1 **Mutant alleles differentially shape cattle complex traits and fitness**

2 Ruidong Xiang[1,2,*], Ed J. Breen[2], Sunduimijid Bolormaa[2], Christy J. Vander Jagt[2], Amanda J.

3 Chamberlain[2], Iona M. Macleod[2], Michael E. Goddard[1,2]

4 [1] *Faculty of Veterinary & Agricultural Science, The University of Melbourne, Parkville 3052,*

5 *Victoria, Australia*

6 [2] *Agriculture Victoria, AgriBio, Centre for AgriBiosciences, Bundoora, Victoria 3083, Australia.*

7 [*]Corresponding Author: ruidong.xiang@unimelb.edu.au

8

9 **Abstract**

10 Classical mutant alleles (MAs), with large effects on phenotype, tend to be deleterious to

11 traits and fitness. Is this the case for mutations with small effects? We infer MAs for 8

12 million sequence variants in 113k cattle and quantify the effects of MA on 37 complex traits.

13 Heterozygosity for variants at genomic sites conserved across 100 vertebrates increase

14 fertility, stature, and milk production, positively associating these traits with fitness. MAs

15 decrease stature and fat and protein concentration in milk, but increase gestation length and

16 somatic cell count in milk (the latter indicative of mastitis). However, the allele frequency of

17 MAs that decrease fat and protein concentration and stature and increase gestation length and

18 somatic cell count is lower than the allele frequency of MAs with the opposite effect. These

19 results suggest bias in the direction of effect of mutation (e.g. towards reduced protein in

20 milk), but selection operating to reduce the frequency of these MAs. Taken together, our

21 results imply two classes of genomic sites subject to long-term selection: sites conserved

22 across vertebrates show hybrid vigour while sites subject to less long-term selection show a

23 bias in mutation towards alleles that are selected against.

## Introduction

24

25    Classical mutations, with a large effect on phenotype, tend to decrease fitness, decrease

26    fitness-related traits and be partially recessive [1-3] (also see the 1st category of mutations

27    defined in [3]). However, the majority of the genetic variance in complex traits is due to

28    mutations of small effect. Do these small-effect mutations show the same characteristics as

29    those classical large-effect mutations? A study in *E. coli* showd that mutations with small

30    effect on fitness tend to be deleterious to protein function [4]. However, how mutations affect

31    complex traits such as body size, health and fertility is unknown.

32    A better understanding of the consequence of mutations not only updates scientific

33    knowledge but also has practical implications. Domestic cattle support humans with food,

34    labour, clothing material and transportation. Today, there are over 4 billion cattle across the

35    world and over ~900 million tonnes of dairy products have been produced annually for

36    human consumption (http://www.fao.org/3/ca8341en/CA8341EN.pdf). When practicing

37    genomic selection, which is widely used in animal breeding [5], it would be an advantage to

38    know a priori whether mutations are more likely to increase or decrease traits of interest.

39    In particular, if a trait is related to fitness, one might expect mutations to be deleterious [2,6].

40    Therefore the first objective of this study is to determine whether mutations, defined as the

41    non-ancestral allele (also known as derived alleles) at segregating sites, tend to increase or

42    decrease individual complex traits and whether this depends on the trait's association with

43    fitness.

44    Traits that are related to fitness typically show inbreeding depression and heterosis caused by

45    directional dominance. That is, fitness decreases with increased inbreeding due to increased

46    homozygosity at loci with recessive deleterious alleles [7]. Conversely, fitness generally

47    increases with heterozygosity [8]. Therefore, directional dominance can be used to link traits to

48    fitness. Here, we introduce a method testing for directional dominance by estimating the

49    effect of heterozygosity at genomic sites on traits of cattle and use this method to identify

50    traits that are associated with fitness. Then, we classify traits showing directional dominance

51    as 'fitness-related traits'.

52    A likely cause of directional dominance is that mutations tend to be deleterious and partially

53    recessive. However, not all sites in the genome affecting a trait may show this pattern. Our

54    second objective is to test the hypothesis that sites, where the same allele has been conserved

55    across vertebrate evolution, are the most likely to show directional dominance. Therefore, we

56    consider conserved sites and other polymorphic sites in this analysis.

57    Cattle presents a unique opportunity for studying the effects of mutation. The cattle family

58    diverged from other artiodactyls up to 30 million years ago [9]. Modern cattle are derived from

59    at least two different subspecies of wild aurochs, i.e., *Bos primigenius primigenius* (Eurasian

60    aurochs) and *Bos primigenius namadicus* (Indian aurochs) which diverged up to 0.5 million

61    years ago [10-17]. Domestication of *Bos p. primigenius* led to the humpless *Bos taurus*

62    subspecies, which has evolved some highly productive breeds for agriculture, such as the

63    famous black-and-white Holstein breed with superior milk productivity. Besides natural

64    selection, dairy cattle breeds experienced very recent and intensive selection for milk

65    production traits [18,19] and stature [20]. Domestication of *Bos p. namadicus* gave rise to the

66    humped *Bos indicus* subspecies which evolved breeds with strong resistance to hot climates,

67    such as Brahman and Gir cattle.

68    In the present study, we use yak, sheep and camel as outgroup species to assign cattle

69    ancestral alleles for 8M sequence variants (at 8M genomic sites). For each of these variants,

70    the alternative to the ancestral allele is the mutant allele (MA). We estimate the effect of the

71    mutant allele at these 8M variable sites on 37 traits of 113k cattle from 4 breeds. We also

72    estimate the effect of heterozygosity on these traits using both conserved sites and all

73    genomic sites.

74    If mutant alleles decrease fitness we expect selection to reduce their allele

75    frequencycompared with mutant alleles that either have no effect or increase fitness.

76    Therefore, we compare the allele frequency of mutant alleles that increase and decrease each

77    trait. We expand the analysis of mutant allele frequency to additional breeds of ancient and

78    modern cattle from the 1000 Bull Genomes database [21,22], which provides validation of our

79    results. Additional analyses of MAs with strong effects on milk production traits [23,24] suggests

80    that the direction of phenotypic effects of these MAs correlates with their direction of effects

81    on the expression of genes in milk cells[4,25].

82

## Results

84    *Directional dominance at sites conserved across 100 vertebrates*

85    To identify traits related to fitness, we have introduced a method to estimate the effect of

86    heterozygosity on 37 traits (described in Supplementary Table 1) recorded in over 100k

87    animals. In total, there were 16,035,443 imputed sequence variants (at 16,035,443 genomic

88    sites) with imputation accuracy $R^2 > 0.4$ and the minor allele frequency (MAF) $> 0.005$

89    available for variant-trait association analysis. A subset of these sequence variants that could

90    be assigned with ancestral alleles was used for analyses related to mutant alleles (described

91    later). For the analysis of the effect of heterozygosity, we fit the average heterozygosity of

92    sequence variants at 317,279 genomic sites conserved across 100 vertebrates ($H'_{cons_j}$) and

93    heterozygosity from variants at the other 15,718,164 sites ($H'_{non-cons_j}$) simultaneously (see

94    Methods). We observed a significant effect of heterozygosity at consserved sites for the yield

95    of protein (Prot), fat (Fat) and milk (Milk), survival (Surv), fertility performance (Fert),

96    stature (Stat) and angularity (related to slimness and milk yield) (Figure 1 and Supplementary

97    Figure 1). For all these traits, heterozygosity at other sites ($H'_{non-cons_j}$) was not significant

98    when fitted together with $H'_{cons_j}$. This directional dominance implies that milk production,

99 fertility, survival and stature show inbreeding depression and heterosis and therefore we

100 classify them as fitness-related traits and this directional dominance for these traits is

101 exclusively explained by genomic sites conserved across vertebrates. To be conserved across

102 vertebrates, mutations at these sites must be deleterious, implying extremely long-term

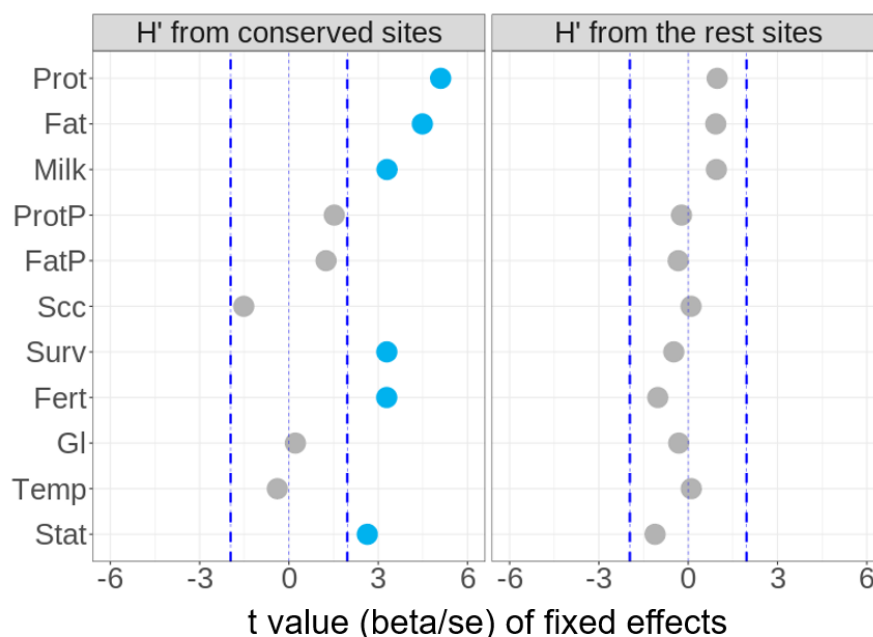103 consistent selection for the ancestral allele at these sites.

