1  **Continent-wide genomic analysis of the African buffalo (*Syncerus caffer*).**

2

3  Andrea Talenti[1,2†], Toby Wilkinson[1,2†], Elizabeth A. Cook[3,4], Johanneke D. Hemmink[1,2,3,4],
4  Edith Paxton[1], Matthew Mutinda[5], Stephen D. Ngulu[6], Siddharth Jayaraman[1], Richard P.
5  Bishop[3], Isaiah Obara[7], Thibaut Hourlier[8], Carlos Garcia Giron[8], Fergal J. Martin[8], Michel
6  Labuschagne[9], Patrick Atimnedi[10], Anne Nanteza[11], Julius D. Keyyu[12], Furaha Mramba[13],
7  Alexandre Caron[14,15,16], Daniel Cornelis[17,18], Philippe Chardonnet[19], Robert Fyumagwa[12],
8  Tiziana Lembo[20], Harriet K. Auty[20], Johan Michaux[21], Nathalie Smitz[22], Philip Toye[3,4],
9  Christelle Robert[1,2,23#], James G.D. Prendergast[1,2#], Liam J. Morrison[1,2#*].

10

11  1. The Roslin Institute, Royal (Dick) School of Veterinary Studies, University of
12     Edinburgh, Midlothian, EH25 9RG, United Kingdom
13  2. Centre for Tropical Livestock Genetics and Health (CTLGH), Roslin Institute,
14     University of Edinburgh, Easter Bush Campus, EH25 9RG, United Kingdom
15  3. International Livestock Research Institute, P.O. Box 30709, Nairobi 00100, Nairobi,
16     Kenya.
17  4. Centre for Tropical Livestock Genetics and Health (CTLGH), ILRI Kenya, P.O. Box
18     30709, Nairobi 00100, Kenya.
19  5. Kenya Wildlife Service, P.O. Box 40241, Nairobi 00100, Kenya.
20  6. Ol Pejeta Conservancy, Private Bag, Nanyuki 10400, Kenya.
21  7. Institute for Parasitology and Tropical Veterinary Medicine, Freie Universität Berlin,
22     Robert-von-Ostertag-Str. 7-13, 14163 Berlin, Germany.
23  8. European Molecular Biology Laboratory, European Bioinformatics Institute,
24     Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SD, United Kingdom.
25  9. Clinomics , Uitzich Road, Bainsvlei, Bloemfontein, 9338, South Africa.
26  10. Uganda Wildlife Authority, Kampala, Uganda.
27  11. College of Veterinary Medicine, Animal Resources and Biosecurity, Makerere
28      University, Kampala, Uganda.
29  12. Tanzania Wildlife Research Institute, Box 661, Arusha, Tanzania.
30  13. Hester Biosciences Africa Limited, P.O. Box 30216, Kibaha, Coastal Region,
31      Tanzania.
32  14. ASTRE, University of Montpellier (UMR), CIRAD, 34090 Montpellier, France.
33  15. CIRAD, UMR ASTRE, RP-PCP, Maputo 01009, Mozambique.
34  16. Faculdade Veterinaria, Universidade Eduardo Mondlane, Maputo, Mozambique.
35  17. CIRAD, Forêts et Sociétés, 34398 Montpellier, France.
36  18. Forêts et Sociétés, University of Montpellier, CIRAD, 34090 Montpellier, France.
37  19. IUCN SSC Antelope Specialist Group co-chair, 92100 Boulogne, France.
38  20. School of Biodiversity, One Health and Veterinary Medicine, College of Medical,
39      Veterinary and Life Sciences, University of Glasgow, Glasgow, United Kingdom.
40  21. Laboratoire de Génétique de la Conservation, Institut de Botanique (Bat. 22),
41      Université de Liège (Sart Tilman), Chemin de la Vallée 4, B4000 Liège, Belgium.
42  22. Royal Museum for Central Africa (BopCo), Leuvensesteenweg 13, 3080, Tervuren,
43      Belgium.
44  23. Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer,
45      University of Edinburgh, Crewe Road South, Edinburgh, EH4 2XU, United Kingdom.

