1      **CRUMBLER: A Tool for the Prediction of Ancestry in Cattle**

2

3      Tamar E. Crum[1], Robert D. Schnabel[1,2], Jared E. Decker[1,2], Luciana CA Regitano[3], and

4                       Jeremy F. Taylor[1]

5

6      [1]Division of Animal Sciences, University of Missouri, Columbia, MO, USA 65211

7      [2]Informatics Institute, University of Missouri, Columbia, MO, USA 65211

8      [3]Embrapa Pecuária Sudeste, São Carlos, SP 13560-970, Brazil

9

10

11      Corresponding author: Jeremy F. Taylor. E-mail: taylorjerr@missouri.edu.

12

13      Short Title: Predicting the Ancestry of Cattle

14

15      E-mail addresses:

16

17      TEC:    tamar.crum@mail.missouri.edu

18      RDS:    schnabelr@missouri.edu

19      JED:    deckerje@missouri.edu

20      LCAR: luciana.regitano@embrapa.br

21      JFT:    taylorjerr@missouri.edu

22

## Abstract

23

24 *Background*

25 In many beef and some dairy production systems, crossbreeding is used to take

26 advantage of breed complementarity and heterosis.  Admixed animals are frequently

27 identified by their coat color and body conformation phenotypes, however, without

28 pedigree information it is not possible to identify the expected breed composition of an

29 admixed animal and in the presence of selection, the actual composition may differ from

30 expectation.  As the roles of DNA and genotype data become more pervasive in animal

31 agriculture, a systematic method for estimating the breed composition (the proportions

32 of an animal's genome originating from ancestral pure breeds) has utility for a variety of

33 downstream analyses including the estimation of genomic breeding values for

34 crossbred animals, the estimation of quantitative trait locus effects, and heterosis and

35 heterosis retention in advanced generation composite animals.  Currently, there is no

36 automated or semi-automated ancestry estimation platform for cattle and the objective

37 of this study was to evaluate the utility of extant public software for ancestry estimation

38 and determine the effects of reference population size and composition and number of

39 utilized single nucleotide polymorphism loci on ancestry estimation. We also sought to

40 develop an analysis pipeline that would simplify this process for members of the

41 livestock genomics research community.

42 *Results*

43 We developed and tested a tool, "CRUMBLER", to estimate the global ancestry of cattle

44 using ADMIXTURE and SNPweights based on a defined reference panel.

45 CRUMBLER, was developed and evaluated in cattle, but is a species agnostic pipeline

46    that facilitates the streamlined estimation of breed composition for individuals with

47    potentially complex ancestries using publicly available global ancestry software and a

48    specified reference population SNP dataset.  We developed the reference panel from a

49    large cattle genotype data set and breed association pedigree information using

50    iterative analyses to identify purebred individuals that were representative of each

51    breed.  We also evaluated the numbers of markers necessary for breed composition

52    estimation and simulated genotypes for advanced generation composite animals to

53    evaluate the precision of the developed tool.

54    *Conclusion*

55    The developed CRUMBLER pipeline extracts a specified subset of genotypes that is

56    common to all current commercially available genotyping platforms, processes these

57    into the file formats required for the analysis software, and predicts admixture

58    proportions using the specified reference population allele frequencies.

59

60    **Background**

61

62    Estimation of the breed composition of individuals with complex ancestries has utility for

63    estimating breed direct and heterosis effects as well as for the estimation of the additive

64    genetic merit of these individuals. It also has value for identifying the breed composition

65    of training populations used for genomic selection and hence the identification of target

66    breeds in which the developed prediction equations may have some relevance.  Visual

67    classification of cattle based on breed characteristics suffers from similar problems as

68    the self-identification of ethnicity in humans [1], as most breed characteristics are

69  determined by alleles at relatively few loci. For example, recent extensive crossing with

70  Angus cattle in the U.S. produces a black hided animal which masks all other solid coat

71  colors found in other breeds and requires only a single dominant allele at the *MC1R*

72  locus. As a consequence, black-hided cattle have a "cryptic" population structure [1,2]

73  and the visual classification of black-hided animals for branded beef programs can

74  result in the marketing of animals with vastly different Angus genome content.

75

76  In the U.S. and many other countries, the breed of an animal is associated with its being

77  registered with a breed association which requires that both parents of the animal be

78  identified and also registered with the association.  For the previous 50 years,

79  parentage has been validated by each breed association using blood or, more recently,

80  DNA typing. Many breed associations have closed herdbooks which means, in theory,

81  that the pedigrees of all animals can be traced back to the animals that founded the

82  breed's herdbook. Other breed associations have open herdbooks, which means that

83  crossbred animals can be registered with the breed if they have been graded up by

84  crossbreeding to purebred status with the expectation that a certain percentage of their

85  genome (e.g., 15/16ths) originates from the respective breed based upon pedigree

86  records and parentage validation. Pedigree errors that occurred prior to, or that were

87  not identified following the implementation of blood typing and DNA testing, lead to

88  admixed animals being incorrectly classified as fullblood and incorrectly identified

89  admixture proportions in purebred animals. The effects of recombination, random

90  assortment of chromosomes into gametes and selection can also lead to considerable

91  variation in the extent of identity by descent between relatives separated by more than a

4

92    single meiosis and can also lead to admixture proportions that differ substantially from

93    expectation based on pedigree.

94

95    Crossbreeding is extensively used in commercial beef production and in other livestock

96    species production systems to capitalize on the effects of breed complementarity and

97    heterosis resulting in herds of females that may have very complex ancestries that

98    frequently use fullblood or purebred bulls sourced from registered breeders. Changes in

99    the decision as to which breed of bull to use can result in large changes in admixture

100   proportions of replacement cows and marketed steers between years and large

101   differences can occur between herds for the same reason. When commercially sourced

102   animals are used to generate resource populations to study the genomics of

103   economically important traits such as feed conversion efficiency [3,4] or bovine

104   respiratory disease [5], the presence of extensive admixture in the phenotyped and

105   genotyped animals may impact the GWAA [3,4] and leads to the training of genomic

106   prediction models in populations for which the breed composition is not understood.  As

107   a consequence, the utility of these models in other industry populations, including the

108   registered breeds in which the majority of genetic improvement is generated is also not

109   understood.

110

111   As the number of genotyped beef animals has increased, the need to classify the breed

112   composition of these animals has necessitated the development of a precise and

113   accurate method for estimating breed composition in cattle based on single nucleotide

114   polymorphism (SNP) data.  Iterative ancestry estimation analyses performed using

115    different software input parameters may identify those that cause output sensitivity and

116    can lead to an interpretation of population structure that is close to the truth [6].  We

117    developed the CRUMBLER analysis pipeline to streamline the genomic estimation of

118    breed composition of crossbred cattle using high-density SNP genotype data, publicly

119    available software, and a reference panel containing genotypes for members of cattle

120    breeds that are numerically important in North America. The CRUMBLER pipeline is

121    species agnostic and could be adapted for breed composition estimation in other

122    species.  CRUMBLER and the reference panel data are available on GitHub

123    (https://github.com/tamarcrum/CRUMBLER).  This pipeline tool is released under the

124    GNU General Public License.

125

126    **Materials and Methods**

127

128    **Genotype data**

129    From among the numerically most important cattle breeds in North America, in terms of

130    their annual numbers of animal registrations, a list was compiled to define the target

131    breeds for reference panel development.  Composite breeds, such as Brangus and

132    Braford, were not included in this list due to lack of available genotype data, but the

133    progenitor Angus, Hereford and Brahman breeds were included. Breeds such as

134    N'Dama, representing African taurine, and Nelore and Brahman, representing *Bos*

135    *taurus indicus* cattle, were included.  We also initially included breeds that were likely to

136    be involved in early crossbreeding of cattle in the U.S. (Texas Longhorn).

137

138    From the 170,544 cattle with high-density SNP genotypes stored within the University of

139    Missouri Animal Genomics genotype database, we extracted genotypes for 48,776

140    animals identified as being registered with one of the numerically important U.S. Breed

141    Associations or belonging to other world breeds. Pedigree data were also obtained for

142    these animals from each of the Breed Associations, where available (Table 1). These

143    individuals had been genotyped using at least one of 9 different genotyping platforms

144    currently used internationally to genotype cattle including the GeneSeek (Lincoln, NE)

145    GGP-90KT, GGP-F250, GGP-HDV3, GGP-LDV1, GGP-LDV3, and GGP-LDV4 assays,

146    the Illumina (San Diego, CA) BovineHD and BovineSNP50 assays, and the Zoetis

147    (Kalamazoo, MI) i50K assay.  The numbers of variants queried by each assay and the

148    number of individuals genotyped using each platform are shown in Table 2.

149

150    **Marker set determination**

151    To maximize the utility of the developed breed assignment tool, we identified the

152    intersection set of SNP markers located on the bovine assays for which we had

153    available genotype data (Table 2).  However, during the process of identifying the

154    animals that would define the breed reference panel, only 16 individuals had been

155    genotyped using the GGP-LDV4 (n=2) and GGP-LDV3 (n=14) assays and no animals

156    had been genotyped using the GGP-LDV1 assay.  To retain as many SNP markers as

157    possible for subsequent analysis, we identified the intersection of markers present on

158    the GGP-90KT, GGP-F250, GGP-HDV3, GGP-LDV3, GGP-LDV4, BovineHD,

159    BovineSNP50 and i50K assays.  This intersection set included 6,799 SNP markers

160    (BC7K).  The intersection of the markers representing 5 assays (GGP-90KT, GGP-

161    F250, GGP-HDV3, BovineHD, and BovineSNP50) was 13,291 markers (BC13K). By

7

162    removing only the 16 individuals from the breed reference panel that had been

163    genotyped on the GGP-LDV3 and GGP-LDV4 assays, we were able to compare

164    ancestry predictions using two marker set densities (BC13K and BC7K).

165

166    **Pipeline**

167    The developed CRUMBLER pipeline integrates the tools and the computational

168    efficiency of publicly available software, PLINK [7,8], EIGENSOFT [9,10] and

169    SNPweights [11] to generate ancestry estimates (Fig. 1).  The pipeline integrates the

170    often cumbersome processes of data reformatting and sequentially processing the data

171    using analytical tools to generate ancestry proportions for targeted individuals based on

172    a curated breed reference panel.

