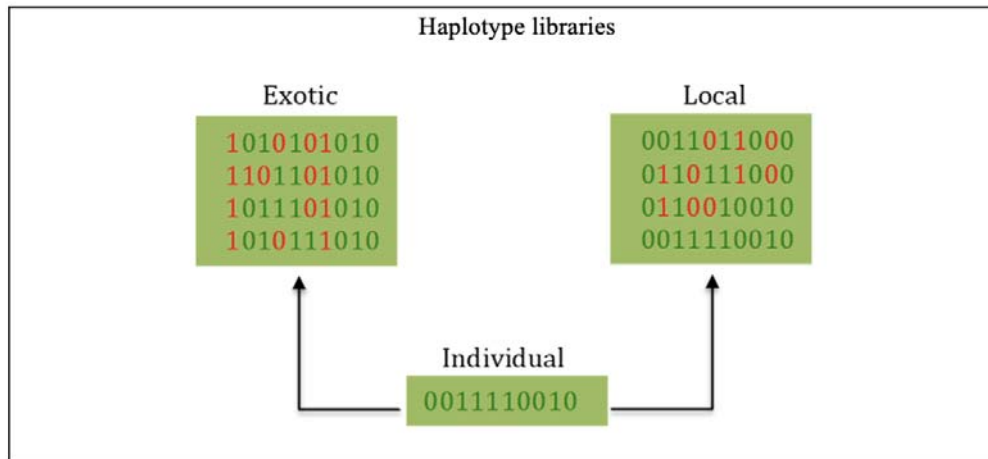199 **Development of allele assignment algorithm**

200 To develop the allele assignment algorithm, we defined 10 different core lengths

201 (Table 1). The alleles of crossbred animals were assigned a breed origin for each core

202 length, and we call this core-based allele assignment. In the core-based assignment

203 each allele could be assigned a breed origin as many as the different core lengths

204 defined. If breed origin assignments of an allele were not the same across the different

205 cores the most frequent breed assignment was considered as a consensus breed origin

206 of an allele.

207

208 **Core-based allele assignment**

209 Haplotype libraries were simulated based on the phased purebred individuals in each

210 population. The assignment algorithm takes phased genotypes for individuals in the

211 crossbred population as inputs, along with haplotype libraries for the indigenous and

212 exotic populations (Fig. 3). To perform allele assignment, we determined whether the

213 exotic or local haplotype contained the best matching haplotype, i.e. the haplotype

214 with the fewest number of markers than the target haplotype. The haplotype is then

215 assigned as originating from that haplotype library. If both haplotype libraries contain

216 an equally good match, then the haplotype is set to missing. For example, in Fig. 3,

217 the haplotype with a core length of 10 SNPs of the individual animal should be

11

218  assigned to the local haplotype as it displays the least error matches with the last core

219  in the local haplotype library.



220

221  **Figure 3 Haplotype libraries based on a core length of ten SNPs.** To assign origin

222  to the haplotype of an individual (bottom genotype sequence), the algorithm searches

223  for the best match in each position in the exotic (top left genotype sequence) and local

224  (top right genotype sequence) haplotypes. In this case, the individual's haplotype

225  should be assigned as a local haplotype because the local haplotype library contains

226  the haplotype with the fewest number of errors, i.e., mismatches (red).

227  **Consensus allele assignment**

228  Allele assignment was compared in each phased genotype and each scenario. Phasing

229  of simulated genotype data was performed in two modes: either using the whole

230  length of a core or by moving it forward 50% of the core length (offset) to define the

231  beginning of a given core (see next section). Assignment was performed across

232  multiple core lengths and two modes of phasing (no offset and offset). Assignment

233  results of each core and mode of phasing were compared and merged across cores to

234  calculate consensus-based assignment. Merging was done by taking a consensus

235  estimate of the breed of origin across multiple cores. The most frequently observed

236 assignment across all the replicates, core lengths, and phasing modes was then taken

237 as the consensus-based assignment.