46

47  †Joint first authors
48  #Joint last authors
49  *Corresponding author: liam.morrison@roslin.ed.ac.uk

50

1

**Abstract**

The African buffalo (*Syncerus caffer*) is a wild bovid with a historical distribution across much of sub-Saharan Africa. Genomic analysis can provide insights into the evolutionary history of the species, and the key selective pressures shaping populations, including assessment of population level differentiation, population fragmentation, and population genetic structure. In this study we generated the highest quality *de novo* genome assembly (2.65 Gb, scaffold N50 69.17 Mb) of African buffalo to date, and sequenced a further 195 genomes from across the species distribution. Principal component and admixture analyses provided surprisingly little support for the currently described four subspecies, but indicated three main lineages, in Western/Central, Eastern and Southern Africa, respectively. Estimating Effective Migration Surfaces analysis suggested that geographical barriers have played a significant role in shaping gene flow and the population structure. Estimated effective population sizes indicated a substantial drop occurring in all populations 5-10,000 years ago, coinciding with the increase in human populations. Finally, signatures of selection were enriched for key genes associated with the immune response, suggesting infectious disease exert a substantial selective pressure upon the African buffalo. These findings have important implications for understanding bovid evolution, buffalo conservation and population management.

**Introduction**

The African buffalo, *Syncerus caffer*, is a key member of the charismatic African megafauna, and was historically distributed across sub-Saharan Africa, inhabiting a diverse range of habitats from dry savannah to montane rainforest. Over the past century the population density and distribution has been much reduced. The population range has also become increasingly fragmented due to man-made pressures, resulting in approximately 70% of the global population being restricted to protected areas [1-3].

The species has been historically divided into varying numbers of subspecies based upon distribution, habitat and morphology, the most recent update of the IUCN Red List recognising *S. caffer caffer* (Eastern and Southern African savannah), *S. c. brachyceros* (Western African savannah), *S. c. aequinoctialis* (Central African savannah), and *S. c. nanus* (Western and Central African forest) [4]. The genetic understanding of population diversity and structure across the species range mostly derives from the application of low resolution tools, such as mitochondrial D-loop sequences, microsatellites and mitogenomes [5-7], with a more recent study using genome-wide single-nucleotide polymorphisms (SNPs) [8]. The two studies to analyse diversity at the genome level, focused on South African *S. c. caffer* animals (n=40) in protected areas [9] and *S. c. caffer* populations (n=59) from East and Southern Africa [10]. These studies have collectively highlighted that the current subspecies classification may not be supported by genetic data, and that there is striking population substructuring within and between the putative subspecies. They have also indicated concerns with respect to low effective population sizes in increasingly isolated populations in some African regions. Improved genetic tools can potentially contribute to conservation management strategies, both in terms of restoring connectivity between relevant populations in order to improve or restore genetic diversity, and avoiding loss of genetic integrity (i.e. maintenance of genetic diversity relevant to local environmental adaptation) through uninformed population mixing (e.g. translocations) [6, 11, 12].

As well as being an iconic species of African wildlife, the African buffalo is the closest bovid relative of domesticated cattle (*Bos taurus taurus & Bos taurus indicus*) in Africa. The African buffalo has co-evolved in Africa with pathogens responsible for important and

99   impactful diseases of cattle such as animal African trypanosomiasis [13] and foot and mouth
100  disease virus (FMD) [14, 15]. For trypanosomiasis, in contrast to the often devastating impact
101  that infection has on cattle, African buffalo are largely tolerant, displaying much less severe
102  clinical signs (e.g. [16, 17]). Additionally, African buffalo are the primary host for the tick-
103  borne protozoan *Theileria parva*, the causative agent of East Coast fever, an often deadly
104  disease in cattle that is asymptomatic in buffalo [18]. These diseases have impeded
105  productivity and the expansion of African pastoralists and their cattle for centuries [19, 20].
106  During the colonial era, European cattle also brought with them diseases then exotic to
107  Africa, such as rinderpest, brucellosis and bovine tuberculosis [21], to which African buffalo
108  are susceptible. African buffalo and cattle co-exist today across many wildlife/livestock
109  interfaces that enhance mutual pathogen transmission [22], and this can result in imposition
110  of strict veterinary controls at these interfaces that often impact local livelihoods and
111  conservation efforts (e.g., [23, 24]). This makes the buffalo particularly interesting in terms
112  of host-pathogen coevolution and potentially providing a route to identifying host genes and
113  pathways relevant to controlling these diseases in livestock.