104



105

106 **Figure 1**. Directional dominance at conserved sites ($H'$) for traits of 104k cows. The beta

107 values and standard errors for each trait were generated using a mixed linear model, fitting $H'$

108 from 317,279 conserved sites (left panel) and $H'$ from the remaining 15,718,164 sites (right

109 panel) together with other fixed effects (e.g., breed). Blue dashed lines indicate t value of -

110 1.96 and 1.96 commonly used to indicate the significance.

111

112 *Assignment of bovine ancestral and mutant alleles*

113 To assign the mutant alleles in cattle, we first determined the alternative, ancestral alleles

114 using artiodactyls, including cattle as the focal species (98 global cattle breeds from the 1000

115 Bull Genomes Project [21,22], Supplementary Table 2) and yak, sheep and camel as outgroup

116 ancestor species (Ensembl 46-mammal sequence data). A probabilistic method [25] was used to

117  assign an ancestral allele for each site mappable between 4 artiodactyl species (see Methods).

118  Out of 42,573,455 equivalent sites between the 4 species, 39,998,084 sites had the ancestral

119  allele assigned with high confidence (probability > 0.8). We compared our results with a

120  previous study using different methods [26]. Of 1,925,328 sites that were assigned ancestral

121  alleles with high confidence in both studies, 1,904,598 (98.7%) sites agreed. However, we

122  have assigned ancestral alleles with high confidence to ~10 times more sites than the previous

123  study due to the use of large sample size and whole-genome sequence data. The full results

124  are publicly available at https://figshare.com/s/dd5985b76a413b56106b.

125

126  *Biases in trait effects between ancestral and mutant alleles*

127  We conducted GWAS of 37 traits using over 16 million imputed sequence variants in bulls

128  (N ~ 9k) and cows (N ~ 104k) separately (see Methods). For 7,910,190 variants where the

129  ancestral allele was assigned, we compared the direction (increase or decrease) of the effect

130  of the mutant alleles (MAs) on the trait (Supplementary Figure 2-3). The same comparison

131  was also performed for variants at the 202,530 out of 317,279 conserved sites where the

132  ancestral alleles could be assigned. Note that for a variant, the effect of a MA is identical to -

133  $1 \times$ the effect of the ancestral allele. We focus the description of effects on MAs, but a MA

134  increasing the trait is identical to an ancestral allele decreasing the trait.

135  Within all analysed variants and conserved variants, for each trait we considered the

136  following three variant categories for systematic comparison: 1) large-effect variants, i.e., p-

137  value of GWAS ($p_{gwas}$) < 5e-8 and the effect direction agreed in both sexes; 2) medium-

138  effect variants, i.e., 5e-8 <= $p_{gwas}$ < 5e-5 and the effect direction agreed in both sexes, and 3)

139  small-effect variants, i.e., 5e-5 <= $p_{gwas}$ < 0.05 and the effect direction agreed in both sexes.

140  Here the effect size refers to the amount of variance explained by variants which is inversely

141  related to the p-value. The use of different effect size is because mutations of small and large

142    effects may be different in their direction of effect. Selecting variants that have the same

143    effect direction between independent GWAS populations [27], such as bulls and cows, helps to

144    eliminate variants with spurious trait associations from the comparison. Based on a previous

145    method [27], the True Discovery Rate by Effect Direction (TDRed) of GWAS between two

146    sexes across 37 analysed traits for the small-, medium- and large-effect variants was 0.8, 0.98

147    and 0.99, respectively.

148    Based on GWAS results of each trait, we calculated the ratio of the number of variants where

149    the MA increased the trait (positive effect) to the number of variants where the MA decreased

150    the trait (negative effect). Across 37 traits and three effect-size groups, MAs showed diverse

151    trait effect patterns (Supplementary Figure 3). Results observed from GWAS were confirmed

152    by BayesR analysis [28], which jointly fits on average 4.3 million variants per trait (See

153    methods and Supplementary Figure 3). Based on jointly estimated effects for a given set of

154    variants, the significance of the effect direction bias was tested using Kolmogorov-Smirnov

155    to estimate the p-value ($p_{ks}$) of the difference in the effect distribution between ancestral and

156    mutant alleles (see Methods). We also tested the significance of bias using LD-clumped ($r^2 <$

157    0.3) [29] variants to calculate the standard error (Supplementary Figure 4).

158    In addition, we checked the direction of effects of MAs which had large positive effects and

159    large negative effects on protein yield, fat yield, milk yield, protein % and fat % on the

160    expression of genes within ±1Mb distance to these MAs (cis eQTL genes, see Methods) in

161    milk cells[23,24]. For 4 out of 5 sets of variants where the mutant allele decreased the trait, we

162    found the mutant allele tended to decrease the expression of cis eQTL genes. For another 4

163    out of 5 sets of variants where mutant alleles increased the trait, the mutant allele tended to

164    increase the expression of cis eQTL genes (Supplementary Table 3). These results suggest

165    correlated direction of effects of MAs on milk production traits and the expression of genes

166    in milk cells.

167

168    In the following text, we focus on 1) MAs within the large- and small-effect categories for

169    milk production traits as these two sets of MAs showed distinct effect direction patterns

170    (Figure 2), and 2) MAs associated with other traits, including those with medium or small

171    effects on somatic cell count (Scc, indicative of mastitis, medium-effect), survival (Surv,

172    small-effect), fertility (Fert, frequency of pregnancy, small-effect), gestation length (Gl,

173    medium effect), temperament (Temp, docility, small-effect) and stature (Stat, medium effect)
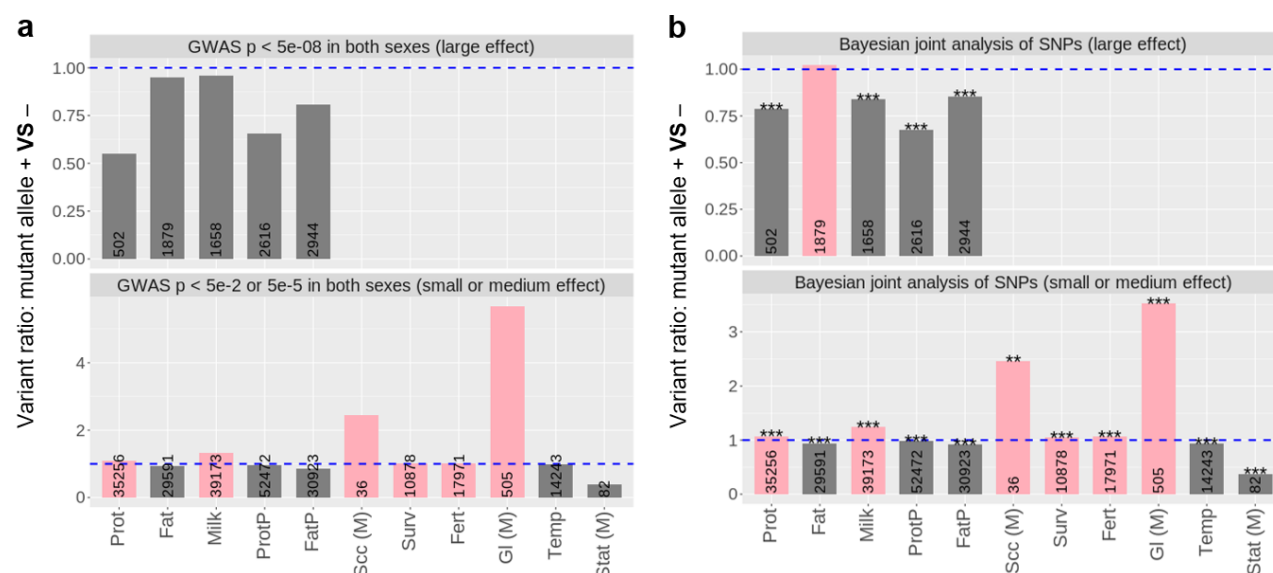
174    (Figure 2).

175



176

177    **Figure 2**. The ratio (y-axis) between the number of variants with mutant alleles increasing

178    the trait (+) and the number of variants with mutant alleles decreasing the trait (−). GWAS

179    effects of mutant alleles are shown for all variants (**a**). BayesR joint effects of mutant alleles

180    from the same variants in (**a**) are shown for all variants (**b**). Pink colour: the majority of

181    variants with mutant alleles tend to increase the trait (taller than the blue-dashed line). Dark

182    grey: the majority of variants with mutant alleles tend to decrease the trait (shorter than the

183    blue-dashed line). Numbers in bars: total number of variants significant at the given

184    threshold. Stars: p-value for the significance of the difference in the distribution of BayesR

185    effects between ancestral and mutant alleles, '*': $p < 0.05$, '**': $p < 0.01$, '***' $p < 0.001$.

186  For somatic cell count (Scc), gestation length (Gl) and stature (Stat), the results are from

187  medium-effect (M) variants and the full results are shown in Supplementary Figure 3.

188

189  The classical model [1-3] predicts that the majority of MAs, or mutations, are deleterious or

190  slightly deleterious. In our study, MAs consistently showed biases towards decreasing protein

191  and fat concentration (Figure 2 and Supplementary Figure 3,4), docility and stature, and

192  towards increasing somatic cell count (an indicator of mastitis) and gestation length. Among

193  these traits only stature showed a significant effect of heterozygosity. For milk yield and

194  protein yield, both of which were classified as fitness-related traits (Figure 1), the bias in the

195  direction of MA depends on the size of the MA effect. Large-effect MAs tended to decrease

196  milk and protein yield whereas small-effect MAs tended to increase them. A possible

197  explanation is that mutation seldom has a large positive effect on milk protein yield or

198  fertility but small positive effect mutations occur and are increased in frequency by natural or

199  artificial selection.