238 To optimise and fine-tune the algorithm's sensitivity, we applied 10 different

239 thresholds for best SNP count of match of haplotypes (Table 2). When the threshold

240 was 0.9, this meant that the breed assignment for the allele needed to be consistent

241 across 90% of the cores, otherwise the assignment was set to missing. To elaborate a

242 threshold of 50%, an allele would have been assigned a breed origin of "A" if the

243 allele had been assigned to breed "A" in more than 50% times of the assignments

244 across all the different core lengths and phasing modes. In every generation, every

245 allele of the crossbred animals was assigned a breed origin in at least 2000 scenarios

246 and results were merged to calculate consensus assignment.

247

248 **Table 2 The different thresholds used for the best count of match of haplotypes**

| Threshold | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| %Matched | 0.50 | 0.55 | 0.60 | 0.65 | 0.70 | 0.75 | 0.80 | 0.85 | 0.90 | 0.95 |

249

250 **Performance of the allele assignment algorithm**

251 To evaluate the performance of the allele assignment algorithm, assignment yield and

252 assignment accuracy were assessed in the following ways:

253    1.  %Correct: the percentage of correctly assigned alleles was computed by

254        comparing the algorithm-derived breed origin with the true breed origin of

255        alleles traced with the "pullIbdHaplo()" function of the AlphaSimR [10].

256    2.  %Incorrect: the percentage of alleles across all scenarios that were incorrectly

257        assigned and was computed by comparing the algorithm-derived breed origin

258    with the true breed origin of alleles traced with the "pullIbdHaplo()" function

259    of the AlphaSimR [10].

260    3. %Unassigned: the percentage of alleles that were not assigned, including

261    missing or unknown breed origin; and

262    4. Accuracy: the accuracy of assigned alleles, calculated as the ratio of correctly

263    and incorrectly assigned alleles. We used the proportion of correctly assigned
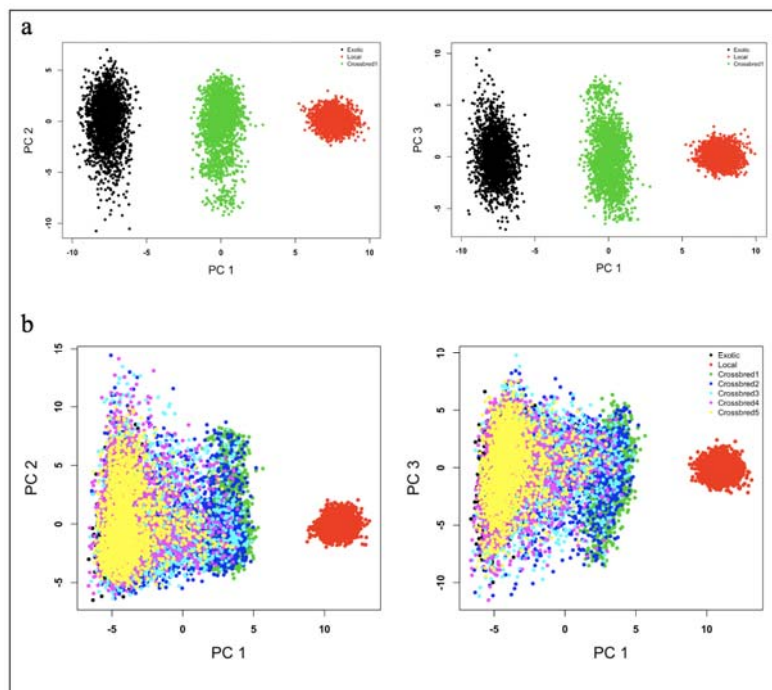
264    alleles as an allele assignment accuracy metric for each core and tail lengths.