114  This study aimed to develop a reference genome for the African buffalo, as a foundation to
115  analyse the population genomic structure across the current distribution of the species in sub-
116  Saharan Africa. Two reference genomes have previously been published, but were generated
117  via short read sequencing, resulting in relatively fragmented final genome assemblies
118  (scaffold N50s of 2.40 Mb and 2.32 Mb, respectively) [25, 26]. Using a combination of long
119  read (PacBio) and Hi-C sequencing, we generated and *de novo* assembled a substantially
120  higher quality and more contiguous reference genome of 2.65 Gb, with a scaffold N50 of
121  69.17 Mb. We then sequenced the genomes of 196 African buffalo samples from across the
122  current species distribution, which enabled the analysis of genetic substructure, admixture
123  between populations, and effective population sizes. We also assessed *S. caffer* genomes for
124  signatures of selection, highlighting genes that may be responsible for environmental
125  adaptation, in particular against diseases important for both buffalo and cattle.

126  **Results**

127  *Assembly statistics*

128  We first generated a *de novo S. c. caffer* reference genome from a male buffalo (OPB4)
129  sampled in Ol Pejeta Conservancy, Kenya, providing the foundation to enable the
130  characterisation of the genetic diversity of African buffalo populations both in terms of their
131  geographic regions and habitats and their current subspecies classification. We applied a deep
132  sequencing strategy, based on a combination of 60x long read (PacBio) and 75x short read
133  (Illumina) reads, to generate a *de novo* reference genome ensuring high per base sequence
134  quality and consensus to achieve good genome contiguity, with an N50 of 69.16Mb. The long
135  reads were assembled using FALCON (Dovetail Genomics) and polished using Arrow.
136  Contigs were then scaffolded using ~393 million 2x150bp Illumina reads of HiC data, using
137  the HiRise software. Gaps in the draft genome were addressed using PBJelly [27]. Finally,
138  Pilon [28] was used for sequential rounds of polishing, each of which was assessed for its
139  resulting assembly quality over previous rounds. The genome following four rounds of
140  polishing displayed the highest assembly statistics, with a total of 3,351 scaffolds, a total
141  length of 2.65 Gb (comparable to 2.72 Gb for the *Bos taurus* genome), a scaffold N50 of
142  69.16 Mb and a quality value (QV) of 35.9, indicating ~1 error every 5,000 bp. The assembly
143  statistics are summarised in Figure 1 and Supplementary Data 1.

144  Previous African buffalo reference genomes, generated by Glanzmann et al. [25] and Chen et
145  al. [26], were based solely on Illumina short read sequencing, which led to highly fragmented
146  assemblies of 442,401 scaffolds with a scaffold N50 of 2.40 Mb, and 150,000 scaffolds with

147 an N50 of 2.30 Mb, respectively. These very fragmented assemblies provided limited scope
148 for downstream analysis of variants and their predicted effects on functional regions, i.e.
149 annotated genes and regulatory regions (a comparison of the three genome assemblies is
150 illustrated in Figure 1).
151