173

174    **PLINK**

175    PLINK PED formatted genotypes are required as input to the pipeline. PLINK

176    (v1.90b3.31) was used for data filtering and formatting.  Genotypes can arise from any

177    of the common bovine genotyping platforms (Table 2), provided that a PLINK

178    compatible MAP file is provided for each assay and data produced using only a single

179    genotyping assay is included in each PED file. The pipeline utilizes the PLINK marker

180    filtering tool (--extract) to extract the user-specified marker subset for ancestry analysis.

181    For analyses of animals genotyped on different genotyping platforms, the marker list

182    representing the intersection of the platforms can be provided to extract the markers

183    that are common to all assays. The pipeline allows multiple input genotype files and

184    uses the PLINK merge genotype files tool (--merge) to combine genotypes into a single

185    file for downstream analysis.

186

187    **EIGENSOFT**

188    The EIGENSOFT convertf package is used to convert all genotypes from PLINK PED

189    format into EIGENSTRAT format which is required by the SNPweights software.  To

190    process the reference panel data, principal component analysis using EIGENSOFT

191    smartpca is used to generate the eigenvalues and eigenvectors that are required to

192    calculate SNP weights using SNPweights.  However, the smartpca package included in

193    EIGENSOFT versions beyond 5.0.2 is not compatible with SNPweights.  SNPweights

194    requires an input variable, "*trace*", to be located in the log file output from the smartpca

195    analysis.  For versions of EIGENSOFT beyond 5.0.2, the source code can be edited to

196    ensure that the log file output is compatible with the SNPweights software (See

197    Supplementary Information).

198

199    **SNPweights**

200    SNPweights implements an ancestry inference model based on genome-wide SNP

201    weights computed using genotype data for an external panel of reference individuals.

202    To obtain SNP weights, the matrix ($g_{ij}$) of reference panel genotypes for SNP **i=1, ..., M**

203    and individual **j=1, ..., N** is normalized by subtracting the mean $\mu_i = N^{-1}\sum_j g_{ij}$ and

204    dividing by the standard deviation $[p_i(1-p_i)]^{0.5}$ for each SNP, where $p_i = \mu_i/2$, to improve

205    the results of the subsequent PCA analysis from which a kinship matrix is generated

206    [15]. A principal component decomposition is then used to generate the eigenvalues

9

207  and corresponding eigenvectors of the kinship matrix [11]. The SNP weights file only

208  needs to be recalculated if the reference panel is changed. EIGENSTRAT formatted

209  target animal genotypes are input into SNPweights, along with the precomputed

210  reference panel SNP weights. The SNP weights are then applied to the target

211  individuals to estimate their ancestry proportions [11].

212

213  **Reference panel development**

214  The definition of a set of reference individuals that define the genotype frequencies at

215  each SNP variant for each reference breed is technically demanding, but vitally

216  important to the process of defining ancestry. This process assumes that selection has

217  not operated to change gene frequencies between target and reference population

218  animals, and that each population is sufficiently large that drift has not impacted allele

219  frequencies. It also assumes that migration between different countries does not

220  influence population allele frequencies when registered animals are imported or

221  exported. FastSTRUCTURE [12] analysis and iterations of animal filtering using

222  SNPweights was performed using the genotypes of candidate reference panel

223  individuals to remove individuals with significant evidence of admixture from the

224  reference breed panel. An overview of the processes and iterations of filtering

225  conducted in the development of this reference panel set is shown in Fig. S1 and Table

226  1.

227

228  **FastSTRUCTURE analysis to identify candidate reference panel individuals**

229    Genotype data for 48,776 individuals produced by one of 8 different genotyping assays

230    were available for fastSTRUCTURE analysis (Table 1) [12]. We initially performed

231    focused fastSTRUCTURE analyses using small numbers of reference breeds including

232    Angus and Simmental, Angus and Gelbvieh, Angus and Limousin, Angus and Red

233    Angus, Red Angus, Hereford, Shorthorn and Salers, Red Angus, Hereford and

234    Shorthorn, and N'Dama, Nelore and Brahman (Figs. S2-S8).  Individuals possessing an

235    ancestry assignment of at least 97% to their designated breed were retained for

236    subsequent analysis (see Supplementary Methods and Table 1).  Following filtering

237    based on fastSTRUCTURE breed assignment, 17,852 individuals representing 19 of the

238    original breeds remained for further analysis (Supplementary Methods and Figs. S2-

239    S8).  All of the Salers animals were removed in this filtering analysis which is consistent

240    with previous work that found that Salers and Limousin were very similar [4]. Variation in

241    reference population sample sizes has been shown to substantially influence the

242    estimation of the number of ancestral populations (K) in ancestry analyses [6,13,14].  To

243    minimize this effect and produce similar sample sizes for each of the reference breeds,

244    we randomly sampled 200 individuals from each reference breed for which at least 200

245    individuals remained after filtering on an ancestry assignment of at least 97%, otherwise

246    all remaining individuals were included for the breed (Table 1).  Following

247    fastSTRUCTURE analysis using K=19 after removal of Salers and using the BC7K

248    marker set, Texas Longhorn was also removed from the reference panel breed list due

249    to the inability to distinguish Texas Longhorn as a distinct population (Figure 2).

250    Further, due to the known common ancestry [15] and similarity between Nelore and

251    Brahman (Figure 2), the breeds were combined to represent *Bos taurus indicus*.

11

252

**SNPweights analyses to refine and validate reference panel members**

254 Random sampling of reference breed individuals was performed to create sample sets

255 containing ≤n individuals per breed, for n = 50, 100, 150 and 200 individuals (Figs. 3 a-b

256 and Figs. S9-S10).  Sampling was performed such that if a reference breed had ≥n

257 candidates then n individuals were randomly sampled, otherwise, all available

258 individuals were sampled. An analysis was performed using the BC7K marker set,

259 SNPweights was used to assign reference breed ancestries to the same sample of

260 individuals that was used to produce the SNP weights for each of the four samples of

261 individuals (Figs. 3 a-b and Figs. S9-S10). In the self-assignment analyses conducted

262 using the reference breed sample sets of ≤100 individuals per breed and ≤50 individuals

263 per breed, 7 individuals were removed due to their estimated breed ancestry being

264 ≤60% to their registry breed (Holstein n=3, Jersey n=1, Japanese Black n=3) (Figs. 3 a-

265 b).

266

**Breeds with open herdbooks**

268 For the Gelbvieh, Limousin, Shorthorn, Simmental, and Braunvieh breeds that have

269 open U.S. herdbook registries, fullblood or 100% ancestry individuals were identified

270 based on pedigree data obtained from the respective breed associations (Table 1).  The

271 term "fullblood" is used to identify cattle for which every ancestor is registered in the

272 herdbook and can be traced back to the breed founders. The term "purebred" refers to

273 animals that have been graded up via crossbreeding to purebred status.  Charolais also

274 has an open herdbook registry in the U.S., however, access to Full French imported

275    Charolais breed members was limited.  As a result, all Charolais individuals identified as

276    purebred in the association registry were retained for downstream analysis, however,

277    these individuals could contain up to 1/32 introgression from another breed.  A random

278    sample of 200 individuals was taken for each breed with more than 200 identified

279    fullblood individuals, otherwise all animals were sampled.  Individuals previously

280    included in the candidate reference panel following preliminary fastSTRUCTURE

281    filtering for the open herd book breeds were removed and replaced with the fullblood

282    individuals.

283

284    **Additional reference panel filtering using SNPweights**

285    After filtering animals identified to not be fullblood based on their pedigree information,

286    we randomly sampled ≤50 individuals per reference breed and utilized SNPweights to

287    estimate weights for each sample and also to estimate breed ancestries for members of

288    the same sample that was used to generate the SNP weights.  Based on these

289    analyses, we created 5 overlapping reference breed sets, each containing individuals

290    with ≥90%, ≥85%, ≥80%, ≥75%, or ≥70% ancestry assignment to their registry breeds

291    (Table 3).

292

293    **Simulated Genotypes**

294    Using the phased BC7K genotypes for the final reference population of 803 individuals

295    (3 Nelore genotyped with the BovineHD assay were removed because they were

296    determined to cause problems for the phasing software), we simulated genotypes for

297    803 individuals each generation (N = 1, 3, 5 and 10) by randomly sampling two

13

298    individuals as parents from generation N-1 and using a Poisson distribution to sample at

299    random a single recombinant chromosome from each parent. The number of

300    recombination events for each sampled chromosome was sampled from a Poisson

301    distribution with mean equal to chromosome length in Mb/100 (i.e. 1.58 Morgans for

302    chromosome 1).  Simulated genotypes were produced for individuals 1 generation

303    removed from the fullblood/purebred reference population animals (i.e., 50% breed A

304    and 50% breed B), 3, 5, and 10, generations, respectively, to evaluate the ability of

305    CRUMBLER to detect large through to small admixture proportions in animals with

306    increasing numbers of breeds represented in their ancestry.  Breed composition

307    estimates for these animals were obtained by tracing the breed of origin of every allele

308    present in each generation N animal. For each marker, we attributed the genomic

309    fragment from the center points of the intervals on each side of each marker to the

310    breed of origin of the two alleles at each marker and summed these across all loci.

311    Finally, we normalized these sums by dividing by the autosomal genome size using

312    UMD3.1 coordinates.

313

314    **Results and Discussion**

315

316    The concept of breed and breed membership is man-made and does not inherently

317    exist in nature.  Moreover, the formation of breeds of cattle is very recent, as cattle

318    domestication began about 10,000 years ago but the formation of herdbooks has

319    occurred only during the last 200-250 years [16]. Nevertheless, the effects of drift and

320    human selection over the last 200 years have caused sufficient divergence among

14

321    breeds that breed differences are identifiable at the molecular level.  Such signals are

322    essential for breed ancestry analyses to be effective in modern admixed animals.

323    Previous work on assigning breed composition in admixed cattle utilized 50K genotype

324    data and a reference panel of 16 breeds, with the basis for reference panel inclusion

325    being breed association registration [17].  However, the continual evolution of

326    genotyping assays has led to content changes resulting in only a relatively small

327    proportion of markers in common among assays. Consequently, there is a need to

328    evaluate whether these markers are sufficient for breed content estimation, leading to

329    their conservation in the design of future assays. Furthermore the development of an

330    analytical pipeline based on these markers would simplify analysis for end-users and

331    the use of a single reference panel would allow the direct comparison of results

332    between applications.