265

# Results

## Genetic structure of the simulated SNP genotype data

Principal component analysis (PCA) of the simulated SNP genotype data separated the crossbreds from the founder breeds (local and exotic breeds). As shown in the PCA plot (Fig. 4a), the first generation of crossbred animals (crossbred1) were positioned in between the founder populations (exotic and local). The PCA plot further revealed the genetic sub-structure from the crossbreeding programme. As we continued the crossbreeding and increased the proportion of exotic genotypes, the crossbreds and the exotic breed were observed to converge into a single cluster (Fig. 4b).



**Figure 4 Plot of principal component analysis of SNP genotypes (PC1 vs. PC2 and PC1 vs. PC3).** Showing the genetic data structure of the founders and the first crossbred cows (a) and of all animals across generations (b).

280

## Allele assignment yield and accuracy

### Allele assignment for each core

283  The average number of alleles in the crossbred animals assigned a breed origin is

284  given in Table 3. The highest average number of unassigned alleles (29 out of 1000

285  SNPs) was observed in the first generation of the crossbred animals (crossbred1). The

286  number of unassigned alleles decreased as the crossbreeding continued and the

287  distance between the local founders and the crossbreds decreased. For example, in

288  crossbred5, where the germplasm is upgraded to almost the exotic breed, 23 out of

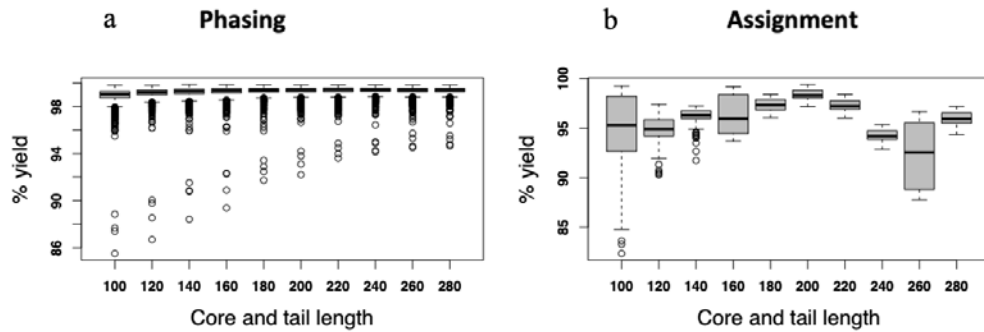289  1000 SNPs were unassigned (Table 3).

290

291  **Table 3 Assignment yield and average number of alleles in crossbred cows**

292  **assigned to local, exotic or not assigned at all**

| Crossbred | Local | Exotic | Unassigned | Assignment Yield |
|-----------|-------|--------|------------|------------------|
| Crossbred1 | 486 | 486 | 29 | 0.95 |
| Crossbred2 | 246 | 730 | 24 | 0.95 |
| Crossbred3 | 123 | 853 | 24 | 0.96 |
| Crossbred4 | 61 | 916 | 24 | 0.96 |
| Crossbred5 | 29 | 947 | 23 | 0.97 |
| Mean | 189 | 786 | 25 | 0.96 |

293

294  The genetic distance and core lengths had a clear effect on the phasing and

295  assignment yield. For longer core lengths (core length of 220-280 SNPs), we

296  observed a more concise and higher phasing yield (Fig. 5a). A core length of 200

297  SNPs was observed to be optimal for allele assignment yield (Fig. 5b). The overall

298  average allele assignment accuracy was 0.99 (Table 4). On average, more than 95%

299  of the assigned alleles in the crossbred animals were correctly assigned, with only less

300  than 2% of incorrectly assigned alleles (Table 4). Both, the incorrectly assigned and

301    unassigned proportion of alleles, either because of missing or ambiguity, were less

302    than 5% (Table 4).