152 *Transcriptome analyses and genome annotation*

153 To enable in depth characterisation of the African buffalo transcriptome and to facilitate the
154 annotation of gene isoforms, we performed full length isoform sequencing (Iso-Seq) across
155 samples from six different tissues (prescapular lymph node, testis, liver, kidney, lung and
156 spleen) collected from the same animal for which the genome was assembled (OPB4). In
157 total 51,521 distinct, high quality isoforms (defined as being supported by at least two full
158 length reads and with >99% base composition accuracy) were detected across these samples
159 (median of 11,520 per tissue, maximum of 27,271 in the testis). Complementing these data,
160 we also generated Illumina RNA-seq data, from the same animal, from eight tissues (heart,
161 prescapular and inguinal lymph nodes, testis, liver, kidney, lung and spleen). All
162 transcriptomic data were deposited to ENA
163 (https://www.ebi.ac.uk/ena/browser/view/GCA_902825105.1) with accession numbers
164 PRJEB36587 and PRJEB36588 for RNA-seq and Iso-Seq, respectively. Together these data
165 have been used to provide a high quality annotation of the buffalo assembly which can be
166 accessed through the Ensembl Rapid Release genome browser:
167 https://rapid.ensembl.org/Syncerus_caffer_GCA_902825105.1.

168 *African buffalo-specific sequence*

169 After aligning the African buffalo genome to eight high quality assemblies of four different
170 Bovidae species (cattle, water buffalo, yak and goat [29-34]), portions of the *S. caffer*
171 genome that did not match any regions in the other assemblies were ascertained. This process
172 identified a total of 24,336,918 intervals, for a total of 145,050,830 bp of sequence not
173 identified in the other eight assemblies. This includes both small variations (e.g. SNPs, small
174 indels), unplaced contigs without alignments to any other genome, and large portions of the
175 genome lacking any alignment.

176 We then refined the region selection by filtering out shorter intervals (<60 bp) and regions
177 defined as too close to a telomere (<10 Kb) or to a gap (<1 Kb), leaving a total of
178 113,654,400 bp in 81,357 fragments longer than 60 bp, which were neither telomeric nor
179 neighbouring an assembly gap. These regions have an average length of 1,397bp (3772.4 bp
180 SD) and a median size of 286bp (min. 61bp, max 308,890bp). The majority of the regions
181 (74,659 fragments accounting for 112,762,919 bp) represent sequence not found in any of the
182 other species genomes considered in the study, whereas the remaining are classified as
183 divergent haplotypes. Of the 113Mb, a total of 64.9Mb (57.1%) are putatively identified as
184 repeats using RED [35]. To rule out the possibility of these novel regions being due to
185 contamination, we confirmed the coverage of these regions was consistent with the rest of the
186 genome, using short-read whole genome sequencing data from 46 samples from the
187 population analysis (see section below; Supplementary Figure 1).

188 HOMER analysis considered 4,286/7,096 sequences with less than 60% of masked
189 nucleotides. These sequences presented 38 motif types enriched (P-value <1e-5), such as the
190 FOSL2/MA0478.1/Jaspar (0.661) motif, a negative regulatory sequence in the
191 differentiation-sensitive adipocyte gene (aP2), a motif identified as a transcriptional enhancer
192 for the Gibbon ape leukaemia virus, and which is also in a region of the human
193 immunodeficiency virus (HIV) [36]. We performed the feature analysis on the annotation
194 generated by Ensembl from the Iso-Seq sequencing data previously described. We identified

195  7,096 annotated genes and 131 pseudogenes overlapping the novel regions, of which 583
196  genes, 194 ncRNA genes and 71 pseudogenes were entirely included in the identified regions
197  (Supplementary Table 1). A total of 317 of 583 genes had at least one biological term
198  annotated. GO terms definitions were fetched using the goatools python package [37]. Out of
199  4,088 terms in the background dataset, 17 (15 GO terms and 2 KEGG pathways) were found
200  significantly enriched. Among the significant terms was the defence response GO term
201  (GO:0006952, FDR-corrected P-value: 0.0189, Supplementary Table 1), described as the
202  response triggered by the presence of a foreign body.