333

334    **Reference panel development**

335    Previously developed cattle reference panels have relied on pedigree accuracy and

336    breed association registration for their definition [17].  Conversely, we used an iterative

337    approach for reference population curation that was able to validate the accuracy of the

338    pedigree information used to identify candidates.  FastSTRUCTURE analyses

339    performed using the candidate individuals for each of the initial 19 reference breeds

340    suggested population subdivision in both the Hereford and Simmental (Figure 2).

341    Pedigree analysis for the Herefords within each subpopulation indicated that the

342    subpopulations comprised animals from the highly inbred USDA Miles City Line 1

343    Hereford population (L1) and other individuals representing broader U.S. Hereford

15

344     pedigrees.  The Miles City L1 Hereford cattle were derived from two bulls, both sired by

345     Advance Domino 13 (AHA registration number 1668403) and 50 Hereford foundation

346     cows.  Since the founding of the L1 Herefords, the migration of germplasm has been

347     unidirectional from L1 into the broader U.S. industry, as the L1 population has been

348     closed since its founding [18].  However, the L1 Herefords have profoundly influenced

349     the U.S. Hereford population.  L1 Herefords do not segregate for recessive dwarfism,

350     which has been a threat to Hereford breeders since the 1950s, and this has led to L1

351     cattle becoming popular in the process of purging herds of the defect [19].  In 1980, the

352     average proportion of U.S. registered Herefords influenced by L1 genetics was 23%.

353     By 2008, this proportion had increased to 81% [18].

354

355     The detected subpopulation division within the Simmental breed (Figure 2) represents

356     the differentiation between purebred and fullblood animals.  For example, progeny of a

357     popular fullblood Simmental sire are present in both subpopulations, however, in one

358     subpopulation the family members are all fullblood and in the other they are all purebred

359     or percentage Simmental animals.  This result supports the need to identify fullblood

360     animals as reference panel breed representatives for breeds with open herdbooks.

361

362     Reference population sample size

363     By randomly sampling individuals from the candidate reference breed set and using

364     SNPweights to assign these individuals to reference populations, we found that

365     reference panel breed sample sizes of ≤50 or ≤100 individuals appeared to capture the

366     diversity within each breed and appropriately determined the ancestry of the tested

16

367    individuals (Fig. 3 a-b).  For each breed, the percent ancestry predicted for the tested

368    reference samples was, on average, 3.86% higher when the SNP weights were

369    estimated using ≤50 individuals per breed than when ≤100 individuals per breed were

370    used (Table 4).  This reflects the increased homogeneity of individuals within each

371    breed and a greater genetic distance between individuals from different breeds as

372    smaller samples of individuals from each breed are used to define the reference panel.

373    Further, due to limitations in the number of genotyped individuals for some breeds

374    (Table 1), as the sample size was increased globally, imbalances were created between

375    the reference panel breed sample sizes which impacted breed composition estimation

376    (Fig. S9-S10). It has previously been shown that the power to detect population

377    structure improves as the reference population sample sizes become more similar

378    [6,14].

379

380    <u>Marker density</u>

381    After the replacement of reference breed individuals with those identified to be fullblood

382    based on pedigree analysis for the open herdbook Gelbvieh, Simmental, Limousin,

383    Braunvieh, Shorthorn, and Charolais breeds, additional self-assignment analyses were

384    conducted to evaluate the effects of marker set size on ancestry prediction.  Breed

385    reference panels were again constructed by randomly sampling ≤50 individuals per

386    breed and SNP weights were calculated using both the BC13K markers and BC7K

387    markers.  The estimated SNP weights were then used to self-assign ancestry to

388    members of the reference panel animals representing the reference breed set.  The

389    ancestry predictions for the reference breed individuals using either the BC7K (Fig. 4a;

390   Fig. S11) or BC13K (Fig. 4b; Fig. S12) marker sets indicate that use of the BC13K

391   marker set did not significantly impact the ancestry predictions.  Consequently, the use

392   of the 6,799 markers common to the 8 commercially available genotyping platforms

393   appears to be sufficient to assign breed ancestry for the majority of animals produced in

394   the U.S.  The CRUMBLER pipeline can accommodate samples genotyped using

395   alternative assays, however, the produced breed composition estimates will be based

396   on the intersection of markers on the assay and the BC7K marker set.

397

398   Assignment thresholds

399   We next examined the effects of reference breed homogeneity on ancestry assignment

400   by identifying reference panel members that had been assigned to their breed of

401   registry using SNPweights with probabilities of ancestry of ≥90%, ≥85%, ≥80%, ≥75%,

402   and ≥70%, respectively (Table 3).  From these individuals, reference breed panels were

403   obtained by randomly sampling ≤50 individuals per breed, until each individual was

404   represented in at least one sample set.  SNP weights were then estimated using the

405   BC7K marker set and ancestry was assigned for these individuals using SNPweights

406   (Figs. 5-6 and Figs. S13-S15). Limiting the reference breed panel members to those

407   individuals with ≥90% ancestry assigned to their breed of registry produced a reference

408   panel that did not represent the extent of diversity within each of the breeds (Fig. 5).  On

409   the other hand, using an ancestry assignment of ≥85% clearly captured greater diversity

410   within each breed (Fig. 6) and maximized the self-assignment of ancestry to the breed

411   of registration (Table 5).

412

413    Reference panel definition

414    To examine whether the specific individuals represented in the reference panel sample

415    influenced the self-assignment of ancestry to the sampled individuals, a second sample

416    of ≤50 distinct individuals per breed was obtained from the individuals with ≥85%

417    assignment to their breed of registration and analyzed with SNPweights (Fig. 7). Fig. 7

418    indicates that the ability to predict ancestry was not influenced by the specific individuals

419    sampled from the set of animals with ≥85% ancestry to their breed of registration.

420

421    Additionally, Figs. 6 and 7 suggest that the use of a reference breed panel constructed

422    by the random sampling of ≤50 individuals per breed from individuals with ≥85% self-

423    assigned ancestry to their breed of registration maintained sufficient within-breed

424    diversity to accurately estimate the ancestry of target individuals. However, these

425    figures also reveal small amounts of apparent introgression from other reference panel

426    breeds within each of the breeds. This does not appear to be an issue of marker

427    resolution since the analyses performed with the BC7K and BC13K marker sets

428    generated similar results (Fig. 4). We conclude that these apparent introgressions are

429    either due to a lack of power to discriminate among breeds using the common markers

430    designed onto commercial genotyping platforms, or represent the presence of common

431    ancestry among the breeds prior to the formation of breed herdbooks ~200 years ago.

432    Molecular evidence for this shared ancestry exists, for example, Hereford and Angus

433    cattle share the *Celtic* polled allele [20] and the segmental duplication responsible for

434    the white anterior, ventral and dorsal coat color pattern occurs only in Hereford and

19

435    Simmental cattle and their crosses [21]. These data clearly indicate that crossbreeding

436    was widespread prior to the formal conceptualization of breeds.

437

438    **Reference panel validation**

439    To evaluate the ability of the selected reference breed panel to identify breed

440    composition, an analysis was conducted for all 170,544 samples in the database which

441    required 60 processor minutes (Fig 6-7).  We extracted animals with pedigree

442    information including fullblood and purebred animals registered with open herdbook

443    breed associations and 2,243 crossbred animals with varying degrees of admixture.

444    Considering the amount of available data, the number of pedigreed admixed animals

445    was very limited and the purebred animals all had similar expected admixture

446    proportions. Consequently, we next simulated genotypes for animals by assuming the

447    random mating of members of the reference breed panel for 1, 3, 5 and 10 generations

448    assuming non-overlapping generations to generate generations of animals with different

449    numbers of breeds and breed proportions represented in their genomes.

450

451    Registered fullblood animals

452    For the Gelbvieh, Limousin, Shorthorn, Simmental, and Braunvieh breeds that have

453    open herdbook registries, fullblood or 100% ancestry individuals were identified based

454    on pedigree data obtained from the respective breed associations (Table 1).

455    CRUMBLER estimates were obtained for these fullblood individuals and the distribution

456    of estimates by breed are in Fig. 8.  For all breeds except Charolais, >50% of the

457    individuals had CRUMBLER estimated percentages of ≥80% to their respective breeds.

458     Average percentage estimates for fullblood Gelbvieh, Limousin, Shorthorn, Simmental,

459     and Braunvieh individuals were 76%, 78%, 83%, 79%, and 85%, respectively (Fig. 8b).

460     However, the number of genotyped imported Full French Charolais animals was limited

461     and so we also analyzed all purebred Charolais individuals which could contain up to

462     $1/32^{nd}$ of their genome introgressed from another breed. The average Charolais breed

463     assignment was 72% and the distribution of estimates was more variable than for the

464     fullblood animals from the other breeds (Fig. 8b).

465

466     <u>Pedigreed crossbred animals</u>

467     Based on pedigree, 2,005 individuals were identified as being primarily Hereford but

468     with varying degrees of Red Angus, Salers, Angus or unknown other breed influence.

469     The analysis results agreed with the pedigree data (Fig. 9a) To investigate the

470     correlations between pedigree and CRUMBLER estimated breed proportions, we

471     removed proportions for breeds that were less than 3% and normalized the remaining

472     values. CRUMBLER estimates were then correlated with the pedigree predicted

473     estimates of the proportion of Hereford in these individuals (Fig. 9c). CRUMBLER

474     tended to underestimate the Hereford proportion as the pedigree estimated Hereford

475     proportion tended to 100%.

476

477     The remaining 238 crossbred individuals were commercial, advanced generation

478     animals with an expected 50% Angus and 50% Simmental ancestry based on pedigree

479     data. Results of the CRUMBLER analysis again support the pedigree data (Fig. 10).

480     The presence of Red Angus ancestry in these animals reveals the inability of the

481    analysis to fully differentiate between Angus and Red Angus, which only diverged in the

482    U.S. in 1954, and also the influence of Red Angus in the U.S. Simmental breed (Fig.

483    S16).

484

485    <u>Simulated genotypes</u>

486    Genomes were simulated using the phased genotypes for 803 individuals from the

487    reference breed panel to contain varying breed numbers and admixture proportions

488    after 1, 3, 5, and 10 generations of random mating with nonoverlapping generations. In

489    generation 1, the admixed individuals were $F_1$ individuals with a 50:50 autosomal

490    genome composition unless both parents were randomly sampled from the same breed.