303    

304    **Figure 5 Effect of core length on assignment yield.** Phasing yield (a) was very high

305    for all core lengths but more concise for longer core lengths (core length of 220-280

306    SNPs). The assignment yield (b) was optimal for a core length of 200 SNPs.

307

308    **Table 4 Percentages of alleles correctly assigned a breed origin (%Correct),**

309    **incorrectly assigned (%Incorrect), missing or unassigned (%Unassigned), and**

310    **accuracy of assignment (Accuracy) for each core-length (Core)**

| Core | %Correct | %Incorrect | %Unassigned | Accuracy |
|------|----------|-----------|-------------|----------|
| 100  | 94.70    | 1.35      | 3.95        | 0.99     |
| 120  | 94.89    | 1.12      | 3.99        | 0.99     |
| 140  | 96.14    | 1.11      | 2.75        | 0.99     |
| 160  | 96.40    | 1.16      | 2.44        | 0.99     |
| 180  | 97.39    | 1.25      | 1.35        | 0.99     |
| 200  | 98.35    | 1.36      | 0.29        | 0.99     |
| 220  | 97.27    | 1.46      | 1.27        | 0.99     |
| 240  | 94.22    | 1.54      | 4.24        | 0.98     |
| 260  | 92.32    | 1.69      | 5.98        | 0.98     |
| 280  | 95.93    | 1.84      | 2.23        | 0.98     |
| Mean | 95.76    | 1.39      | 2.85        | 0.99     |

311

312  **Consensus allele assignment across all cores**

313  On consensus, the average percentage of incorrectly assigned alleles was nearly zero

314  (Table 5). The overall mean consensus-based assignment accuracy (accuracy = 1,

315  Table 5) was higher than the overall mean core-based assignment accuracy (accuracy

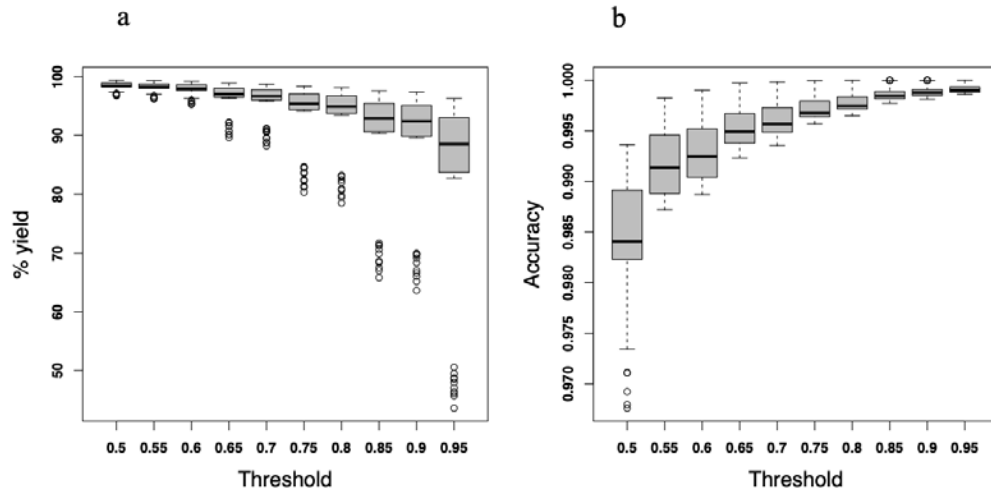316  = 0.99, Table 4).

317

318  **Table 5 Consensus-based percentages of alleles correctly assigned (%Correct),**

319  **incorrectly assigned (%Incorrect), missing or unassigned (%Unassigned) a breed**

320  **origin, and accuracy of assignment (Accuracy) across all core-lengths and**

321  **generation for each threshold**

| Threshold | %Correct | %Incorrect | %Unassigned | Accuracy |
|-----------|----------|------------|-------------|----------|
| 0.50 | 98.40 | 1.60 | 0.00 | 0.98 |
| 0.55 | 98.16 | 0.78 | 1.06 | 0.99 |
| 0.60 | 97.84 | 0.66 | 1.50 | 0.99 |
| 0.65 | 96.26 | 0.42 | 3.32 | 1.00 |
| 0.70 | 95.79 | 0.36 | 3.86 | 1.00 |
| 0.75 | 93.56 | 0.25 | 6.19 | 1.00 |
| 0.80 | 92.94 | 0.20 | 6.86 | 1.00 |
| 0.85 | 89.10 | 0.12 | 10.78 | 1.00 |
| 0.90 | 88.41 | 0.10 | 11.48 | 1.00 |
| 0.95 | 81.67 | 0.08 | 18.26 | 1.00 |
| Mean | 93.21 | 0.46 | 6.33 | 1.00 |