200

201  Also, there was a slight majority of small-effect MAs which tended to increase fertility and

202  survival, both of which were positively related to fitness (Figure 1). The effects of these sets

203  of MAs is partially due to pleiotropy, i.e., the effect of these MAs on multiple traits

204  (Supplementary Table 4). For instance, while small-effect MAs increasing milk yield

205  decreased fat yield, protein % and fat %, they also increased protein yield. Also, while small-

206  effect MAs increasing fertility increased gestation length, they also increased stature.

207

208  The simplest explanation for the bias in the direction of MA effects is that it is due to a bias

209  in the direction of mutation. For instance, that mutation more often leads to a decrease in

210  fat% rather than an increase. However, it is also possible that mutations that decrease fat%

211  are selected and therefore more likely to be discovered than mutations that increase fat%.

212    Below we exclude this possibility by comparing the allele frequency at variants where the

213    MA increases or decreases the trait.

214

215    *Allele frequency of mutant alleles in modern and ancient cattle*

216    Across all variable sites, the allele frequency of MAs was lower than the allele frequency of

217    ancestral alleles (Supplementary Figure 5). Also, the frequency of MAs at conserved sites

218    (0.27) was lower than the frequency of MAs across all sites (0.32). This is consistent with the

219    selection for the ancestral allele which is necessary to maintain conservation of the same

220    allele across vertebrates.

221    We grouped variants based on their mutant allele reducing (MAs−) or increasing the trait

222    (MAs+) and compared their allele frequency in over 110k Holstein, Jersey, crossbreds and

223    Australian Red bulls and cows (Figure 3a,b). To account for LD, we estimated the error of

224    MA frequency based on LD-clumped ($r^2 < 0.3$) [29] variants. As an external validation, we also

225    considered this analysis in a selection of 7 subspecies/breeds of 1,720 ancient and modern

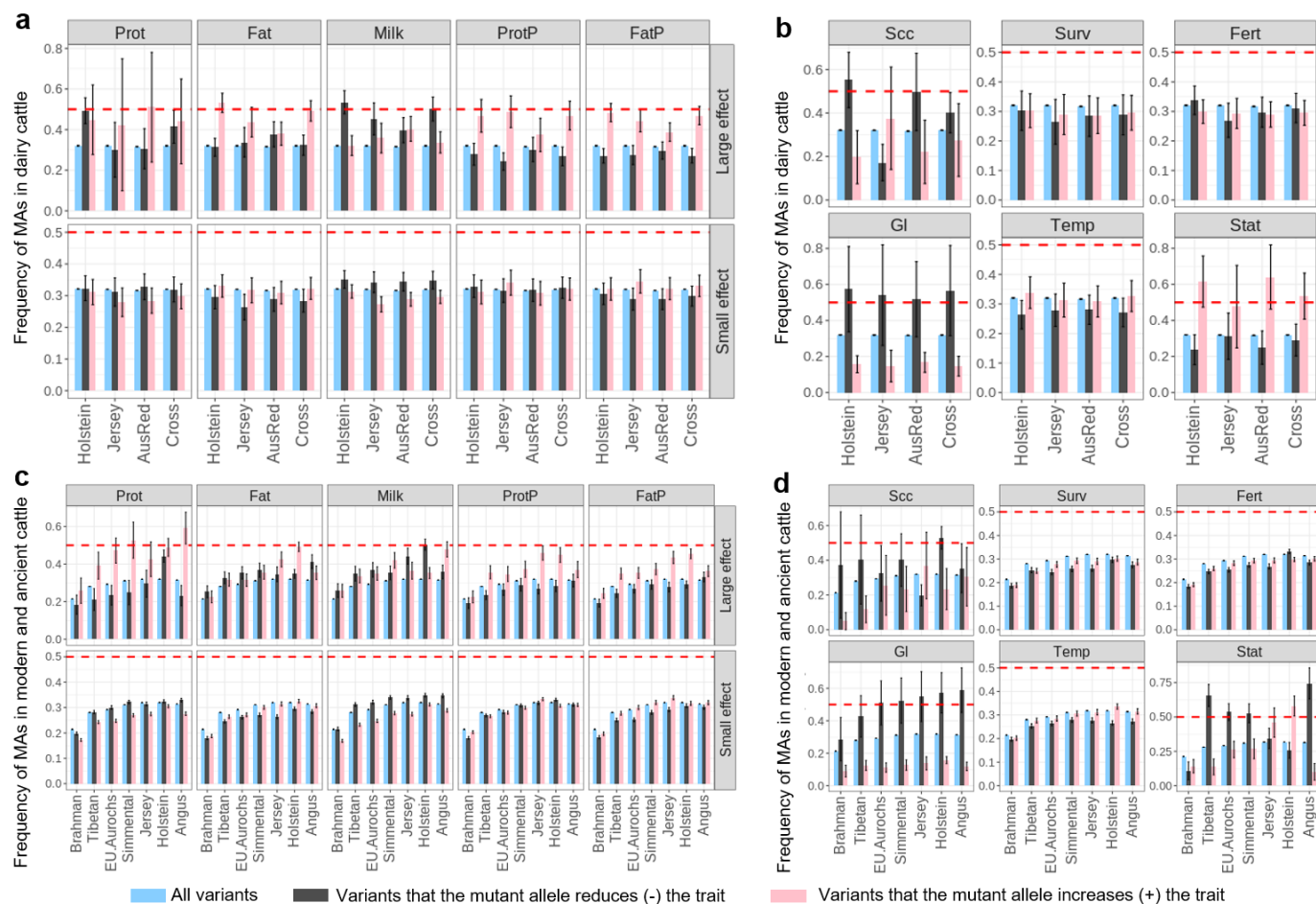226    cattle from the 1000 Bull Genomes Project [21,22] (Figure 3c,d).

227

**Figure 3**. The allele frequency of mutant alleles (MAs) in cattle. The average frequency of variants associated with different traits is shown with standard error bars based on LD clumped variants. All variants include the 7.9M variants where mutant alleles were assigned. Red dashed line represents the frequency of 0.5. In the dairy cattle section (**a** and **b**), 90,627 Holstein, 13,465 Jersey, 3,358 Australian Red (AusRed) and 4,649 crossbreds were used. In the ancient and modern cattle (**c** and **d**), 210 Brahman, 25 Tibetan, 10 Eurasian Aurochs, 242 Simmental, 95 Jersey, 840 Holstein and 287 Angus were used. For panels **b** and **d**, results for survival (Surv), fertility (Fert) and temperament (Temp) were from small-effect MAs while results for somatic cell count (Scc), gestation length (Gl) and stature (Stat) were from medium-effect MAs.

For fat%, protein%, docility and stature MAs that increase the trait had higher allele frequency than MAs that decrease the trait. For somatic cell count and gestation length, the reverse is true. That is, MAs increasing somatic cell count and gestation length had lower allele frequency than MAs that decreased the trait (Figure 3). Thus, although MAs more

244   commonly decreased fat% than increased it, the allele frequency was higher at sites where the

245   MAs increased fat%. This implies that selection acts against MAs that decrease fat% or

246   favours MAs that increase fat%. Consequently, the higher incidence of MAs that decrease

247   fat% cannot be due to selection favouring them but must be due to mutation more often

248   resulting in an allele that decreases fat% than increases it. Comparing results in Figure 2 and

249   3 shows that this is the usual pattern – the more common direction of effects of mutation

250   generates alleles that are selected against and hence have a reduced allele frequency.

251   For other traits, the results are less clear-cut. For milk yield, the majority of MAs of large

252   effect tended to decrease the trait (Figure 2). Interestingly, these large-effect milk-decreasing

253   MAs, which were deleterious, had a higher frequency than those MAs increasing milk yield

254   (Figure 3). On the other hand, the majority of MAs of small effect tended to increase the milk

255   yield (Figure 2). Yet, these small-effect MAs that increase milk were at a lower frequency

256   than MAs that decrease milk yield (Figure 3). Interpretation of these results is helped by

257   remembering that milk yield is negatively correlated with fat% and protein% (Supplementary

258   Table 4).

259

260   *Selection of trait-associated mutant alleles in modern and ancient cattle*

261   The above results for MA frequency at trait-associated variants imply selection. The selection

262   could be consistent across breeds which would limit the divergence of allele frequency

263   between breeds or it could be different between breeds leading to divergence in allele

264   frequency. We compared the average of Wright's fixation index ($\overline{F_{ST}}$), for MA+ variants and

265   MA− variants calculated using dairy cattle (Figure 4a,b) and ancient and modern cattle

266   (Figure 4c,d). To account for LD, we estimated the error of $\overline{F_{ST}}$ based on LD-clumped ($r^2 <$

267   0.3) [29] variants.