203  **Population genetics**

204  To better understand African buffalo genetic diversity, we generated short read sequencing
205  data for a further 195 animals deriving from across the continental range of the species (at a
206  coverage of 15x for 146 samples, and 30x for 50 samples; Table 1 & Figure 2A; for full
207  sample list and metadata see Supplementary Table 2). This included samples from the
208  currently described four subspecies; *S. c. caffer*, *S. c. nanus*, *S. c. brachyceros* and *S. c.*
209  *aequinoctialis* (Table 1 & Figure 2A), and two putative *S. c. nanus* and *S. c. aequinoctialis*
210  hybrids (based upon morphology and geography at time of sampling - labelled as
211  'intermediates'). Together, these samples derived from 21 sites/localities or protected areas
212  across 12 different countries. We performed phylogenetic and population analyses including
213  only samples with a high call rate (>85%), and analysing only the biallelic polymorphic
214  SNPs (minor allele frequency >5%), as well as only considering unrelated individuals
215  (samples fourth degree or greater).

217  As can be seen in Figure 2B the genetic relationships between the samples largely mirrors
218  their geographic origin, with the first principal component (PC1) reflecting differentiation
219  between samples from Eastern/Southern and Western Africa, which corresponds to a split
220  between the Western/Central African subspecies (*S. c. aequinoctialis*, *S. c. brachyceros* and
221  *S. c. nanus*) and Eastern/Southern African *S. c. caffer*. The second component (PC2)
222  correlates with differentiation between *S. c. caffer* samples from the Northern part of the
223  subspecies' range (Kenya, Tanzania, Uganda) compared to *S. c. caffer* samples from
224  Southern Africa. Notably, there was a clear signature of geography within the *S. c. caffer*
225  data, with each geographic sub-population forming a distinct cluster in the PCA and a cline
226  observed from Uganda to Kenya and Tanzania in the North, through Mozambique to samples
227  from Botswana and Zimbabwe, and finally South Africa in the South. In Western/Central
228  Africa, *S. c. aequinoctialis*, *S. c. brachyceros* and *S. c. nanus* sub-populations also formed
229  separate clusters, although the *S. c. nanus* and *S. c. brachyceros* populations clustered closely
230  together. Samples were initially grouped by sub-species and country of sampling. However,
231  based on PCA results, the Tanzania and Kenya, and Botswana and Zimbabwe samples were
232  grouped together, reflecting their geographic proximity. This resulted in nine subgroups for
233  downstream analyses; referred to hereafter as *S. c brachyceros*, *S. c. nanus*, *S. c.*
234  *aequinoctialis*, intermediate (putative hybrids between *S. c. nanus*, *S. c. aequinoctialis*), *S. c.*
235  *caffer* Uganda, *S. c. caffe*r Kenya/Tanzania, *S. c. caffer* Mozambique, *S. c. caffer*
236  Zimbabwe/Botswana and *S. c. caffer* South Africa. Population sample sizes post-filtering
237  ranged from 2 for the *S. c. nanus* spp to 48 for the *S. c. caffer* from Tanzania (see Table 1),
238  leaving a total of 163 samples for the phylogenetic analyses.

239  In order to explore the relationship between these populations further, and to mitigate the
240  different sample size between subpopulations resulting in over-representation of population-
241  specific variation in the dataset as far as possible [38], we downsampled the larger groups to
242  15 representative samples (for those with less, all samples were included). Since the *S. c.*
243  *brachyceros* population had a total of 16 samples, we did not perform any reduction on this

244 population. This resulted in a subset of 95 individuals to be considered for the population
245 genetic analyses (see Supplementary Table 2 for samples included in these analyses). As
246 shown in the principal components analysis (PCA) pre- and post-reduction (Supplementary
247 Figure 2), the general structure of the sample was not affected by the subsampling.

248 Bootstrapped admixture analyses identified three clusters as a parsimonious solution for the
249 number of subpopulations (i.e. lowest iteration number and reduced increase in CV error;
250 Supplementary Figure 3), representing the East, West and Southern African high level
251 groupings. At K=9 the clusters recapitulate the nine sub-groupings described above (see
252 Supplementary Figure 4 for admixture results at multiple K). The same high-level structure is
253 reflected in the 100-bootstrap identity-by-state phylogenetic tree (Figure 3B).