491    CRUMBLER estimates of breed composition using the simulated genotypes were

492    strongly correlated with the simulated compositions, especially for generations 1 and 3

493    (Fig 11).  As the number of generations increased, the number of breeds represented in

494    the simulated genomes tended to increase and the proportion of the genome originating

495    from any one breed tended to decrease and the correlation between the simulated

496    proportions and CRUMBLER estimates also decreased. Nevertheless, by generation 10

497    44% of animals had their genome proportions estimated with a correlation of at least

498    70%. In the U.S. commercial crossbreeding does not usually involve the use of more

499    than 3-4 breeds of cattle and while the number of generations of crossbreeding may

500    very well be 10 or perhaps more, many generations will involve the mating of animals

501    with similar genome ancestries and the proportions for each breed will be much greater

502    than present in the generation 10 animals in Fig 11. Consequently, the achieved

503    accuracies are likely to be closer to the generation 3 or 5 results where 99% and 68% of

504 animals, respectively, had their genome proportions estimated with a correlation of

505 greater than 80%.

506

507 <u>Advanced generation composite animals</u>

508 The ancestry model assumes that neither drift or selection has acted to alter the allele

509 frequencies from those created by the initial admixture proportions. We examined

510 CRUMBLER estimates of breed composition for advanced generation members of the

511 Brangus (n=11,362), Beefmaster (n=3,832) and Santa Gertrudis (n=2,010) composite

512 breeds where selection has had the opportunity to change breed composition from

513 expectations at breed formation.  Brangus individuals are expected to be ⅝ Angus and

514 ⅜ Brahman, Beefmaster individuals ¼ Hereford, ¼ Shorthorn and ½ Brahman, and

515 Santa Gertrudis ⅝ Shorthorn and ⅜ Brahman, respectively.  These breeds use mating

516 strategies that produce individuals that are expected to possess these proportions for

517 registration within each of the respective breed's herdbook.  However, registerable

518 animals are ultimately advanced generation composites and so drift, meiotic sampling of

519 parental chromosomes and selection are all expected to create individual variation in

520 these ancestry proportions.  CRUMBLER results for these advanced generation

521 composites, also known as the American breeds, are shown in Figure 12.  Table 6

522 contains the average breed proportion estimates assigned to each of these breeds by

523 CRUMBLER and their standard deviations across the animals analyzed for each breed.

524 In every instance, CRUMBLER underestimates the expected proportions for each of the

525 American breed populations, however, the ancestral breeds clearly dominate the

526 assignments (Table 6).  Interestingly, on average, CRUMBLER estimated proportions of

23

527    Holstein ancestry for advanced generation Beefmaster and Brangus animals (Figure 12

528    and Table 6).  These American breeds do not contain any Holstein introgression and

529    they do not contain ancestry from a "Ghost Population", a population that is not present

530    in the reference set, which would lead to a breed assignment to a reference breed that it

531    most closely resembled [6].  We speculate that this effect is caused by selection

532    creating a deviation in allele frequencies from those found in the founder breeds which

533    the model explains by an introgression from a distantly related breed, in this case,

534    Holstein. Stratifying these genotyped animals according to the number of generations

535    from foundation fullblood animals and examining the extent of estimated Holstein

536    introgression, which would be expected to increase with generation number, would

537    enable this to be tested, but we did not have access to the necessary data. However,

538    this hypothesis is supported by the fact that the Santa Gertrudis had the least estimated

539    Holstein introgression and the breed has published estimates of additive genetic merit

540    for many fewer years than the Beefmaster or Brangus.

541

542    **Admixture**

543    We also tested the ADMIXTURE software [22] for ancestry estimation and integration

544    into the CRUMBLER pipeline using the same reference breed panel that was developed

545    for use with SNPweights.  ADMIXTURE uses maximum likelihood estimation to fit the

546    same statistical model as STRUCTURE, however, STRUCTURE does not allow the

547    specification of individuals of known descent to be used as a reference panel [22].

548    ADMIXTURE allows a supervised analysis, in which the user can specify a reference

549    set of individuals, by specifying the "--supervised" flag and requires an additional file

24

550     with a ".pop" suffix to specify the genotypes of the reference population individuals [22].

551     Unlike SNPweights, the reference population individuals' genotypes must be provided in

552     a genotype file for each analysis.

553

554     We first conducted an ADMIXTURE analysis in which we self-assigned ancestry for the

555     animals in the reference breed set formed with ≤50 individuals per breed from the

556     individuals that had ≥85% assignment to the breed of registration (Fig. 13). The results

557     shown in Fig. 13 are similar to those in Fig. 6 for the same reference panel, albeit with

558     perhaps less evidence of background introgression. We next conducted an analysis

559     using the reference panel used in Fig. 13 merged with data for the 2,005 high

560     percentage crossbred Herefords animals.  The results shown in Fig. 14, reveal a

561     significant change in the ancestry proportions estimated for the reference panel

562     Guernsey, Gelbvieh and Romagnola individuals between the two analyses which used

563     exactly the same reference panel, but differed only in the number of individuals for

564     which ancestry was to be estimated.  This suggests that ADMIXTURE may use the

565     target individuals to update information provided by the reference panel individuals

566     specified in the ".pop" file.  Consequently, the ADMIXTURE estimated ancestry

567     proportions appear to be context dependent and may vary based on the other

568     individuals included in the analysis.

569

570     Moreover, the order in which the target individuals appear in the genotype input file also

571     appears to affect ADMIXTURE estimates of ancestry proportions for the target

572     individuals.  Fig. 15 shows the results of an ADMIXTURE analysis in which the target

25

573  individuals were identical to those shown in Fig. 14, but for which the order of the

574  reference individuals and the 2,005 Hereford crossbred individuals was reversed in the

575  input files. In Fig. 14, the reference individuals appear before the 2,005 Hereford

576  crossbred individuals in the input file, whereas in Fig. 15, the 2,005 Hereford crossbred

577  individuals appeared before the reference individuals in the input file. The results reveal

578  a significant change in ancestry proportions for Guernsey and Gelbvieh, but the

579  Romagnola now appear to be non-admixed.  Finally, we performed an ADMIXTURE

580  analysis for these animals in which the order of animals in the input genotype file was

581  completely randomized (Fig. 16).  Following analysis, the individuals were sorted to

582  generate Fig. 16. Again, the ancestry proportions for the Guernsey, Gelbvieh and

583  Romagnola individuals suggest these breeds to be admixed.

584

585  STRUCTURE and ADMIXTURE are widely used for characterizing admixed populations

586  [6], however, we have not found any reports in the literature that indicate that the

587  software is sensitive to the input order of individuals.  However, we suspect that the

588  majority of users would have no need or motivation to run the software with permuted

589  data input files. Nevertheless, because of these inconsistencies between results, we

590  chose to not use ADMIXTURE for ancestry estimation within the CRUMBLER pipeline.

591

592  **Broader application using additional commercially available assays**

593  To broaden the spectrum of data from different commercially available assays that can

594  be evaluated, an additional intersection of markers was obtained using 11 commercially

595  available bovine assays including the GGP-90KT, GGP-F250, GGP-HDV3, GGP-LDV3,

596     GGP-LDV4, BovineHD, BovineSNP50, i50K, Irish Cattle Breeding Federation (Cork,

597     Ireland) IDBv3, and GeneSeek (Lincoln, NE) BOVG50v1 assays.  The intersection SNP

598     set included 6,363 SNPs (BC6K).  A SNPweights self-assignment analysis using the

599     reference set of individuals with ≥85% assignment to their breed of registration was

600     conducted to assess the effects of the reduction in number of markers used for ancestry

601     assignment.  The ancestry proportions assigned based on the BC6K marker set (Fig.

602     17) did not differ appreciably from those obtained using the BC7K marker set (Fig. 6).

603     This result indicates the utility of CRUMBLER and the reference panel breed set across

604     the spectrum of commercially available genotyping platforms.

605

## Conclusions

607     The determination of a set of reference population breeds and individuals that define

608     allele and genotype frequencies at each variant for each of the breeds is arguably the

609     most important, yet technically difficult step in the process of ancestry estimation.  We

610     employed several iterations of filtering to remove recently admixed individuals and

611     identify a relatively homogeneous set of individuals that nevertheless represented the

612     variation that might be expected among individuals within a breed.  Once determined,

613     the reference panel genotype data need only be processed once to obtain SNP weights

614     removing the need to share genotype data for reference individuals in subsequent

615     studies [11]. The upfront development of an external reference breed panel capitalizes

616     on the rich ancestry information available in large available datasets, and relatedness,

617     variation in sample sizes and diversity among the target individuals does not affect the

618     inference of ancestry [11].

27

619

620    In cattle, the visual evaluation of breed characteristics is a poor method for evaluating

621    the ancestry of individuals.  Breed association pedigrees can be used to estimate

622    expected breed compositions, however, the random assortment of chromosomes into

623    gametes and selection can lead to ancestry proportions that differ from those expected

624    based upon pedigree. Moreover, the vast majority of commercial beef cattle in the U.S.

625    have no or very limited pedigree information and since these animals are frequently

626    used for genomic research [3–5], there is a need for a tool that can routinely provide

627    ancestry estimates for downstream use in GWAA or other genetic studies.

628

629    We tested ADMIXTURE and SNPweights and found that results from ADMIXTURE

630    appear to depend on the ancestry and order of appearance of individuals within the

631    genotype input file. We therefore developed an analysis pipeline, CRUMBLER, based

632    upon PLINK, EIGENSOFT and SNPweights to automate the process of ancestry

633    estimation. The developed bovine pipeline utilizes the 6,799 SNPs present on 8

634    commercially utilized bovine SNP genotyping assays and results using these SNPs are

635    consistent with results obtained when 13,291 SNPs were used. From an available

636    48,776 genotyped individuals, we also developed a reference panel of 806 individuals

637    sampled from 17 breeds to have ≤50 individuals per breed that had ≥85% assignment

638    to their breed of registration. This panel appears to allow the robust estimation of the

639    ancestry of advanced generation admixed animals, however, all breeds share some

640    common ancestry which predates the recent development of breed association

641    herdbooks [16,23].

642

643    CRUMBLER is not limited to application in cattle and with the provision of suitable

644    reference breed allele frequencies can be applied to other species for ancestry

645    estimation.  CRUMBLER pipeline scripts and reference panel breed SNP weights are

646    available on GitHub (https://github.com/tamarcrum/CRUMBLER).