268

269 ▮ All variants  ▮ Variants that the mutant allele decreases (-) the trait  ▮ Variants that the mutant allele increases (+) the trait
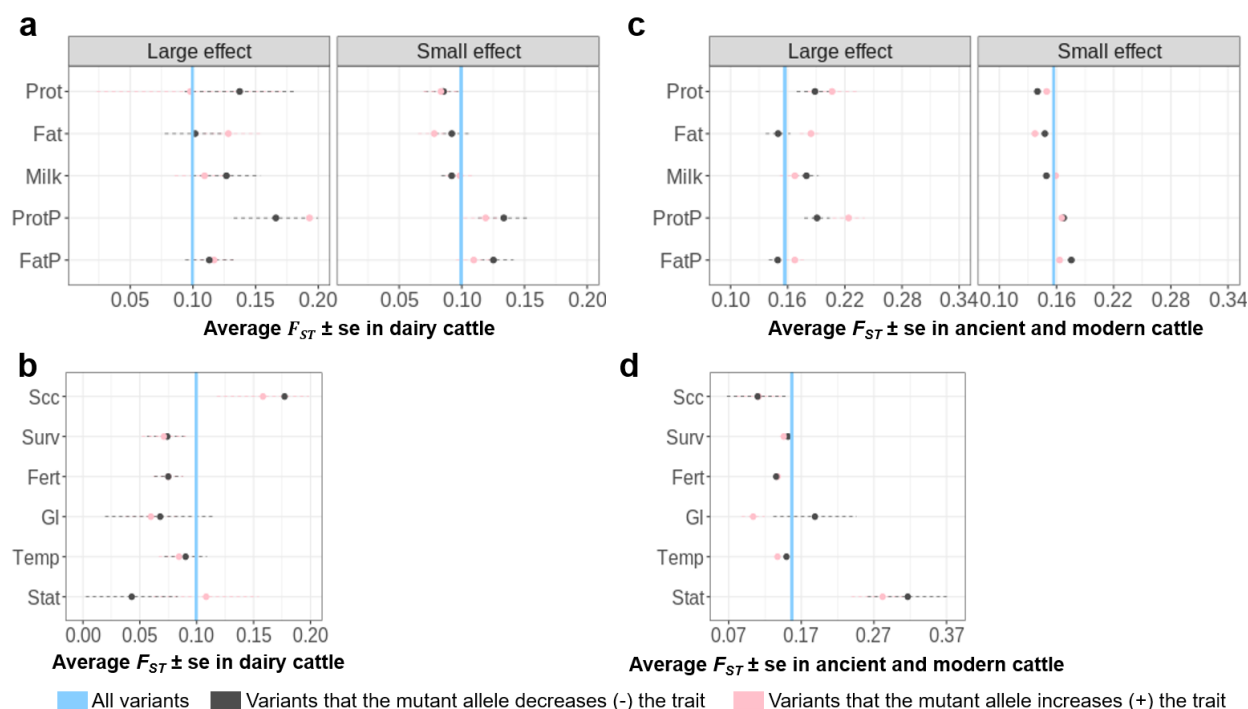
270 **Figure 4**. Selection (average Wright's fixation index $\overline{F_{ST}}$) of variants with mutant alleles that

271 increase or decrease the trait in dairy cattle (**a,b**) and ancient and modern cattle (**c,d**). The $\overline{F_{ST}}$

272 is shown as dots with its standard error bars estimated using LD clumped variants. The blue

273 line represents the $\overline{F_{ST}}$ for 7.9M variants analysed ($0.1\pm4.3$e-05) in dairy cattle in **a** and **b**; and

274 $\overline{F_{ST}} = 0.157\pm5$e-05 in ancient and modern cattle in **c** and **d**). For panels **b** and **d**, results for

275 survival (Surv), fertility (Fert) and temperament (Temp) were from small-effect MAs while

276 results for somatic cell count (Scc), gestation length (Gl) and stature (Stat) were from

277 medium-effect variants.

278

279 In general, variants associated with milk production traits (including somatic cell count,

280 Figure 4a) showed higher than average $F_{ST}$ among dairy breeds implying divergent selection,

281 while variants associated with other traits, including survival and fertility, tended to have

282 below-average $F_{ST}$ indicating convergent selection (Figure 4b). $\overline{F_{ST}}$ for gestation length was

283 below average especially for MA+, probably due to selection against mutations that increase

284 gestation length in all breeds (Figure 4d).

285 Among ancient and modern cattle, $\overline{F_{ST}}$ is high for both MA+ and MA− variants for stature

286 indicating divergent selection for height (Figure 4b). The allele frequency of MAs decreasing

287   height was the least frequent in Holstein cattle and was the most frequent in Tibetan cattle

288   living at high altitude and Angus cattle selected for beef production (Figure 3d). This

289   suggests that the direction of selection could vary across cattle breeds under different

290   environmental conditions and/or artificial selection.

291

## Discussion

293   For some traits (e.g. survival, fertility) we expect that an increase in the trait leads to an

294   increase in fitness. It is these traits which typically show heterosis and inbreeding depression

295   due to directional dominance. The simplest explanation for these observations is that

296   mutations at sites affecting the trait tend to reduce the trait and be partially recessive.

297   However, our results show that it is not all sites affecting these traits that show directional

298   dominance but only those where the same allele is highly conserved across vertebrates. This

299   result explains why the mutations tend to lead to a decrease in the trait - long-term selection

300   has nearly fixed the favourable allele and so any mutation will cause a decrease in the trait

301   and in fitness. We partially confirm this explanation by finding that mutations for these traits

302   (milk and protein yield, stature but not fertility and survival) do tend to decrease the trait

303   although, for milk and protein yield, it is only mutations of large effect for which the effects

304   tend to be negative. This long term selection cannot be directly on traits involving lactation

305   since the same allele is conserved in vertebrates other than mammals.

306   For other traits we expect that an intermediate value leads to the highest fitness. For instance,

307   too high or too low a fat% in milk might be detrimental to the fitness of the mother or the

308   infant or both. These traits do not typically show inbreeding depression or heterosis. The

309   fittest allele might vary between species and environments. Therefore one might expect that

310   mutations are equally likely to increase or decrease the trait. However, that is not what we

311   found: for fat% and protein% mutations tend to decrease the trait whereas for SCC and

312    gestation length they tend to increase the trait. We hypothesise that at some of the genomic

313    sites affecting these traits selection has been consistent enough in mammals, or at least in

314    cattle, so that mutations cause a decrease in fat% and protein% and an increase in mastitis or

315    SCC and gestation length (leading to difficulty calving). This hypothesis is supported by our

316    finding that selection decreases the allele frequency of these mutations. This low allele

317    frequency is not only seen in dairy cattle but in beef breeds and *Bos indicus* breeds.

318    The findings on individual traits can be unexpected due to pleiotropy. That is, mutations

319    affect multiple traits. There are mutations at *DGAT1* and *GHR* loci that increase milk yield

320    but decrease fat% and protein% (Supplementary Figure 6). These are only at appreciable

321    frequency in domesticated cattle, especially breeds artificially selected for milk volume.

322    Their low allele frequency in other breeds and species suggest that natural selection acts

323    against the mutation thus increasing fat% and protein% but decreasing milk yield. Similarly,

324    there is a negative genetic correlation between milk yield and fertility so mutations that

325    increase milk yield might be favoured despite their negative effect on fertility. MAs

326    decreasing fertility tended to be most frequent in the Holstein breed (Figure 3d), perhaps

327    because these alleles tended to increase milk yield and stature.

328    For milk, fat and protein yield the results differ between mutations of large and small effects.

329    Mutations with a large effect on milk protein yield more often decrease protein yield than

330    increase it perhaps because the physiology supporting milk protein synthesis has been

331    optimised in part at least. Mutations with a small effect on protein yield are almost equally

332    likely to increase or decrease yield perhaps because natural selection favours an intermediate

333    level of milk protein yield because too high a yield drains the cow of nutrients needed for

334    maintenance and reproduction.

335    Effects of MAs on phenotypes might be mediated by their effects on gene expression. Based

336    on cis eQTL data [23], we found that MAs with large effects on milk production traits had

337    direction of effects that were correlated with their direction of effects on gene expression in

338    milk cells. This result shows that the effect direction of MAs on gene expression may also

339    show systematic biases and this may be related to their effects on phenotypic traits. Future

340    studies with larger sample size and more tissues for eQTL mapping may update our

341    understanding of the MA effects on molecular phenotypes.

342    The selection which we have observed affecting the frequency of mutations of positive and

343    negative effect could be both natural selection acting over a long period before and since the

344    domestication of cattle, and artificial selection acting over the last 10,000 years and, more

345    intensely, over the last ~100 years in dairy cattle. Artificial selection may differ between

346    breeds and generate high Fst between breeds. For fat%, protein% and stature at least one

347    class of mutation is more common than random mutations and the overall $F_{ST}$ between breeds

348    tended to be high. Our analysis also highlighted some specific breeds. For example, the

349    selection of variants associated with somatic cell count led to high $F_{ST}$ among dairy cattle but

350    low $F_{ST}$ in our other breeds. Holstein cattle have been selected to be tall [20] and this is

351    reflected in the low frequency of MAs decreasing stature in Holstein. On the other hand, the

352    high frequency of MAs decreasing stature in Tibetan cattle (Figure 4d) may be due to its

353    adaptation to high altitude [30].

354    Although mutation is biased in its effect on some traits, the bias is small for most traits. That

355    is, mutations decreasing protein yield are only slightly more common than mutations that

356    increase protein yield. Also, although conserved sites explain directional dominance and are

357    enriched for polymorphisms affecting complex traits [31], they do not explain the majority of

358    the genetic variance. That is, there are many sites affecting traits, such as milk yield and

359    stature, at which the allele carried varies between species implying that the fittest allele varies

360    depending on the environment and the background genotype of the species.