322

323  **Effect of admixture level and thresholds on assignment yield and**

324  **accuracy**

325  The threshold level had the opposite effect on assignment yield and accuracy (Fig. 6).

326  Increasing the threshold decreased the assignment yield and increased the accuracy,

327  whereas a less stringent threshold generated higher assignment yields. Increasing the

328  threshold stringency further reduced the assignment yield (Fig. 6a). On the contrary

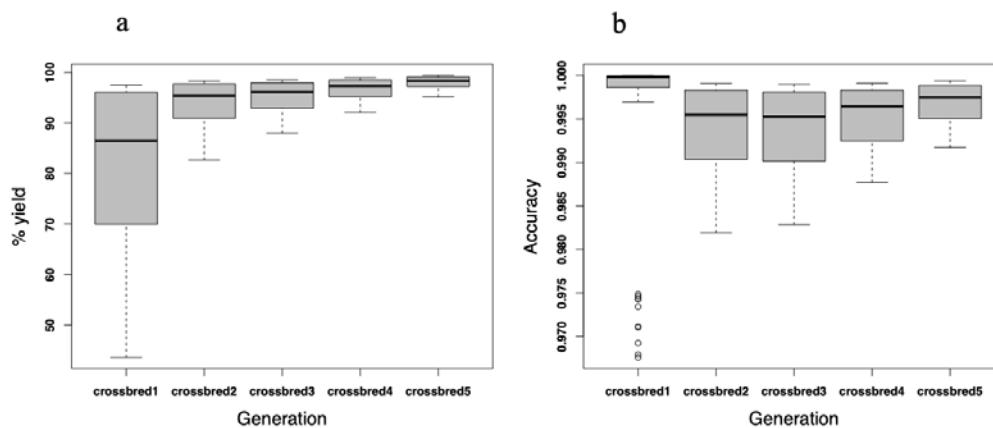329  and as expected, the less stringent threshold reduced the accuracy (Fig. 6b).

330

**Figure 6 Percentage of allele assignment yield (a) and accuracy (b) of assignment.** Using the consensus-based allele assignment algorithm as a function of threshold level

334

The effect of admixture level on assignment yield and accuracy was not as clear as that of threshold level. However, the assignment yield appeared to increase from the first to the later generations of crossbreds (Fig. 7a). On the other hand, the higher threshold stringency decreased the assignment yield (Fig. 7b).



339

340 **Figure 7 Percentage of allele assignment yield (a) and accuracy (b) of**

341 **assignment.** Using the consensus-based allele assignment algorithm as a function of

342 crossbreeding (admixture) level.


# Discussion

344 In low- and middle-income countries (LMICs), such as those in Eastern Africa, a

345 large proportion of dairy production is carried out by smallholders who keep fewer

346 than 10 cattle [14]. These cattle are mostly crosses between indigenous African breeds

347 and exotic dairy breeds, with little phenotypic or pedigree data available [14]. Despite

348 the need and efforts to increase the productivity of those dairy cattle, it has not been

349 possible to implement conventional breeding programmes in these populations. In

350 populations with no or poor pedigree and phenotype records, genomic selection and

351 other novel methods, such as an efficient algorithm to assign the breed origin of

352 alleles in those crossbred animals, are of interest. To evaluate the performance and

353 adaptability of the crossbreds in the LMICs, methods to accurately identify the breed

354 origin of alleles on both the individual level and the individual genetic variant are

355 important. Such methods could also provide ways to predict the effectiveness of

356 foreign germplasm in a low-input production system [4]. For the smallholder farmers

357 in Eastern Africa, providing methods to assign a breed origin of alleles would enable

358 better choice of exotic bulls to introduce and which local bulls to use to sustainably

359 harness the genetics of local adaptation traits of the indigenous breeds and the high

360 milk yield potential of exotic dairy breeds.