647 **Additional files**

648 **Supplementary Information (PDF).** This file contains the source code changes in

649 SMARTPCA within versions of EIGENSOFT beyond 5.0.2 to enable compatibility with

650 SNPweights.

651

652 **Supplementary Methods (PDF).** This file describes the preliminary fastSTRUCTURE

653 analyses conducted on subsamples of breeds in the development of the reference

654 breed panel.

655

656 **Supplementary Figures (PDF).** This file contains Supplementary Figures S1-S16.

657

658 **Fig. S1 An overview of the processes and iterations of filtering conducted in the**

659 **development of the reference panel.**

660

661 **Fig. S2 Preliminary FastSTRUCTURE analysis of candidate Angus and Simmental**

662 **reference population animals.**

663

664 **Fig. S3 Preliminary fastSTRUCTURE analysis of candidate Angus and Gelbvieh**

665 **reference population animals.**

666

667 **Fig. S4 Preliminary fastSTRUCTURE analysis of candidate Angus and Limousin**

668 **reference population animals.**

669

670 **Fig. S5 Preliminary fastSTRUCTURE analysis of candidate Angus and Red Angus**

671 **reference population animals.**

672

673 **Fig. S6 Preliminary fastSTRUCTURE analysis of candidate Red Angus, Hereford,**

674 **Shorthorn and Salers reference population animals.**

675

676 **Fig. S7 Preliminary fastSTRUCTURE analysis of candidate Red Angus, Hereford**

677 **and Shorthorn reference population animals.**

678

679 **Fig. S8 Preliminary fastSTRUCTURE analysis of candidate N'Dama, Nelore and**

680 **Brahman reference population animals.**

681

682 **Fig. S9 SNPweights self-assignment analysis for the reference sample set**

683 **containing ≤200 individuals per breed analyzed using the BC7K marker set.**

684

685 **Fig. S10 SNPweights self-assignment analysis for the reference sample set**

686 **containing ≤150 individuals per breed analyzed using the BC7K marker set.**

687

688 **Fig. S11 SNPweights self-assignment analysis for the reference sample set**

689 **containing ≤50 individuals per breed analyzed using the BC7K marker set.**

690

691 **Fig. S12 SNPweights self-assignment analysis for the reference sample sets**

692 **containing ≤50 individuals per breed analyzed using the BC13K marker set.**

693

**Fig. S13 SNPweights self-assignment analysis for the reference sample set with**

**≥80% ancestry to breed of registry and ≤50 individuals per breed using the BC7K**

**marker set.**

697

**Fig. S14 SNPweights self-assignment analysis for reference sample set with ≥75%**

**ancestry to breed of registry and ≤50 individuals per breed using the BC7K**

**marker set.**

701

**Fig. S15 SNPweights self-assignment analysis for the reference sample set with**

**≥70% ancestry to breed of registry and ≤50 individuals per breed using the BC7K**

**marker set.**

705

**Fig. S16 SNPweights self-assignment analyses using a reference panel with ≤50**

**individuals per breed and sampling from the individuals with ≥85% assignment to**

**their breed of registry but with (a) Red Angus or (b) Angus excluded from the**

**reference panel.**

710

**Authors' contributions**

TC conceived the study and managed the project. TC, RS, JD, and JT contributed to

defining the research questions and analytical approaches and interpretation of the

results. TC programmed the CRUMBLER pipeline and carried out the data analyses.

TC and JT drafted the manuscript.  LR provided the Nelore samples but did not have

716    involvement in the scientific direction. All authors read and approved the final

717    manuscript.

718

**Competing interests**

720    The authors declare that they have no competing interests.

721

**Availability of data and materials**

723    Project Name: CRUMBLER

724    Project Home Page: https://github.com/tamarcrum/CRUMBLER

725    Programming Language: Python

726    Other Requirements: PLINK, EIGENSOFT, and SNPweights

727    License: GNU GPL

728

**Ethics approval and consent to participate**

730    Not applicable.

731

**Funding**

737

**Consent for Publication**

739    Not applicable.

740

741    **Acknowledgements**

742    Not applicable.

743

## References

744

745  1.  Burnett MS, Strain KJ, Lesnick TG, de Andrade M, Rocca WA, Maraganore DM.

746     Reliability of self-reported ancestry among siblings: implications for genetic

747     association studies. Am J Epidemiol. 2006 Mar 1;163(5):486–92.

748  2.  Pritchard JK, Stephens M, Donnelly P. Inference of population structure using

749     multilocus genotype data. Genetics. 2000 Jun;155(2):945–59.

750  3.  Saatchi M, Beever JE, Decker JE, Faulkner DB, Freetly HC, Hansen SL, et al.

751     QTLs associated with dry matter intake, metabolic mid-test weight, growth and feed

752     efficiency have little overlap across 4 beef cattle studies. BMC Genomics. 2014

753     Nov 20;15:1004.

754  4.  Seabury CM, Oldeschulte DL, Saatchi M, Beever JE, Decker JE, Halley YA, et al.

755     Genome-wide association study for feed efficiency and growth traits in U.S. beef

756     cattle. BMC Genomics. 2017 May 18;18(1):386.

757  5.  Neibergs HL, Seabury CM, Wojtowicz AJ, Wang Z, Scraggs E, Kiser JN, et al.

758     Susceptibility loci revealed for bovine respiratory disease complex in pre-weaned

759     holstein calves. BMC Genomics. 2014 Dec 22;15:1164.

760  6.  Lawson DJ, van Dorp L, Falush D. A tutorial on how not to over-interpret

761     STRUCTURE and ADMIXTURE bar plots. Nat Commun. 2018 Aug 14;9(1):3258.

762  7.  Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al.

763     PLINK: a tool set for whole-genome association and population-based linkage

764     analyses. Am J Hum Genet. 2007 Sep;81(3):559–75.

765     8.  Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-

766         generation PLINK: rising to the challenge of larger and richer datasets.

767         Gigascience. 2015 Feb 25;4:7.

768     9.  Patterson N, Price AL, Reich D. Population structure and eigenanalysis. PLoS

769         Genet. 2006 Dec;2(12):e190.

770     10. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal

771         components analysis corrects for stratification in genome-wide association studies.

772         Nat Genet. 2006 Aug;38(8):904–9.

773     11. Chen C-Y, Pollack S, Hunter DJ, Hirschhorn JN, Kraft P, Price AL. Improved

774         ancestry inference using weights from external reference panels. Bioinformatics.

775         2013 Jun 1;29(11):1399–406.

776     12. Raj A, Stephens M, Pritchard JK. fastSTRUCTURE: Variational Inference of

777         Population Structure in Large SNP Data Sets. Genetics. 2014 Jun 1;197(2):573–

778         89.

779     13. Wang J. The computer program structure for assigning individuals to populations:

780         easy to use but easier to misuse. Mol Ecol Resour. 2017 Sep;17(5):981–90.

781     14. Puechmaille SJ. The program structure does not reliably recover the correct

782         population structure when sampling is uneven: subsampling and new estimators

783         alleviate the problem. Mol Ecol Resour. 2016 May;16(3):608–27.

784    15. Sanders JO. History and Development of Zebu Cattle in the United States. J Anim

785        Sci. 1980 Jun 1;50(6):1188–200.

786    16. Decker JE, McKay SD, Rolf MM, Kim J, Molina Alcalá A, Sonstegard TS, et al.

787        Worldwide patterns of ancestry, divergence, and admixture in domesticated cattle.

788        PLoS Genet. 2014 Mar;10(3):e1004254.

789    17. Kuehn LA, Keele JW, Bennett GL, McDaneld TG, Smith TPL, Snelling WM, et al.

790        Predicting breed composition using breed frequencies of 50,000 markers from the

791        US Meat Animal Research Center 2,000 Bull Project. J Anim Sci. 2011

792        Jun;89(6):1742–50.

793    18. Leesburg VLR, MacNeil MD, Neser FWC. Influence of Miles City Line 1 on the

794        United States Hereford population. J Anim Sci. 2014 Jun;92(6):2387–94.

795    19. McCann LP. battle of bull runts. 1974; Available from: http://agris.fao.org/agris-

796        search/search.do?recordID=US201300539215

797    20. Wiedemar N, Tetens J, Jagannathan V, Menoud A, Neuenschwander S,

798        Bruggmann R, et al. Independent polled mutations leading to complex gene

799        expression differences in cattle. PLoS One. 2014 Mar 26;9(3):e93435.

800    21. Whitacre L. Structural variation at the KIT locus is responsible for the piebald

801        phenotype in Hereford and Simmental cattle. 2014; Available from:

802        http://search.proquest.com/openview/45eba5fa3c5757a2c4c2ab18af1a8a98/1?pq-

803        origsite=gscholar&cbl=18750&diss=y

804  22.  Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in

805      unrelated individuals. Genome Res. 2009 Sep;19(9):1655–64.

806  23.  Decker JE, Pires JC, Conant GC, McKay SD, Heaton MP, Chen K, et al. Resolving

807      the evolution of extant and extinct ruminants with high-throughput phylogenomics.