361  The sequence variants associated with a complex trait are not necessarily causal but likely to

362  be in high LD with the causal variants. This tends to dilute the signal that might be

363  discovered if causal variants were used. However, variants in high LD may share a similar

364  evolutionary history and therefore show some of the same characteristics. We used BayesR

365  which jointly fits variants and LD-clumping to account for LD. However, we acknowledge

366  that we cannot completely remove the effects of LD on our results. Therefore, future studies

367  with even larger sample sizes, e.g., ~1 million, may update our results.

368  Genomic selection [32], used in the breeding of livestock and crops, estimates the genetic value

369  of individuals for traits of interest from the alleles they carry at genetic markers such as

370  SNPs. The equation predicting genetic value uses the effect of each SNP on the trait

371  estimated in a training population. The best methods treat the SNP effects as random

372  variables drawn from a prior distribution. To date it has been assumed that the effects of a

373  mutation are equally likely to be positive or negative on the trait but, if it was known that one

374  direction of effect was more likely, this could be built into the prior distribution resulting in

375  an increase in the accuracy with which genetic value is predicted.

376  In conclusion, our results support a new hypothesis which provides a new picture of the

377  effects of mutation and selection on mammalian complex traits. Directional dominance,

378  which causes heterosis and inbreeding depression, is characteristic of loci where mutations

379  decrease the trait and fitness and this pattern has been consistent over the evolution of

380  vertebrates. More recent selection, although not causing directional dominance, leads to a

381  bias in the direction of mutation because the mutation results in an allele which is less fit than

382  the ancestral allele and tends to affect a complex trait in a consistent direction. This

383  hypothesis, if supported by future research, adds to our understanding of the evolution of

384  complex traits and has practical value in the artificial selection of livestock and crops.

385

## Methods

**Data preparation for calling bovine ancestral alleles.** The assignment of bovine ancestral alleles was based on a model comparison of alleles from cattle with alleles from outgroups of yak (*Bos grunniens*), sheep (*Ovis aries*) and camel (*Camelus dromedarius*). According to the evolutionary relationships reported previously [9], among ruminants, yak is an outgroup species closely related to cattle, while sheep is less closely related to cattle than yak. Goat is equivalent to sheep in its relationship to cattle, but we chose sheep in the current study. Camel without the rumen is distantly related to cattle as they are artiodactyls. For the cattle species, we used whole-genome sequence data of 98 individuals from Run 7 of the 1000 Bull Genomes Project [21,22]. Each of those above mentioned 98 individuals represents a breed collected by the consortium. Only those whole-genome sequence samples with coverage > 10x were selected and if multiple individuals were found for a breed, the whole-genome sequence sample with the highest coverage was chosen. Both *Bos taurus* and *Bos indicus* subspecies were included (Supplementary Table 2). The pre-processing of sequence reads and alignment of sequence data is done by project partners using the standard 1000 Bull Genomes Project pipeline: http://www.1000bullgenomes.com/. Only BAM files from 1000 Bull Genomes partners are collected and processed by the consortium. The latest published data from the 1000 Bull Genomes Project (1832 samples) can be found at https://www.ebi.ac.uk/eva/?eva-study=PRJEB42783. The details of variant calling can be found in [33]. Briefly, Genome Analysis Toolkit (GATK v.3.8) [34] was used for variant calling. Variants from the GATK VQSR (Variant Quality Score Recalibration) 99.90 to 100.00 Tranche for SNP and INDEL were excluded, and Beagle v.4.0 [35] was used to impute sporadic missing. Whole-genome sequence data in VCF format for these 98 cattle, as a subset from the 1000 Bull Genomes Project database, was generated for further analysis.

410     For the outgroup species (to determine ancestral alleles), we used whole-genome sequence

411     data of 46 mammals stored in the Multiple Alignment File generated by Ensembl EPO

412     pipeline (http://asia.ensembl.org/info/genome/compara/multiple_genome_alignments.html). The

413     46-mammal EPO Multiple Alignment File was downloaded. Then, the software WGAbed

414     (https://henryjuho.github.io/WGAbed/) from python v2.7 was used to retrieve sequence data

415     for cattle, yak, sheep and camel in bed file format. Only sites with sequence data available in

416     at least one outgroup species were kept. Using the cattle coordinates in the 4-species

417     WGAbed files, the sequence data of the outgroup species were matched with the 98 cattle. As

418     a result, 42,573,455 sites found in the 98 cattle and in at least one outgroup species were

419     found. Sequence data on these 42,573,455 sites across 4 species were used to determine the

420     bovine ancestral alleles.

421     **Probabilistic determination of bovine ancestral alleles.** We used the method proposed by

422     Keightley et al [25] with the model choice of The Kimura two-parameter (K2) which accounts

423     for allele frequency of the focal species to determine the probability of an allele being

424     ancestral at each available site. The method was implemented in estsfs [25] and the K2 model

425     was chosen due to its equivalent accuracy to other models but better computation efficiency.

426     As described above, the sequence data of three outgroup species were used. The order of

427     phylogenetic tree topology was cattle → yak → sheep → camel. As requested by the

428     software, allele counts of A, C, G and T were determined for the focal species (cattle) and for

429     out species at each available site. For cattle, the total allele count for each site was 196 (98

430     ×2). For each outgroup species, the total allele count for each site was up to 1. Missing

431     sequence data in the outgroup species were treated as 0 counts. For each site, estsfs produced

432     a probability ($P_{ancs}$) of the major allele in the focal species being ancestral. We then

433     determined alleles which were major at a site with $P_{ancs} > 0.8$ or those alleles which were

434     minor at a site with $P_{ancs} < 0.2$ to be ancestral. For those sites where the major or minor

435    alleles could not be determined but the $P_{ancs} > 0.8$ or $< 0.2$, the cattle allele with the highest

436    frequency in the 3 out species was assigned ancestral. The rest of the sites were determined as

437    ambiguous where no clear ancestral alleles could be determined. The detailed results of

438    ancestral alleles for those 42,573,455 sites across 4 species and the probability of the alleles

439    being ancestral or ambiguous is publicly available at:

440    https://figshare.com/s/dd5985b76a413b56106b.

441    **Sequence variants under conserved sites across 100 vertebrate species**. The variant

442    selection followed previous procedures [31]. Briefly, conservation was determined by the

443    criteria of PhastCon score [36] $> 0.9$ based on the sequence data of those 100 species. The

444    choice of 0.9 as the cutoff was arbitrary. However, since PhastCons score ranges from 0 to 1,

445    this cutoff kept relatively highly conserved sites. Also, in a previous study [31], cattle variants

446    from sites with PhastCon score $> 0.9$ were highly enriched for the heritability of cattle traits.

447    The conserved sites were primarily determined using the human genome coordinates (hg38)

448    and were lifted over to the bovine genome ARS-UCD1.2 using the LiftOver software

449    (https://genome.ucsc.edu/cgi-bin/hgLiftOver) with a lift-over rate $> 92\%$. In total, 317,279

450    variants in the current study were assigned as the conserved variants.

451    **Animals and phenotypes for variant-trait association analysis**. Data was collected by

452    farmers and processed by DataGene Australia (http://www.datagene.com.au/) for the official

453    May 2020 release of National breeding values. No live animal experimentation was required.

454    Phenotype data was based on trait deviations for cows and daughter trait deviations for bulls.

455    Daughter trait deviations were the average trait deviations of a bull's daughters and all

456    phenotypes were pre-corrected for known fixed effects, with processing done by DataGene.

457    Phenotype data used included a total of 8,949 bulls and 103,350 cows from DataGene.

458    Holstein (6,886♂ / 87,003♀), Jersey (1562♂ / 13,353♀), cross-breed (36♂ / 5,037♀) and

459    Australian Red dairy (265♂ / 3,379♀) breeds were included. In total, 37 traits related to milk

460     production, mastitis, fertility, temperament and body conformation were studied

461     (Supplementary Table 1). Larger trait values of fertility (Fert), ease of birth (Ease),

462     temperament (Temp), milking speed (MSpeed), likeability (Like) meant poor performances,

463     so to assist the interpretability of the study, we have corrected the trait direction so that larger

464     values of Fert, Ease, Temp, MSpeed and Like meant increased fertility performance (calving

465     frequency), labour ease, docility, milking speed and the overall preference as a dairy cow

466     (Supplementary Table 1). This correction only affected the reported effect direction of the

467     mutant allele.

468     **Genotype data for association analysis**. The genotypes used in the current study included a

469     total of 16,035,443 imputed bi-allelic sequence variants with Minimac3 [37,38] imputation

470     accuracy $R^2 > 0.4$ and the minor allele frequency (MAF) > 0.005 in both sexes. Most bulls

471     were genotyped with a medium-density SNP array (50K: BovineSNP50 Beadchip, Illumina

472     Inc) or a high-density SNP array (HD: BovineHD BeadChip, Illumina Inc) and most cows

473     were genotyped with a low-density panel of approximately 6.9k SNPs overlapping with the

474     standard-50K panel. The low-density genotypes were first imputed to the Standard-50K panel

475     and then all 50K genotypes were imputed to the HD panel using Fimpute v3 [31,39]. Prior to

476     sequence imputation, the HD genotypes were converted to forward sequence format. Then,

477     all HD genotypes were imputed to sequence using Minimac3 with Eagle (v2) to pre-phase

478     genotypes ([38,40]). The reference set for imputation included sequences of 3090 *Bos taurus*

479     animals from Run7 of the 1000 Bull Genomes Project [21] aligned to the ARS-UCD1.2

480     reference bovine genome (https://www.ncbi.nlm.nih.gov/assembly/GCF_002263795.1/) [22,41].