361

362 Different genomic tools and algorithms [5,9] have been developed to assign a breed

363 origin to alleles in crossbred pig populations without needing pedigree records. Using

364 simulated genotype data, we have developed an algorithm to assign alleles a breed

365 origin in a dairy cattle breeding programme that would represent haphazardly

366 admixed local cows and imported exotic bulls as commonly practised in LMICs. As

367 shown in Fig. 1, we used the exotic bulls as a source of purebred genotype data to

368 cross with the admixed local cows for five subsequent generations. The simulated

369 genotypes of exotic purebred and local admixed breeds were phased and the origins of

370 haplotypes and associated alleles of the newly created crossbred cows were assigned a

371 breed origin. In agreement with earlier studies in crossbred pig populations [5,9], our

372 results demonstrated that alleles of admixed crossbred cattle populations could be

373 accurately assigned a breed origin without the need for pedigree records.

374

375 The assignment of alleles to a breed origin was performed according to haplotypes

376 defined by different core lengths. In a simulation study, Vandenplas et al. [5] assessed

377 the impact of core length and observed higher assignment yield for haplotypes of

378 longer core lengths. While this appears to be supported in our results, a core- and tail-

379 length of 200 SNPs was observed as the optimal length for maximum assignment

380 yield. Similarly, the impact of genetic relationship on assignment yield is comparable

381 to values reported in simulated and empirical studies. Using simulated data,

382 Vandenplas et al. [5], showed that a greater distance between breeds favourably

383 affected the percentage of allele assigned, which is consistent with the highest

384 percentage of allele assignment yield observed in crossbred5 (97%, Table 3) that are

385 relatively distant to the local pure breeds.

386

387 The accuracy of allele assignment, both in the core-based (0.99, Table 3) and

388 consensus-based (1.00, Table 4), across all scenarios was very high. This allele

389 assignment accuracy is better than the results obtained from simulated (0.98) and

390 empirical (88.57- 92.45) data [9]. The performance of the current algorithm is better

391 than reported allele assignment accuracies of 96% using STRUCTURE 2.2 and 85%

392 using GENECLASS 2 reported by Negrini et al. [15]. The relative performance

393 improvement could be attributed to the optimization process of developing the current

394 allele assignment algorithm. For example, the breed origin of alleles in crossbred

395 animals was determined after an allele assignment was evaluated for every core and

396 haplotype library in different scenarios to reach a consensus assignment. The choice

397 of threshold for best SNP match in haplotypes can also affect the algorithm's

398 assignment yield and accuracy. Instead of using fixed allele frequency and best SNP

399 matches to assign a breed origin to alleles, the observed expected trade-offs between

400 assignment yield and accuracy (Fig. 6) have been optimized. When the best SNP

401 match counts in haplotypes are too low, there will be a high assignment yield but low

402 accuracy and vice versa. In the current simulated genotype data, the best SNP match

403 count threshold of 50-60% appeared optimal.

404

405 Despite some suggestions to use haplotype instead of allele to reduce the effects of

406 incorrect allelic assignments [5], the current algorithm was able to assign a breed

407 origin to alleles as accurate as the assignment of a breed origin to haplotypes. The

408 developed algorithm can be used to determine a breed origin of alleles in genomic

409 predictions with models where breed-specific effects are required [16,17]. The

410 developed algorithm can also be used in modelling breeding programmes of admixed

411 populations. Accurate breed identification, on both the level of the individual and of

412 the individual genetic variant is critical to achieving sustainable breed improvement.