808      Proc Natl Acad Sci U S A. 2009 Nov 3;106(44):18644–9.

809

**Table 1. Genotype data for 48,776 registered individuals from 20 breeds were used to establish the reference population.**

| Breed | No. Registered Individuals | No. FullBlood Individuals[a] | No. Individuals Assigned to Breed[b] | Sampled Individuals[c] | No. Individual After Pedigre and SNPweights[c] |
|---|---|---|---|---|---|
| Angus | 5552 | 5552 | 485 | 200 | 200 |
| Hereford | 969 | 969 | 348 | 200 | 200 |
| Limousin | 2734 | 321 | 367 | 200 | 200 |
| Charolais | 1542 | 1489 | 1542 | 200 | 200 |
| Simmental | 15858 | 337 | 1583 | 200 | 196 |
| Japanese Black | 97 | 97 | 97 | 97 | 94 |
| Braunvieh | 148 | 69 | 148 | 148 | 69 |
| Gelbvieh | 12835 | 51 | 6000 | 200 | 51 |
| Romagnola | 37 | 37 | 37 | 37 | 37 |
| Salers | 68 | 68 | 0 | 0 | 0 |
| Texas Longhorn | 45 | 45 | 45 | 0 | 0 |
| Shorthorn | 291 | 178 | 166 | 166 | 178 |
| Red Angus | 1377 | 1377 | 124 | 124 | 124 |
| Holstein | 5816 | 5816 | 5816 | 200 | 197 |
| Jersey | 119 | 119 | 119 | 119 | 118 |
| Brown Swiss | 92 | 92 | 92 | 92 | 90 |
| Guernsey | 30 | 30 | 30 | 30 | 30 |
| N'Dama | 98 | 98 | 59 | 59 | 59 |
| Brahman | 127 | 127 | 86 | 86 | 50 |
| Nelore | 941 | 941 | 708 | 200 | 50 |
| Total | 48776 | 17813 | 17852 | 2558 | 2143 |

813    [a]Number of registered animals determined by pedigree analysis to be fullblood for breed
814    associations with open herdbooks.
815    [b]Number of registered animals assigned to their identified breed with $P \geq 0.97$ by
816    fastSTRUCTURE in preliminary analyses and retained for subsequent analyses.
817    [c]A random sample of 200 individuals was obtained for breeds with >200 individuals after
818    fastSTRUCTURE analysis and all individuals were sampled for breeds with ≤200 per breed and
819    the data were again analyzed by fastSTRUCTURE with K=19 after removal of the Salers.
820    [d]Animals that were determined to not be fullblood by pedigree analysis and animals
821    assigned with $P \leq 0.60$ by SNPweights to their breed of registry were removed.
822

823 **Table 2. The number of variants queried by each assay and the number of**
824 **individuals from the 20 reference breeds genotyped using each assay.**
825

| Assay | No. of Variants | No. of Registered Individuals |
|---|---|---|
| BovineSNP50 | 58336 | 20485 |
| BovineHD | 777962 | 2303 |
| GGP-F250 | 227234 | 3068 |
| GGP-90KT | 76999 | 4407 |
| GGP-LDV3 | 26504 | 6065 |
| GGP-HDV3 | 139977 | 3630 |
| GGP-LDV4 | 30105 | 8653 |
| GGP-LDV1 | 8762 | 165 |
| Zoetis i50K | 59825 | 0 |
| ICBF IDBv3 | 53450 | 0 |
| BOVGv1 | 47843 | 0 |
| Total | | 48776 |

826

**Table 3. Number of individuals for each reference breed assigned to their breed of registration by minimum ancestry threshold.**

| Breed | Breed Assignment Probability | | | | |
|---|---|---|---|---|---|
| | ≥90% | ≥85% | ≥80% | ≥75% | ≥70% |
| Angus | 51 | 136 | 184 | 199 | 200 |
| Hereford | 58 | 136 | 184 | 200 | 200 |
| Limousin | 93 | 127 | 144 | 162 | 173 |
| Charolais | 52 | 92 | 119 | 132 | 147 |
| Simmental | 21 | 43 | 81 | 103 | 121 |
| Japanese Black | 52 | 73 | 78 | 83 | 86 |
| Braunvieh | 37 | 57 | 63 | 65 | 68 |
| Gelbvieh | 23 | 31 | 39 | 43 | 43 |
| Romagnola | 10 | 25 | 32 | 36 | 37 |
| Shorthorn | 34 | 98 | 159 | 170 | 177 |
| Red Angus | 48 | 88 | 110 | 120 | 123 |
| Holstein | 39 | 119 | 172 | 193 | 196 |
| Jersey | 52 | 77 | 91 | 108 | 116 |
| Brown Swiss | 38 | 64 | 73 | 82 | 86 |
| Guernsey | 12 | 22 | 29 | 30 | 30 |
| N'Dama | 27 | 45 | 59 | 59 | 59 |
| Brahman | 15 | 40 | 50 | 50 | 50 |
| Nelore | 32 | 50 | 50 | 50 | 50 |
| Total | 694 | 1323 | 1717 | 1885 | 1962 |

42

831 **Table 4. Ancestry proportion statistics for the self-assignment of reference panel**
832 **members from samples of ≤50 or ≤100 individuals from the candidate reference**
833 **breed individuals.**

834

| Breed | Min % (≤50) | Avg % (≤50) | Max % (≤50) | Min % (≤100) | Avg % (≤100) | Max % (≤100) |
|---|---|---|---|---|---|---|
| Angus | 86.22 | 90.40 | 95.54 | 78.49 | 87.05 | 94.13 |
| Hereford | 79.75 | 90.08 | 95.05 | 73.41 | 87.39 | 96.81 |
| Limousin | 69.52 | 88.53 | 98.16 | 18.36 | 86.40 | 98.81 |
| Charolais | 78.14 | 90.19 | 99.82 | 48.93 | 77.46 | 93.96 |
| Simmental | 81.06 | 90.37 | 97.66 | 61.36 | 73.05 | 88.11 |
| Japanese Black | 81.44 | 90.00 | 97.07 | 24.51 | 86.50 | 98.95 |
| Braunvieh | 71.59 | 89.46 | 98.61 | 65.46 | 88.36 | 98.70 |
| Gelbvieh | 73.03 | 76.27 | 81.63 | 60.92 | 74.59 | 80.33 |
| Romagnola | 75.05 | 87.18 | 96.66 | 74.79 | 85.99 | 95.12 |
| Shorthorn | 84.42 | 88.69 | 94.54 | 70.71 | 85.27 | 96.35 |
| Red Angus | 79.00 | 89.60 | 96.33 | 68.07 | 86.83 | 97.38 |
| Holstein | 85.82 | 90.30 | 97.51 | 62.95 | 86.97 | 97.81 |
| Jersey | 78.55 | 89.28 | 95.93 | 61.23 | 86.54 | 97.18 |
| Brown Swiss | 80.10 | 89.22 | 96.40 | 61.68 | 86.02 | 98.42 |
| Guernsey | 79.53 | 89.19 | 95.85 | 77.40 | 88.31 | 94.36 |
| N'Dama | 80.67 | 89.25 | 96.90 | 78.91 | 87.78 | 95.67 |
| *Bos taurus indicus* | 87.83 | 91.91 | 97.75 | 81.43 | 89.79 | 97.60 |

835

**Table 5. Average predicted ancestry and variance in predicted ancestry for candidate reference breed individuals when filtered on minimum predicted ancestry.**

| Breed | Avg % (70%) | Var (70%) | Avg % (75%) | Var (75%) | Avg % (80%) | Var (80%) | Avg % (85%) | Var (85%) | Avg % (90%) | Var (90%) |
|---|---|---|---|---|---|---|---|---|---|---|
| Angus | 86.50 | 0.21 | 87.95 | 0.19 | 87.33 | 0.22 | 88.86 | 0.13 | 72.34 | 0.97 |
| Hereford | 86.99 | 0.22 | 87.09 | 0.23 | 87.48 | 0.19 | 88.25 | 0.13 | 84.62 | 0.43 |
| Limousin | 86.77 | 0.55 | 89.03 | 0.44 | 87.92 | 0.38 | 88.48 | 0.43 | 80.62 | 1.19 |
| Charolais | 80.18 | 2.16 | 85.03 | 1.77 | 86.28 | 0.99 | 88.56 | 0.52 | 81.54 | 0.76 |
| Simmental | 72.73 | 0.89 | 78.45 | 0.58 | 83.81 | 0.36 | 89.65 | 0.15 | 87.82 | 0.50 |
| Japanese Black | 87.85 | 0.52 | 88.04 | 0.39 | 88.46 | 0.27 | 88.74 | 0.21 | 80.06 | 0.61 |
| Braunvieh | 87.01 | 0.37 | 87.84 | 0.36 | 87.33 | 0.38 | 88.71 | 0.21 | 80.47 | 1.24 |
| Gelbvieh | 86.68 | 0.41 | 87.10 | 0.43 | 87.52 | 0.34 | 88.43 | 0.34 | 83.31 | 1.25 |
| Romagnola | 86.16 | 0.33 | 86.37 | 0.32 | 87.16 | 0.32 | 86.22 | 0.29 | 86.38 | 1.16 |
| Shorthorn | 85.97 | 0.26 | 87.03 | 0.22 | 86.80 | 0.14 | 87.38 | 0.07 | 83.00 | 0.70 |
| Red Angus | 86.41 | 0.53 | 87.08 | 0.48 | 87.40 | 0.35 | 87.46 | 0.23 | 23.37 | 0.66 |
| Holstein | 86.44 | 0.27 | 87.82 | 0.21 | 87.54 | 0.13 | 88.77 | 0.12 | 79.71 | 0.61 |
| Jersey | 87.01 | 0.46 | 86.93 | 0.44 | 87.86 | 0.24 | 87.98 | 0.27 | 80.52 | 0.71 |
| Brown Swiss | 86.22 | 0.47 | 86.73 | 0.51 | 88.24 | 0.26 | 88.11 | 0.20 | 82.23 | 0.70 |
| Guernsey | 86.46 | 0.23 | 87.64 | 0.19 | 87.50 | 0.25 | 88.02 | 0.51 | 80.43 | 2.36 |
| N'Dama | 87.76 | 0.19 | 87.91 | 0.21 | 87.89 | 0.15 | 89.25 | 0.17 | 86.40 | 0.52 |
| *Bos taurus indicus* | 87.68 | 0.07 | 88.24 | 0.09 | 87.55 | 0.11 | 88.53 | 0.09 | 84.89 | 0.38 |
| Average | 85.58 | 0.48 | 86.84 | 0.41 | 87.30 | 0.30 | 88.32 | 0.24 | 78.69 | 0.87 |