254 Comparison of the genetic diversity between all pairs of populations (as represented by the
255 $F_{ST}$ statistic) highlights that this is largely a function of physical distance, i.e. the diversity
256 observed between two populations increases broadly linearly with increasing distance
257 between them (Figure 3C, Mantel test r: 0.65, p=0.0018; underlying $F_{ST}$ data detailed in
258 Supplementary Table 3). However, sub-structure in this isolation-by-distance analysis is
259 observed. After excluding the *S. c. caffer* Hluhluwe-Umfolozi and *S. c. nanus* populations,
260 the relationship is even stronger, and variation in the $F_{ST}$ values between the remaining
261 groups can potentially largely all be explained by the distances between them (red line in
262 Figure 3C, Mantel test r: 0.96, p=0.0013). This is consistent with the idea that these African
263 buffalo have historically formed large continuous groups of populations with differentiation
264 between populations simply reflecting the reduced mating probability with increasing
265 distance. *S. c. nanus*, the forest buffalo, shows an unusually steep increase in differentiation
266 relative to other populations (blue line in Figure 3C). This could be for a variety of reasons,
267 including geographical barriers reducing the gene flow between this group and the others
268 analysed. Animals found at the same location should exhibit little differentiation, and
269 consistent with this, the intercept of the slopes is not significantly different from 0 in these
270 comparisons (both linear regression intercept P>0.4), i.e. when comparing *S. c. nanus* to other
271 populations or the non- *S. c. nanus* and non-*S. c. caffer* Hluhluwe-Umfolozi populations to
272 each other. However, this is not the case for comparisons involving the South African *S. c.*
273 *caffer* Hluhluwe-Umfolozi population. Under the assumption of a simple linear relationship
274 between genetic differentiation and geographic distance, the predicted level of diversity at a
275 distance of 0 km is significantly higher than 0 (green line in Figure 3C, linear regression
276 intercept $P=2.7 \times 10^{-4}$). This suggests, that unlike in the other population comparisons, there is
277 elevated differentiation between this population and others, above and beyond that expected
278 from their geographic distance apart. This may reflect an isolation event with respect to the
279 Hluhluwe-Umfolozi population.

280 EEMS analysis (Figure 4A) adds to this picture of continental gene flow, with the Congo
281 river basin likely representing a significant barrier of migration, particularly between
282 Western/Central African *S. c. nanus* and *S. c. caffer* populations in Eastern Africa. The data
283 also suggest that the Rift Valley potentially presents a geographical barrier to gene flow
284 within the African buffalo.

285 The Relate software and genome-wide genealogies were used to estimate population-specific
286 population sizes over time for the largest buffalo groupings (*S. c. caffer* only: Uganda,
287 Tanzania/Kenya and Zimbabwe/Botswana – grouped as defined by PCA, admixture and
288 phylogenetic analyses; Figures 2 and 3). As shown in Figure 4 there has been a sharp
289 reduction in the estimated effective population sizes across these groups in the last
290 approximately 10,000 years, broadly mirroring the expansion of human effective population
291 sizes over a similar time-period (Figure 4B). There were not sufficient numbers in all

292    individual populations for robust $N_e$ analyses, but for the populations that did have sufficient
293    numbers, contemporary $N_e$ estimates were ~1,300, 2,000 and 3,000 for Uganda,
294    Tanzania/Kenya and Zimbabwe/Botswana, respectively. These data suggest that the effective
295    population sizes of these Eastern and Southern African *S. c. caffer* are above the levels of
296    conservation concern. Coelescence estimates are shown in Supplementary Figure 5.

297    However, analysis of all populations highlights that the *S. c. nanus* and South African *S. c.*
298    *caffer* Hluhluwe-Umfolozi samples have high levels of homozygosity ($F_{ROH}$ of 0.29 and 0.36
299    compared to a range of 0.12 to 0.21 for the other populations; Figure 3D). This is consistent
300    with the known extreme bottlenecks experience by the Hluhluwe-Umfolozi buffalo
301    population [9]; the *S. c. nanus* samples derive from Lekedi NP in Gabon, and we are unaware
302    of historical population-level data that would inform of bottlenecks – while the homozygosity
303    analysis is obviously on individual genomes, with this population we would caution
304    overinterpretation as we only have data from two individuals.