481     The accuracy of the sequence data for individual animals in the 1000 Bull Genomes Project is

482     routinely checked against their own high-density SNP array genotypes and the concordance

483     has been above 95% [33]. The empirical accuracy of imputation to sequence using the 1000

484     Bull Genomes project has been routinely tested for dairy breeds: for example, in Holsteins

485  the average correlation between imputed and real sequence variants was 0.92 to 0.95 using

486  Run5 of the 1000 Bull Genomes project (N= 1577)[42]. Therefore, we believe our imputed data

487  is more accurate: first because the number of reference animals has almost doubled and

488  second because in our study we impose a Minimac3 $R^2$ filter to remove poorly imputed

489  variants. A Minimac3 $R^2$ threshold of 0.4 was used because our in-house tests demonstrate

490  that this is approximately equivalent to an empirical imputation accuracy (correlation) of

491  0.85.

492  **Genome-wide association studies**. The above mentioned traits were analysed one trait at a

493  time independently in each sex with linear mixed models using GCTA [43]:

494  $$\mathbf{y} = \mathbf{mean} + \mathbf{breed} + \mathbf{bx} + \mathbf{a} + \mathbf{error} \quad (equation\ 1)$$

495  where $\mathbf{y}$ = vector of phenotypes for bulls or cows, **breed** = three breeds for bulls, Holstein,

496  Jersey and Australian Red and four breeds for cows (Holstein, Jersey, Australian Red and

497  MIX); $\mathbf{bx}$ = regression coefficient $b$ on variant genotypes $\mathbf{x}$; $\mathbf{a}$ =  random polygenic effects

498  $\sim N(0, \mathbf{G}\sigma_g^2)$ where $\mathbf{G}$ = genomic relatedness matrix based on all variants and $\sigma_g^2$ = random

499  polygenic variance; **error** = the vector of random residual effects $\sim N(0, \mathbf{I}\sigma_e^2)$ , where I =

500  the identity matrix and $\sigma_e^2$ the residual variance.The purpose of fitting breeds as fixed

501  effects together with the GRM in the model was to have strong control of the population

502  structure which may cause spurious associations between variants and phenotype. The

503  construction of GRM followed the default setting (--make-grm) in GCTA[43]:

504  (https://cnsgenomics.com/software/gcta/#MakingaGRM).

505

506  **Bayesian mixture model analysis**. In the above-described GWAS, sequence variants, many

507  of which are in high LD, were analysed one at a time. In order to assess variant effects and

508  account for LD, we fitted selected variants jointly in BayesR [28]. For each trait, variants that

509    showed the same sign between bulls and cows (regardless of p-value) and could be assigned

510    with an ancestral allele were analysed with BayesR. Across 37 traits, the number of variants

511    analysed ranged from 3,961,180 to 4,737,492. To reduce the computational burden of

512    BayesR, we estimated the joint effects of these variants for each trait in bulls. BayesR models

513    the variant effects as mixture distribution of four normal distributions including a null

514    distribution, $N(0, 0.0\sigma^2{}_g)$, and three others: $N(0, 0.0001\sigma^2{}_g)$, $N(0, 0.001\sigma^2{}_g)$,

515    $N(0, 0.01\sigma^2{}_g)$, where $\sigma^2{}_g$ was the additive genetic variance for the trait. The starting value

516    of $\sigma^2{}_g$ for each trait was estimated using GREML implemented in MTG2 [44] with a single

517    genomic relationship matrix made of all 16M sequence variants. The statistical model used in

518    the single-trait BayesR was:

$$\mathbf{y} = \mathbf{Wv} + \mathbf{Xb} + \mathbf{e} \text{ (equation 2)}$$

520    where $\mathbf{y}$ was a vector of phenotypic records; $\mathbf{W}$ was the design matrix of marker genotypes;

521    centred and standardised to have a unit variance; $\mathbf{v}$ was the vector of variant effects,

522    distributed as a mixture of the four distributions as described above; $\mathbf{X}$ was the design matrix

523    allocating phenotypes to fixed effects; $\mathbf{b}$ was the vector of fixed effects of breeds; $\mathbf{e}$ = vector

524    of residual errors. As a result, the effect $v$ for each variant jointly estimated with other

525    variants were obtained for further analysis.

526    **The difference in effect distribution between ancestral and mutant alleles**. For an

527    analysed variant, one allele is ancestral and then the other is mutant. If there is a bias in effect

528    direction in ancestral alleles or mutant alleles in a given set of variants, the effect distribution

529    of the ancestral and mutant alleles would be different. We tested if the distribution of the

530    effect of ancestral alleles estimated from BayesR was significantly different from that of

531    mutant alleles using the two-sample Kolmogorov-Smirnov test implemented by ks.test() in R

532    v3.6.1. The coding was ks.test(a,m) where a was the vector of variant effects based on the

533    ancestral alleles and m was a vector of variant effects based on the mutant alleles. To be more

534     conservative, we also tested the significance of biases using LD-clumped ($r^2 < 0.3$ within

535     1Mb windows) variants with small, medium and large effects using default settings in

536     plink1.9 [29].

537     **Heterozygosity of individuals at conserved sites**. It is widely accepted that higher genomic

538     heterozygosity is linked to gene diversity, therefore, fitness. However, it is not clear at which

539     set of genes or variants heterozygosity is more related to fitness. Also, the simple estimation

540     of heterozygosity, i.e., assigning allele counts of 0 or 2 as homozygous and 1 as

541     heterozygous, leads to biases as the estimation is not independent of additive effects

542     (illustrated later). Our previous work showed conserved sites across 100 vertebrate species

543     significantly contribute to trait variation [31,45] and it is also logical to assume mutations at

544     conserved sites tend to have strong effects on fitness. Therefore, we firstly partitioned the

545     genome into 317,279 conserved and 15,718,164 non-conserved variants. Then, we re-

546     parameterised the genotype allele count for each variant commonly used to model the

547     dominance deviation, so that the estimation of dominance deviation is independent of the

548     additive effects. We focused on cows because their traits were largely measured on

549     themselves, contrasting to bull traits which were based on their daughters' traits. We

550     estimated the variant-wise sum of the re-parameterised allele count value for dominance

551     deviation which was later termed as $z'_{D_i}$ for each variant $i$ in cows. The sum was averaged by

552     the number of variants and this average value based on re-parameterised dominance allele

553     count for the individual $j$ was termed as $H'_j$ to represent the individual heterozygosity. We

554     estimated the individual heterozygosity from conserved sites ($H'_{cons_j}$) and non-conserved

555     sites ($H'_{non-cons_j}$) and these computations are specified in the following text.

556     According to quantitative genetics theory [46-48], the genetic value ($G'$) of an individual can be

557     partitioned into the mean (μ), additive genetic value (A) arising from additive effect ($a$) and

558     dominance genetic value (D) arising from dominance deviation ($d$). At a single locus, let the

559     allele frequency of the three genotype classes of AA, AB and BB be $p^2$, $2pq$ and $q^2$,

560     respectively. In a simple genetic model, the genetic value can be decomposed as:

561     $$G' = \mu + A + D + e = \mu + x_{A_i}a + z_{D_i}d + e \text{ (equation 3)}$$

562     Where $x_{A_i}$ was the allele count for genotype AA, AB and BB for locus or variant $i$ which

563     were usually coded as 0, 1, 2, respectively, to represent the additive component, and $z_{D_i}$ was

564     usually coded as 0, 1, 0, for genotype AA, AB and BB for variant $i$, respectively, which

565     differentiates the homozygous and heterozygous to represent the dominance component.

566     Therefore, in the simplest form, the genome-wide heterozygosity of the individual $j$ can be

567     calculated as:

568     $$H_j = \left.\sum_i^N z_{D_i}\middle/ N\right. \text{ (equation 4)}$$

569     where $H_j$ is the simple genome-wide heterozygosity of individual j, $N$ is the total number of

570     variants. Note that such calculation of $H_j$ can also be used to derive inbreeding coefficient,

571     where $I_j = (\sum_i^N 2p_i q_i) \times H_j$. $I_j$ was the inbreeding coefficient for the j$_{th}$ individual.

572     In equation 3, however, due to the non-zero correlation between $x_A$ and $z_D$ under Hardy-

573     Weinberg equilibrium (HWE), the estimation of $a$ and $d$ is not independent, i.e.,

574     $cov(x_A, z_D) = 2p(1 - p)(1 - 2p) \neq 0$ under HWE. This then resulted in the estimation of

575     $H_j$ not being independent of the additive components. Therefore, we proposed to re-

576     parameterise this model to estimate $a$ and $d$ independently.