**Table 6. Average breed ancestry percentages assigned to American Breed individuals.**

| Breed | Avg. Ancestry Beefmaster % (± st. dev) | Avg. Ancestry Brangus % (± st. dev) | Avg. Ancestry Santa Gertrudis % (± st. dev) |
|---|---|---|---|
| Angus | 3.29 (± 4.27) | 32.15 (± 8.96) | 4.90 (± 4.48) |
| Hereford | 16.13 (± 2.83) | 2.03 (± 2.93) | 2.50 (± 4.05) |
| Limousin | 1.40 (± 2.28) | 1.73 (± 2.56) | 1.29 (± 2.19) |
| Charolais | 6.89 (± 3.97) | 2.07 (± 3.79) | 5.26 (± 3.42) |
| Simmental | 2.65 (± 3.12) | 1.16 (± 2.92) | 0.40 (± 1.40) |
| Japanese Black | 0.53 (± 3.46) | 0.10 (± 0.63) | 0.22 (± 0.89) |
| Braunvieh | 0.63 (± 1.64) | 0.33 (±1.29) | 0.59 (± 1.63) |
| Gelbvieh | 3.19 (± 3.30) | 3.14 (± 3.67) | 2.59 (± 3.20) |
| Romagnola | 1.05 (± 1.94) | 0.54 (± 1.39) | 0.68 (± 1.57) |
| Shorthorn | 15.36 (± 4.72) | 5.86 (± 3.42) | 37.71 (± 5.46) |
| Red Angus | 3.66 (± 3.57) | 13.60 (± 3.95) | 1.18 (± 3.46) |
| Holstein | 6.22 (± 6.73) | 4.53 (± 4.82) | 0.89 (± 2.83) |
| Jersey | 0.73 (± 1.65) | 0.52 (± 1.37) | 0.26 (± 1.08) |
| Brown Swiss | 1.05 (± 2.14) | 1.28 (± 2.26) | 0.73 (± 1.81) |
| Guernsey | 1.53 (± 2.20) | 0.17 (± 0.81) | 1.50 (± 2.14) |
| N'Dama | 0.52 (± 1.35) | 0.19 (± 0.87) | 0.16 (± 0.76) |
| *Bos taurus indicus* | 27.32 (± 4.84) | 23.09 (± 6.73) | 30.50 (± 4.52) |

**Fig. 1** Flow diagram of the breed composition pipeline.

**Fig. 2** FastSTRUCTURE results for a random sample of ≤200 individuals per breed from the pool of 17,852 potential reference individuals at K=19. Breed identification is shown below each colored block and each animal is represented as a vertical line within the block.

a.



b.

**Fig. 3** SNPweights self-assignment analysis results for reference panel sample sets consisting of: (a) ≤100 individuals per breed, or (b) ≤50 individuals per breed. Seven individuals were filtered for ≤60% ancestry to their breed of registry (Holstein n=3, Jersey n=1, Japanese Black n=3).

a.



b.

**Fig. 4** SNPweights self-assignment of ancestry for candidate reference breed individuals following evaluation of open herdbook breeds using: (a) the BC7K, or (b) the BC13K marker panels. Reference breed panels were constructed by random sampling ≤50 individuals per breed and SNP weights were estimated using the BC7K and BC13K marker sets.

**Fig. 5** Reference breed panel constructed by the random sampling of ≤50 individuals per breed from individuals with ≥90% ancestry was self-assigned to reference breed ancestry using the BC7K marker set.

**Fig. 6** Reference breed panel constructed by the random sampling of ≤50 individuals per breed from individuals with ≥85% ancestry was self-assigned to reference breed ancestry using the BC7K marker set.

**Fig. 7** Reference breed panel constructed by the independent random sampling of a second sample of ≤50 individuals per breed from individuals with ≥85% ancestry after eliminating individuals represented in the first sample was self-assigned to reference breed ancestry using the BC7K marker set.

a.



b.



**Fig. 8** (a) Distribution by breed of SNPweights ancestry assignment results for 2,408 registered fullblood animals from open herd book breeds. (b) Pictorial representation of CRUMBLER estimates for 2,408 registered fullblood animals from open herd book breeds.

**Fig. 9** (a) SNPweights ancestry results for 2,005 crossbred Hereford individuals with *a-priori* breed composition estimates determined by pedigree. (b) Breed assignment reference breed key. (c) Hereford SNPweights estimated proportions using CRUMBLER are plotted against the pedigree estimates. Data point color indicates the breed for which SNPweights assigned the highest proportion for each individual.

**Fig. 10** (a) SNPweights ancestry results for 238 crossbred individuals with *a-priori* breed composition estimates of 50% Angus and 50% Simmental based on a reference panel with ≤50 individuals per breed sampled from individuals with ≥85% assignment to their breed of registry. (b) Breed assignment for the crossbred individuals can be determined using this reference breed key.

**Fig. 11** Genotypes were simulated for the indicated number of generations of random mating, with generation 1 (G1) animals being 50:50 proportion except when two parents from the same breed were mated. SNPweights results were obtained using CRUMBLER pipeline parameters correlations between these estimates and the known simulated breed compositions were produced and the proportion of individuals within each correlation class is indicated.

a.



b.



**Fig. 12** (a) SNPweights ancestry results using CRUMBLER pipeline for 11,362 Brangus, 3,832 Beefmaster, and 2,010 Santa Gertrudis individuals. (b) Breed assignment for these advanced generation composite animals can be determined using this reference breed key.