413 In the current simulation study, we developed an algorithm, which assigns haplotypes

414  in crossbred dairy cows to the haplotypes of likely constituent breeds, i.e. either to

415  exotic or local breeds. With high accuracy of assigning the breed of origin to alleles,

416  we may be able to introduce a resilient or adaptive haplotype into the crossbred cows.

417  In livestock, we infer haplotypes from multigenerational pedigrees from which tracing

418  of breed origin of alleles can be challenging. With the developed algorithm, alleles in

419  crossbred animals could be accurately assigned a breed of origin without the need for

420  a multigenerational pedigree.

421

422  It's important to acknowledge that the African dairy cattle populations are

423  characterized by extensive crossbreeding involving many breeds of Taurine and

424  Indicine origin. This broad genetic diversity may challenge the accurate estimation of

425  SNP effects despite the accurate assignment of breed origin of alleles. While the BOA

426  method relies on the recent local ancestry for each SNP marker allele, it ignores

427  deeper ancestry, which is important for estimating SNP marker effects across many

428  breeds with different genomic histories. Furthermore, the BOA method does not take

429  full advantage of linkage information (correlation between nearby SNP markers) and

430  does not fully reflect the underlying genomic history of a study population [18].

431  Future studies developing algorithms and methods that consider the BOA and the

432  genomic history of individuals and that would work for any level of crossbreeding

433  and admixture in a population will be needed.

434

## Conclusions

436  The developed algorithm assigns a breed origin to alleles with an accuracy of 99% in

437  admixed animals from a crossbreeding programme designed to mimic breeding

438  programmes in the LMICs. The algorithm is straightforward in its application and

439 does not require prior knowledge of pedigree and relationships between crossbred and

440 purebred animals, making it relevant and applicable in breeding programmes

441 practised in LMICs. However, it should be noted that the algorithm was developed

442 and tested on simulated data. Further studies are required to test and apply the

443 algorithm on real data.

444 **List of abbreviations**

445 ADGG: African Dairy Genetic Gains

446 BOA: Breed Origin of Allele

447 CTLGH: Centre for Tropical Livestock Genetics and Health

448 LMICs: Low- and middle-income countries

449 MaCS : Markovian Coalescent Simulator

450 SNP: Single Nucleotide Polymorphism

451 # Declarations

452 **Ethics approval and consent to participate**

453 Not applicable

454 **Consent for publication**

455 Not applicable

456 **Availability of data and materials**

457 The scripts for data simulation and algorithm development are available [See

458 Additional file 2, Script S2], [See Additional file 3, Script S3] and [See Additional

459 file 4, Script S4].

460 # Competing interests

24

461     RCG and JMH are now employed by Bayer Crop Science.

# 474   Authors' contributions

475     BW, RCG, IH and JMH designed the study. BW performed the analyses and drafted

476     the manuscript. IH has substantively revised the manuscript, addressed all the

477     comments from co-authors and submitted the manuscript. GG and JMH supervised

478     the study and contributed to the manuscript. All authors read and approved the final

479     manuscript.