577     According to Falconer [47] at this locus, the additive effects can be derived using the regression

578     of genetic value on the number of A alleles, where $A'_{AA} = 2q \times \alpha$, $A'_{AB} = (p - q) \times \alpha$ and

579     $A'_{BB} = -2p \times \alpha$. $A'$ is the re-parameterised additive genetic value and $\alpha$ is the allele

580     substitution effect: $\alpha = a + (p - q)d$. Because the dominance deviation is the difference

581     between the genetic value and the mean plus the additive value, the dominance effects can be

582 derived as $D'_{AA} = -2p^2 \times d$, $D'_{AB} = 2pq \times d$ and $D'_{BB} = -2q^2 \times d$. $D'$ is the re-

583 parameterised dominance genetic value. Therefore, equation 3 can be re-parameterised as:

584 $$G' = \mu + A' + D' + e = \mu + x'_{A_i}\alpha + z'_{D_i}d + e \text{ (equation 5)}$$

585 Where $x'_A$ was coded as $2q$, $p - q$ and $-2p$ for genotype of AA, AB and BB of variant $i$,

586 respectively, to represent the additive component and $z'_D$ was coded as $-2p^2$, $2pq$, $-2q^2$ for

587 genotype of AA, AB and BB of variant $i$, respectively, to represent the dominance

588 component. Such re-parametrisation has the following features: 1) The covariance between

589 the additive and dominance effects is zero; 2) the variance of the additive effects gives the

590 additive variance; and 3) The variance of the dominance deviations gives the dominance

591 variance. Equation 5 then leads to:

592 $$H'_j = {\sum_i^N z'_{D_i}}\Big/{N} \text{ (equation 6)}$$

593 Where $H'_j$ was the re-parameterised genome-wide heterozygosity for individual $j$, $z'_D$ was

594 $-2p^2$, $2pq$, $-2q^2$ for genotype of AA, AB and BB of variant $i$ and N was the total number of

595 variants. We then applied equation 6 to conserved and non-conserved variants to estimate

596 individual heterozygosity from conserved sites ($H'_{cons_j}$) and non-conserved sites

597 ($H'_{non-cons_j}$). We then fitted $H'_{cons_j}$ and $H'_{non-cons_j}$ as fixed effects together with the fixed

598 effects of breed jointly in GREML similar to equation 1. The difference was that there is no

599 fixed effect of variants but more fixed effects due to the fitting of $H'_{cons_j}$ and $H'_{non-cons_j}$.

600 The GREML analysis used the implementation with MTG2 [44].

601 **Mutant allele frequency and $F_{ST}$ in different breeds/subspecies**. Two sets of data were

602 used for this analysis. The first dataset was the Australian dairy cattle (8,949 bulls and

603 103,350 cows, Holstein, Jersey, Australian Red and crossbreds) used for GWAS as described

604 above. The second data set used for the analysis of mutant allele frequency and $F_{ST}$ was the

605 curated whole-genome sequence data of 1,720 cattle from the 1000 Bull Genomes database

606  (Run 7) [21,22], which we refer to as modern and ancient cattle. Samples that met the quality

607  criteria of the 1000 Bull Genomes project were selected and they included 210 Brahman, 25

608  Tibetan, 10 Eurasian Aurochs, 242 Simmental, 95 Jersey, 843 Holstein and 295 Angus.

609  Genome sequences from 6 Gir and 12 Nellore cattle from the 1000 Bull Genomes database

610  were also analysed to support the results of mutant allele frequency of *Bos indicus*.

611  Additional information on these 1720 animals including related accession numbers (if

612  available) can be found in Supplementary Data 1. The ancient genome data were part of the

613  project of Verdugo et al 2019 [15] who processed and published the original data (PRJEB31621

614  at European Nucleotide Archive). These data were collected by Run 7 of the 1000 Bull

615  Genomes Project and processed by its standard pipeline

616  (http://www.1000bullgenomes.com/).

617  Sequence data at 7,910,190 variants assigned with mutant alleles were retrieved for these

618  animals to make a plink (v1.9) binary genotype file. The A1 allele of the plink genotypes was

619  set to the mutant allele and its frequency was calculated using the '--freq' function for

620  different selections of populations and variant sets. Average mutant allele frequency and the

621  standard error were calculated for different selections of variants, e.g., variants with mutant

622  alleles increasing or decreasing traits. Standard errors for frequency and $F_{ST}$ (described

623  below) were all estimated using LD-clumped variants in the same procedure in plink [29] as

624  described above. For variants associated with milk production traits, i.e., the yield of milk

625  protein, fat and milk and percentage of protein and fat, we selected variants with large

626  (GWAS p-value < 5e-8 in both sexes) and small (GWAS p-value < 5e-2 and p-value > 5e-5

627  in both sexes) effects to focus on. For other trait-associated variants, the group with the

628  largest effects available were selected for this comparison. For example, for stature, there

629  were no variants with p-value < 5e-8 in both sexes, we then selected the medium-effect

630  variants (GWAS p-value < 5e-5 and p-value > 5e-8 in both sexes). For fertility, there was no

631     variants with p-value < 5e-5 in both sexes, we then selected the small-effect variants (GWAS

632     p-value < 5e-2 and p-value > 5e-5 in both sexes) for the comparison. Average mutant allele

633     frequency and the standard error were also calculated for all 7.9M variants analysed as the

634     baseline. The analysis procedure for allele frequency on the Australian dairy cattle was

635     applied to these 1000 Bull Genomes individuals.

636     With the same plink binary genotype file described above and the population structure for

637     dairy cattle (4 dairy breeds) and for ancient and modern cattle (7 breeds/subspecies), GCTA

638     [43] was used to calculate the $F_{ST}$ value with the method described in Weir [49] with the option of

639     '--fst' and '--sub-pop'. The average $F_{ST}$ value with standard errors was then calculated for

640     different selections of variants in the same fashion for selecting variant groups to compare the

641     mutant allele frequency as described above.

642     **cis eQTL in milk cells**. This analysis was based on 105 Holstein cattle who had RNA-seq

643     data in milk cells described and published previously (NCBI SRA SRP111067) [23,24]. The raw

644     reads of these data were aligned to the ARS-UCD1.2 reference bovine genome using STAR[50]

645     and the quality check followed what was described previously [23]. FeatureCount [51] was used to

646     extract gene counts and the voom [52] normalised counts were used in the following analyses.

647     The normalised gene expression was analysed as phenotypes in the same GWAS model as

648     equation 1 using GCTA, except that there were no breed effects (all animals are Holstein) but

649     were other fixed effects of Experiment, Days in Milk, 1st PC and 2nd PC extracted from the

650     expression count matrix. Variants analysed were those that had large positive effects and

651     large negative effects ($p_{gwas} < 5e-8$) on protein yield, fat yield, milk yield, protein % and fat

652     %. For these variants, the normalised expression of genes within ±1Mb distance to them were

653     analysed as phenotype. In other words, the analysis focused on cis eQTL genes for these

654     large-effect variants were analysed. When GWAS results of gene expression were obtained

655     (cis eQTL), the effect allele was mapped to the ancestral allele to determine the effects of

656  MAs. For quantifying the number of eQTL for each effect direction of MAs, only the SNPs

657  with the smallest p-value were considered.

658

## 659  **Data availability**

660  Our predictions of cattle ancestral alleles for those 42,573,455 sites have been made publicly

661  available at: https://figshare.com/s/dd5985b76a413b56106b. Multiple alignment data used to

662  determine cattle ancestral alleles are publicly available via Ensembl EPO pipeline

663  (http://asia.ensembl.org/info/genome/compara/multiple_genome_alignments.html).

664  Australian farmers and DataGene Australia (http://www.datagene.com.au/) are owners and

665  custodians of the raw phenotype and genotype data of Australian dairy animals. Access to

666  these data for research requires permission from DataGene under a Data Use Agreement. The

667  DNA sequence data as part of the 1000 Bull Genomes Consortium [20-22] are available to

668  consortium members and the membership is open. Sequence data of 1832 samples from the

669  1000 Bull Genome Project have been made publicly available at:

670  https://www.ebi.ac.uk/eva/?eva-study=PRJEB42783. The gene expression data is publically

671  available (NCBI SRA SRP111067). In addition: 1. The summary data of the effect direction

672  and effect category of those 7.9M sequence variants for which the ancestral alleles can be

673  assigned is published at https://figshare.com/s/ef020d948523c31c0e67; 2. The allele frequency

674  of mutant alleles of those 7.9M sequence variants for which the ancestral alleles can be

675  assigned for the Holstein and Jersey cattle from the 1000 Bull Genome Project is published at

676  https://figshare.com/s/20154b1d8e60e012e532; 3. The coordinates of conserved sites analysed in

677  the manuscript is published at: https://figshare.com/s/df9d3662f8f7fb8e72da. Other supporting

678  data are shown in the supplementary materials of the current manuscript.

679

## 680  **Code availability**

681    The probability of ancestral allele assignment used the software published by [25]. The linear

682    mixed model used GCTA [43] and MTG2 [44]. The Bayesian analysis used BayesR [53]. The R

683    code of estimating heterozygosity across conserved sites will be made public upon

684    publication.

685

686    **References:**

687    1      Ohta, T. Slightly deleterious mutant substitutions in evolution. *Nature* **246**, 96-98 (1973).

688    2      Eyre-Walker, A., Keightley, P. D., Smith, N. G. & Gaffney, D. Quantifying the slightly
689           deleterious mutation model of molecular evolution. *Molecular Biology and Evolution* **19**,
690           2142-2149 (2002).

691    3      Eyre-Walker, A. & Keightley, P. D. The distribution of fitness effects of new mutations.
692           *Nature Reviews Genetics* **8**, 610-618 (2007).

693    4      Mehlhoff, J. D. *et al.* Collateral fitness effects of mutations. *Proceedings of the National*
694           *Academy of Sciences* **117**, 11597-11607 (2020).

695    5      Meuwissen, T., Hayes, B. & Goddard, M. Genomic selection: A paradigm shift in animal
696           breeding. *Animal frontiers* **6**, 6-14 (2016).

697    6      Chen, J., Glémin, S. & Lascoux, M. From drift to draft: how much do beneficial mutations
698           actually contribute to predictions of Ohta's slightly deleterious model of molecular
699           evolution? *Genetics* **214**, 1005-1018 (2020).