**Fig. 13** Self-assignment of ancestry for the animals in the reference breed set formed with ≤50 individuals per breed from the individuals that had ≥85% assignment to their breed of registration using ADMIXTURE.
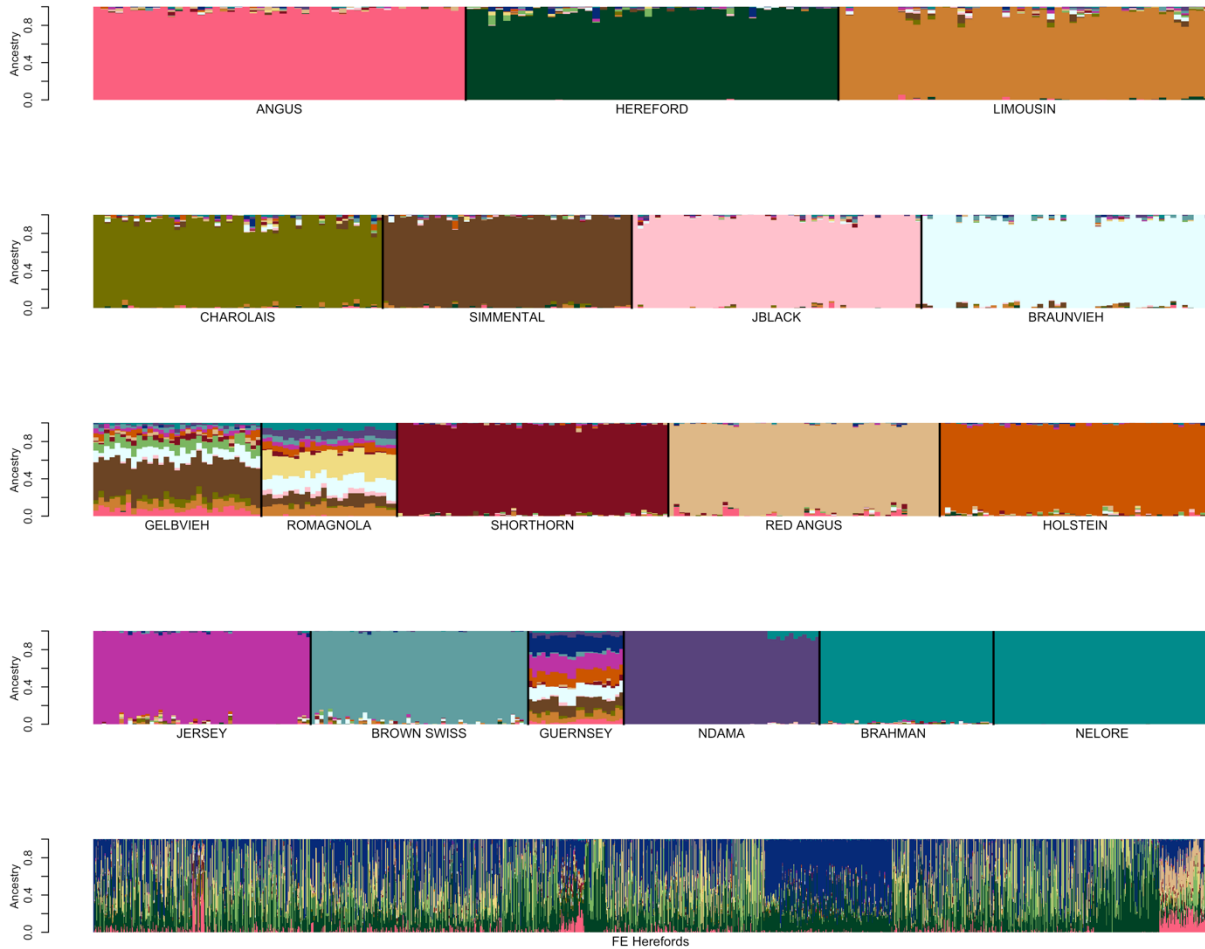
**Fig. 14** ADMIXTURE analysis conducted using the same data as shown in Figure 13 (first four rows), merged with an additional 2,005 high percentage crossbred Hereford target individuals (last row). Here, the 2,005 Hereford crossbred individuals appear after the reference individuals in the input genotype file.
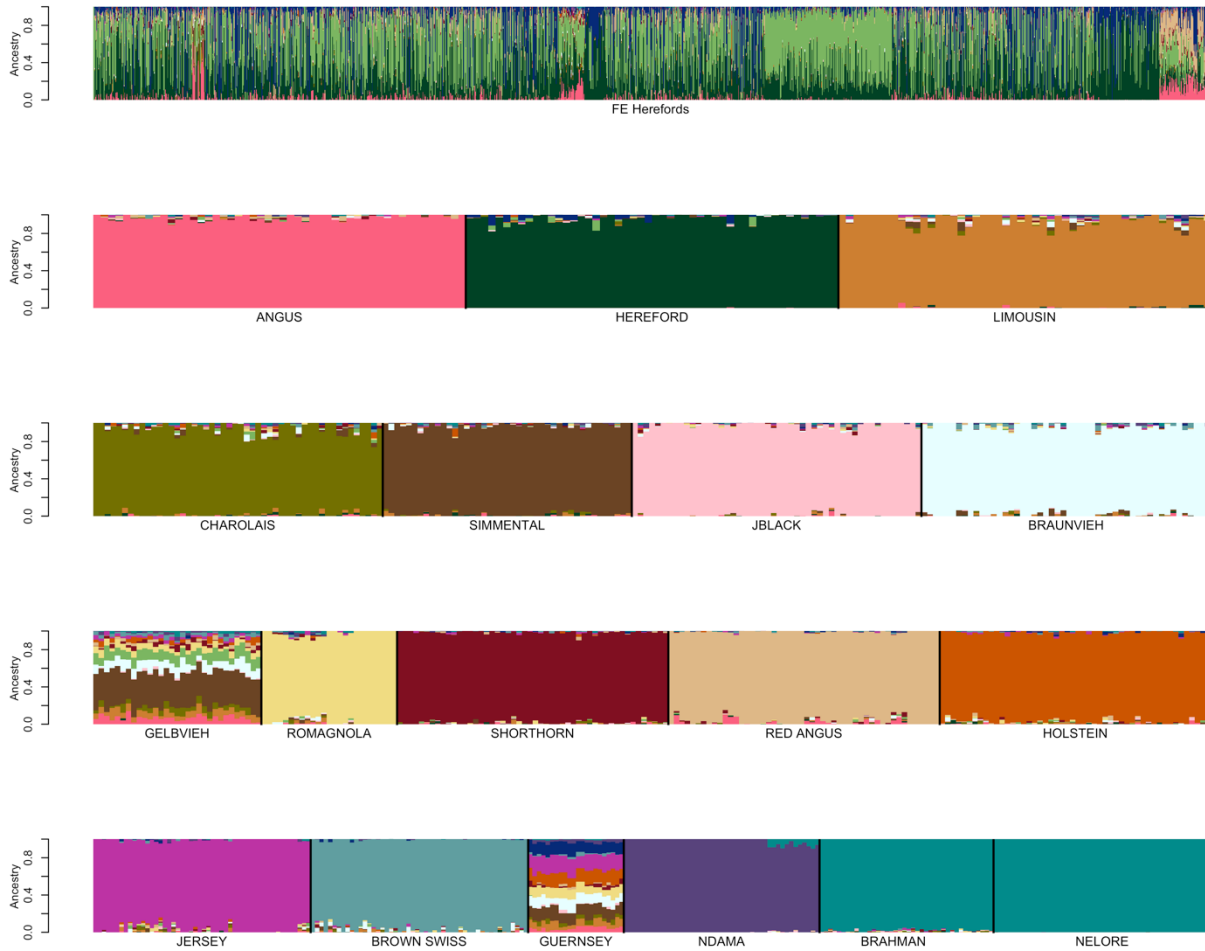
**Fig. 15** ADMIXTURE analysis conducted using the same data as shown in Fig. 14. Here, the 2,005 Hereford crossbred individuals appear before the reference individuals in the input genotype file. The first row represents the 2005 Hereford crossbred samples. Rows 2 to 5 show the reference panel individuals.
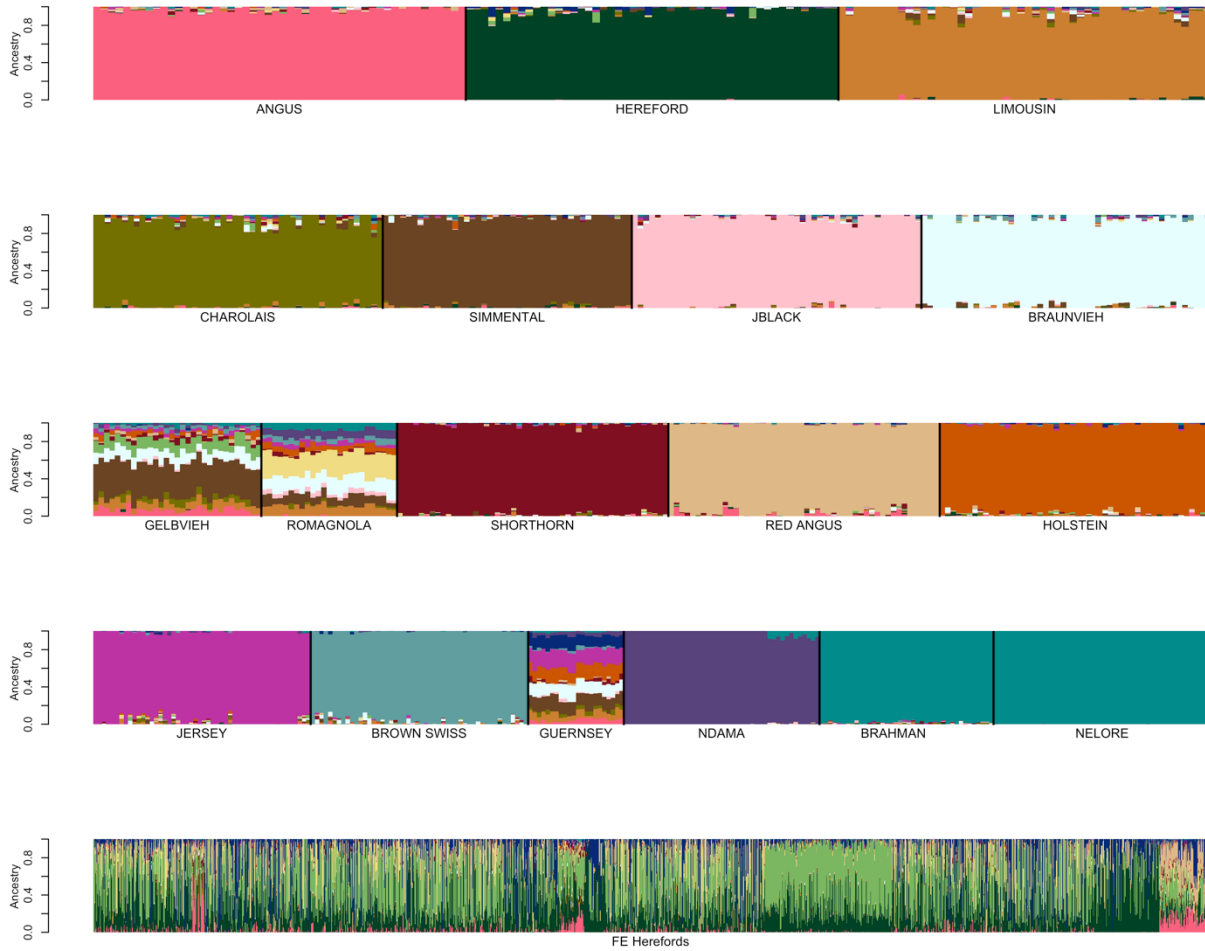
**Fig. 16** ADMIXTURE analysis conducted using the same data as shown in Figs. 14 and 15, but with the order of the individuals in the input genotype file randomized. The animals were sorted following analyses to generate this figure where the first four rows represent the reference panel individuals, the fifth row shows the 2,005 Hereford crossbred animals.
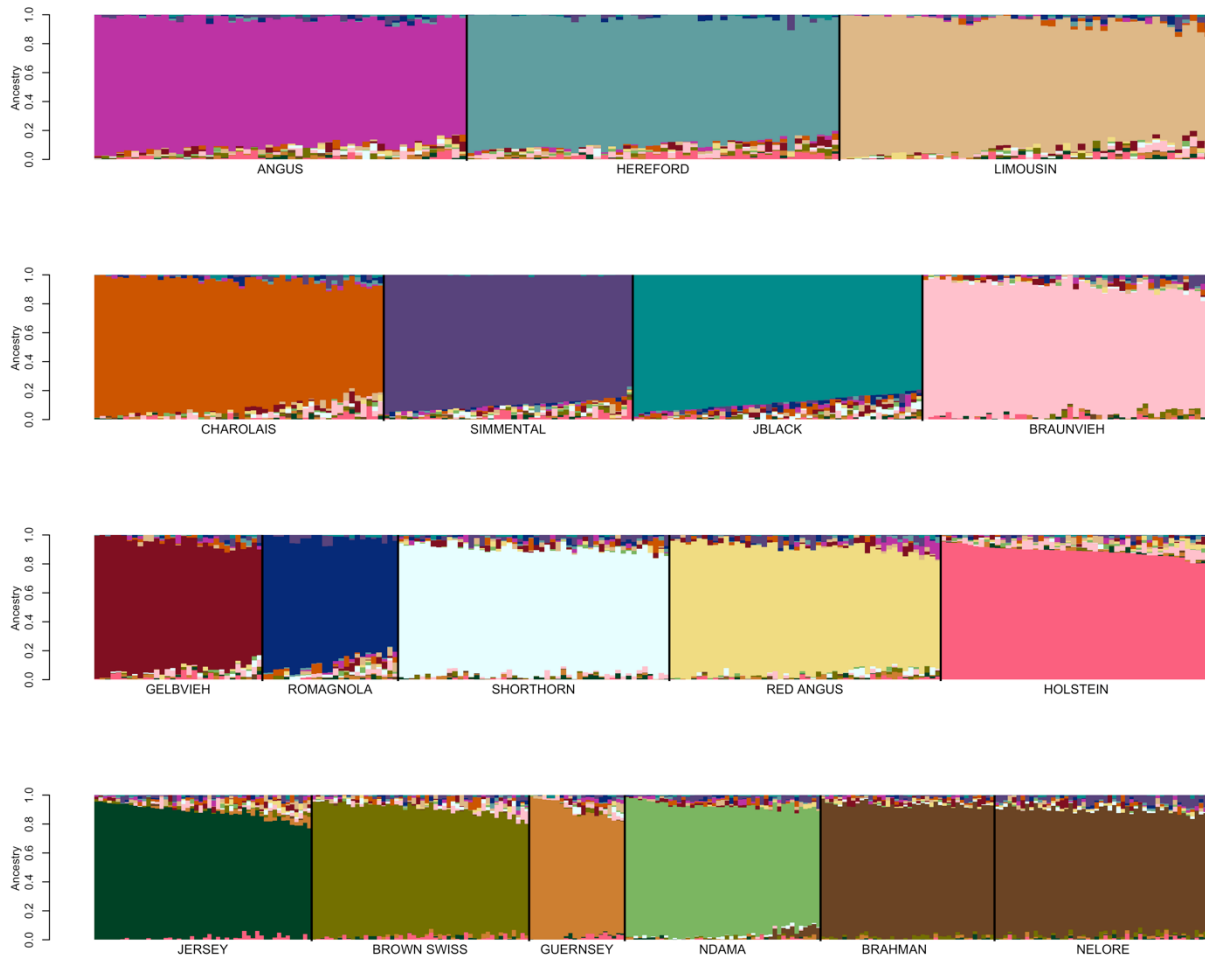
**Fig. 17** Reference breed panel constructed by the random sampling of ≤50 individuals per breed from individuals with ≥85% ancestry was self-assigned to reference breed ancestry using the BC6K marker set.