484

## References

486  1. Leroy G, Baumung R, Boettcher P, Scherf B, Hoffmann I. Review: sustainability of
487  crossbreeding in developing countries; definitely not like crossing a meadow. Animal.
488  2016;10**:**262-73.
489
490  2. Marshall K, Gibson JP, Mwai O, Mwacharo JM, Haile A, Getachew T, et al.
491  Livestock genomics for developing countries - African examples in practice. Front
492  Genet. 2019; 10:297.
493
494  3. VanRaden PM, Cooper TA. Genomic evaluations and breed composition for
495  crossbred US dairy cattle. Interbull Bull. 2015; 49: 19-23.
496
497  4. Kuehn LA, Keele JW, Bennett GL, McDaneld TG, Smith TPL, Snelling WM, et al.
498  Predicting breed composition using breed frequencies of 50,000 markers from the US
499  Meat Animal Research Center 2,000 Bull Project. J Anim Sci. 2011; 89**:**1742-50.
500
501  5. Vandenplas J, Calus MPL, Sevillano CA,Windig JJ, Bastiaansen JWM. Assigning
502  breed origin to alleles in crossbred animals. Genet Sel Evol. 2016; 48:61.
503
504  6. Sevillano CA, Vandenplas J, Bastiaansen JWM, Bergsma R, Calus MPL. Genomic
505  evaluation for a three-way crossbreeding system considering breed-of-origin of
506  alleles. Genet Sel Evol. 2017; 49:75.
507
508  7. Guillenea A, Lund MS, Evans R, Boerner V, Karaman E. A breed-of-origin of
509  alleles model that includes crossbred data improves predictive ability for crossbred
510  animals in a multi-breed population. Genet Sel Evol. 2023; 55:34.
511
512  8. Eiríksson JH, Karaman E, Su G, Christensen OF. Breed of origin of alleles and
513  genomic predictions for crossbred dairy cows. Genet Sel Evol. 2021; 53:84.
514
515  9. Sevillano CA, Vandenplas J, Bastiaansen JWM, Calus MPL. Empirical
516  determination of breed-of-origin of alleles in three-breed cross pigs. Genet Sel Evol.
517  2016; 48:55.
518
519  10. Gaynor RC, Gorjanc G, Hickey JM. AlphaSimR: an R package for breeding
520  program simulations. G3. 2021; 11:2.
521
522  11. Chen GK, Marjoram P, Wall JD. Fast and flexible simulation of DNA sequence
523  data. Genome Res. 2009; 19**:**136-42.
524
525  12. R: a language and environment for statistical computing. Vienna: R Foundation
526  for Statistical Computing; 2017.
527
528  13. Hickey JM, Kinghorn BP, Tier B, Wilson JF, Dunstan N, van der Werf JHJ. A
529  combined long-range phasing and long haplotype imputation method to impute phase
530  for SNP genotypes. Genet Sel Evol. 2011; 43:12.

531

532 14. Brown A, Ojango J, Gibson J, Coffey M, Okeyo M, Mrode R. Short
533 communication: genomic selection in a crossbred cattle population using data from
534 the dairy genetics East Africa project. J Dairy Sci. 2016; 99:7308-12.

535

536 15. Negrini R, Nicoloso L, Crepaldi P, Milanesi E, Colli L, Chegdani F, et al.
537 Assessing SNP markers for assigning individuals to cattle populations. Anim Genet.
538 2009; 40**:**18-26.

539

540 16. Christensen OF, Madsen P, Nielsen B, Su GS. Genomic evaluation of both
541 purebred and crossbred performances. Genet Sel Evol. 2014; 46:23.

542

543 17. Duenk P, Calus MPL, Wientjes YCJ, Breen VP, Henshall JM, Hawken R, et al.
544 Estimating the purebred-crossbred genetic correlation of body weight in broiler
545 chickens with pedigree or genomic relationships. Genet Sel Evol. 2019; 51:6.

546

547 18. Fan C, Mancuso N, Chiang CWK. A genealogical estimate of genetic
548 relationships. Am J Hum Genet. 2022; 109: 812-24.

549

## Additional files

550 **Additional files**

551 Additional file 1 Script S1

552 File format: .txt

553 Title: R Scripts to simulate the genotypes.

554 Description: The Additional file 1 describes the scripts to simulate genotypes of the

555 crossbred cattle.

556

557 Additional file 2 Script S2

558 Title: R scripts to phase the simulated genotypes

559 File Format: .txt

560 Description: Additional file 2 describes the scripts to phase the simulated genotypes.

561

562 Additional file 3 Script S3

563 Title: R scripts to assign a breed of origin to alleles of crossbred cattle

564 File Format: .txt

565 Description: Additional file 3 describes the scripts to assign a breed of origin to

566 alleles of crossbred cattle

567

568 Additional file 4 Script S4

569 Title: R scripts for consensus allele assignment

570 File Format: .txt

571 Description: Additional file 4 describes the Scripts for consensus allele assignment.