700    7      Keller, L. F. & Waller, D. M. Inbreeding effects in wild populations. *Trends in Ecology &*
701           *Evolution* **17**, 230-241, doi:https://doi.org/10.1016/S0169-5347(02)02489-8 (2002).

702    8      Turelli, M. & Ginzburg, L. R. Should individual fitness increase with heterozygosity? *Genetics*
703           **104**, 191-209 (1983).

704    9      Jiang, Y. *et al.* The sheep genome illuminates biology of the rumen and lipid metabolism.
705           *Science* **344**, 1168-1173, doi:10.1126/science.1252806 (2014).

706    10     Loftus, R. T., MacHugh, D. E., Bradley, D. G., Sharp, P. M. & Cunningham, P. Evidence for two
707           independent domestications of cattle. *Proceedings of the National Academy of Sciences* **91**,
708           2757-2761 (1994).

709    11     Bradley, D. G., MacHugh, D. E., Cunningham, P. & Loftus, R. T. Mitochondrial diversity and
710           the origins of African and European cattle. *Proceedings of the National Academy of Sciences*
711           **93**, 5131-5135 (1996).

712    12     Troy, C. S. *et al.* Genetic evidence for Near-Eastern origins of European cattle. *Nature* **410**,
713           1088-1091 (2001).

714    13     Chen, S. *et al.* Zebu cattle are an exclusive legacy of the South Asia Neolithic. *Molecular*
715           *biology and evolution* **27**, 1-6 (2010).

716    14     Utsunomiya, Y. *et al.* Genomic clues of the evolutionary history of Bos indicus cattle. *Animal*
717           *genetics* **50**, 557-568 (2019).

718    15     Verdugo, M. P. *et al.* Ancient cattle genomics, origins, and rapid turnover in the Fertile
719           Crescent. *Science* **365**, 173-176 (2019).

720    16     Upadhyay, M. *et al.* Genetic origin, admixture and population history of aurochs (Bos
721           primigenius) and primitive European cattle. *Heredity* **118**, 169-176 (2017).

722    17     Park, S. D. *et al.* Genome sequencing of the extinct Eurasian wild aurochs, Bos primigenius,
723           illuminates the phylogeography and evolution of cattle. *Genome biology* **16**, 234 (2015).

724   18   García-Ruiz, A. *et al.* Changes in genetic selection differentials and generation intervals in US
725        Holstein dairy cattle as a result of genomic selection. *Proceedings of the National Academy*
726        *of Sciences* **113**, E3995-E4004 (2016).
727   19   Coffey, E., Horan, B., Evans, R. & Berry, D. Milk production and fertility performance of
728        Holstein, Friesian, and Jersey purebred cows and their respective crosses in seasonal-calving
729        commercial farms. *Journal of Dairy Science* **99**, 5681-5689 (2016).
730   20   Bouwman, A. C. *et al.* Meta-analysis of genome-wide association studies for cattle stature
731        identifies common genes that regulate body size in mammals. *Nature genetics* **50**, 362
732        (2018).
733   21   Daetwyler, H. D. *et al.* Whole-genome sequencing of 234 bulls facilitates mapping of
734        monogenic and complex traits in cattle. *Nature genetics* **46**, 858 (2014).
735   22   Hayes, B. J. & Daetwyler, H. D. 1000 Bull Genomes Project to Map Simple and Complex
736        Genetic Traits in Cattle: Applications and Outcomes. *Annual review of animal biosciences* **7**,
737        89-102, doi:10.1146/annurev-animal-020518-115024 (2019).
738   23   Xiang, R. *et al.* Genome variants associated with RNA splicing variations in bovine are
739        extensively shared between tissues. *BMC Genomics* **19**, 521, doi:10.1186/s12864-018-4902-
740        8 (2018).
741   24   Chamberlain, A. *et al.* in *11th world congress on genetics applied to livestock production*
742        *(WCGALP). Auckland, New Zealand: Volume Molecular Genetics.*  254.
743   25   Keightley, P. D. & Jackson, B. C. Inferring the probability of the derived vs. the ancestral
744        allelic state at a polymorphic site. *Genetics* **209**, 897-906 (2018).
745   26   Rocha, D., Billerey, C., Samson, F., Boichard, D. & Boussaha, M. Identification of the putative
746        ancestral allele of bovine single-nucleotide polymorphisms. *Journal of Animal Breeding and*
747        *Genetics* **131**, 483-486 (2014).
748   27   Xiang, R., van den Berg, I., MacLeod, I. M., Daetwyler, H. D. & Goddard, M. E. Effect direction
749        meta-analysis of GWAS identifies extreme, prevalent and shared pleiotropy in a large
750        mammal. *Commun Biol* **3**, 88, doi:10.1038/s42003-020-0823-6 (2020).
751   28   Erbe, M. *et al.* Improving accuracy of genomic predictions within and between dairy cattle
752        breeds with imputed high-density single nucleotide polymorphism panels. *Journal of dairy*
753        *science* **95**, 4114-4129 (2012).
754   29   Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
755        datasets. *Gigascience* **4**, 7, doi:10.1186/s13742-015-0047-8 (2015).
756   30   Chen, N. *et al.* Whole-genome resequencing reveals world-wide ancestry and adaptive
757        introgression events of domesticated cattle in East Asia. *Nature Communications* **9**, 1-13
758        (2018).
759   31   Xiang, R. *et al.* Quantifying the contribution of sequence variants with regulatory and
760        evolutionary significance to 34 bovine complex traits. *Proceedings of the National Academy*
761        *of Sciences* **116**, 19398-19408, doi:10.1073/pnas.1904159116 (2019).
762   32   Meuwissen, T., Hayes, B. & Goddard, M. Prediction of total genetic value using genome-wide
763        dense marker maps. *Genetics* **157**, 1819-1829 (2001).
764   33   Daetwyler, H. *et al.* in *Proc Assoc Adv Anim Breed Genet.*  201-204.
765   34   DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-
766        generation DNA sequencing data. *Nature genetics* **43**, 491-498 (2011).
767   35   Browning, S. R. & Browning, B. L. Rapid and accurate haplotype phasing and missing-data
768        inference for whole-genome association studies by use of localized haplotype clustering.
769        *American journal of human genetics* **81**, 1084-1097, doi:10.1086/521987 (2007).
770   36   Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
771        genomes. *Genome research* **15**, 1034-1050 (2005).
772   37   Fuchsberger, C., Abecasis, G. R. & Hinds, D. A. minimac2: faster genotype imputation.
773        *Bioinformatics* **31**, 782-784 (2014).

774  38  Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and accurate
775      genotype imputation in genome-wide association studies through pre-phasing. *Nature*
776      *genetics* **44**, 955 (2012).
777  39  Sargolzaei, M., Chesnais, J. P. & Schenkel, F. S. A new approach for efficient genotype
778      imputation using information from relatives. *BMC Genomics* **15**, 478, doi:10.1186/1471-
779      2164-15-478 (2014).
780  40  Loh, P.-R. *et al.* Reference-based phasing using the Haplotype Reference Consortium panel.
781      *Nature genetics* **48**, 1443 (2016).
782  41  Rosen, B. D. *et al.* De novo assembly of the cattle reference genome with single-molecule
783      sequencing. *GigaScience* **9**, giaa021 (2020).
784  42  Pausch, H. *et al.* Evaluation of the accuracy of imputed sequence variant genotypes and their
785      utility for causal variant detection in cattle. *Genetics Selection Evolution* **49**, 1-14 (2017).
786  43  Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex
787      trait analysis. *The American Journal of Human Genetics* **88**, 76-82 (2011).
788  44  Lee, S. H. & Van der Werf, J. H. MTG2: an efficient algorithm for multivariate linear mixed
789      model analysis based on genomic information. *Bioinformatics* **32**, 1420-1422 (2016).
790  45  Xiang, R. *et al.* Genome-wide fine-mapping identifies pleiotropic and functional variants that
791      predict many traits across global cattle populations. *Nature Communications* **12**, 860,
792      doi:10.1038/s41467-021-21001-0 (2021).
793  46  Fisher, R. A. XV.—The correlation between relatives on the supposition of Mendelian
794      inheritance. *Earth and Environmental Science Transactions of the Royal Society of Edinburgh*
795      **52**, 399-433 (1919).
796  47  Falconer, D. S. & Mackay, T. F. C. *Introduction to quantitative genetics*.  (Longman, 1996).
797  48  Lynch, M. & Walsh, B. *Genetics and Analysis of Quantitative Traits*.  (Sinauer Associates,
798      1998).
799  49  Weir, B. S. & Ott, J. Genetic data analysis II. *Trends in genetics* **13**, 379 (1997).
800  50  Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,
801      doi:10.1093/bioinformatics/bts635 (2013).
802  51  Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for
803      assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930 (2014).
804  52  Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. Voom: precision weights unlock linear model
805      analysis tools for RNA-seq read counts. *Genome Biology* **15**, R29 (2014).
806  53  Moser, G. *et al.* Simultaneous discovery, estimation and prediction analysis of complex traits
807      using a Bayesian mixture model. *PLoS genetics* **11**, e1004969 (2015).

808

## Acknowledgements

## Author contributions

M.E.G. and R.X conceived the study. R.X. performed all analyses. E.J.B. contributed to the BayesR analysis. I.M.M., assisted with data curation. S.B., C.J.J., and A.J.C. contributed to the imputation of sequence variants. R.X. and M.E.G. wrote the paper. R.X., M.E.G., E.J.B. and I.M.M. revised the paper. All authors read and approved the final manuscript.

## Competing Interests

The authors declare no competing interests.