305    **Selective sweeps**

306    African buffalo are exposed to a range of different environmental pressures across their
307    distributional range, including a range of pathogens that also impact domesticated bovids
308    such as cattle. To investigate selective sweeps between and within the nine population
309    groupings we calculated the XP-EHH and $P_R$ Relate Selection Test statistics [39, 40]. Due to
310    being more susceptible to artefactual results deriving from smaller sample sizes than the XP-
311    EHH statistic, the calculation of the $P_R$ statistic was restricted to just the populations with
312    more than 20 samples after filtering for relatedness (i.e. the Uganda, Zimbabwe/Botswana
313    and Tanzanian/Kenyan populations). These two tests are complementary in that whereas the
314    XP-EHH statistic tests for differences in haplotype homozygosity between populations, $P_R$
315    characterises the speed of spread of particular genomic lineages within a population, relative
316    to others. Supplementary Table 4 summarises the results of these two tests. In total, 73 loci of
317    elevated XP-EHH levels overlapping a gene were identified in at least one population
318    comparison, and 34 $P_R$ significant loci were detected in one of the three studied populations.
319    Of the XP-EHH loci, 9 also overlapped a significant $P_R$ peak (Supplementary Table 4). These
320    9 loci spanned 11 genes, with several having strong links to immune response, including
321    putative killer cell immunoglobulin-like receptor like protein KIR3DP1 (LOC102402296), T
322    cell receptor beta variable 5-1-like (LOC112577699), the major histocompatibility complex
323    gene TRIM26 and N-acetylneuraminic acid phosphatase (NANP). The latter is involved in
324    sialic acid synthesis which in turn is linked to immune response modulation, and NANP has
325    also been observed to be under recent positive selection in both humans and cattle [41, 42].
326    Two further of these nine genes linked to both XP-EHH and $P_R$ peaks in African buffalo were
327    also previously linked to recent positive selection in water buffalo [42], namely myeloid-
328    associated differentiation marker-like (LOC102403696) and tyrosine-protein phosphatase
329    non-receptor type substrate 1-like (SIRPA-like) gene (LOC102396916). LOC102396916 was
330    associated with significant $P_R$ peaks in both the Uganda and Tanzania/Kenyan populations
331    and also elevated XP-EHH scores in the South African *S. c. caffer* vs intermediate and *S. c*
332    *aequinoctialis* populations (Figure 5; Supplementary Figure 6). SIRPA is an
333    immunoglobulin-like cell surface receptor for CD47 (a cell surface protein that is involved in
334    the promotion/regulation of cellular proliferation) and has been associated with a range of
335    infectious diseases, including *Theileria annulata* infection in cattle [43] (*T. annulata* being
336    the causative agent of tropical theileriosis across North Africa and Asia, and is closely related
337    to *Theileria parva* found in Eastern Africa). This gene has previously been identified to be
338    associated with selective sweeps between water buffalo breeds (elevated XP-CLR statistics
339    between Mediterranean and Jaffrabadi, and Pandharpuri and Banni water buffalo breeds
340    [42]). Characterisation of this gene's expression profile in the water buffalo expression atlas

7

341  highlighted that it falls within a macrophage-specific cluster of genes [44]. Together these
342  results therefore point towards this gene being a potentially important target of selection
343  across bovids due to its role in immune response. Consequently, five of these nine genes
344  under putative selection in African buffalo show strong links to immune response, with two
345  of the remaining genes being uncharacterised and their function being unknown.

346

347  **Discussion**

348  *African buffalo genome*

349  The genome generated in this study represents a substantial improvement on current genomic
350  resources available for *S. caffer*, with greater contiguity and much improved assembly and
351  annotation – this, and the allied gene expression datasets, will hopefully serve as useful
352  resources for the bovid and African buffalo research communities. The genome assembly is
353  currently at the scaffold rather than chromosomal level, and so karyotype and features such as
354  centromeres remain undefined, and the genome also contains Y chromosome and
355  mitochondrial sequences that have not been completely resolved. There is therefore clearly
356  scope for further improvement of the reference genome. An interesting finding was the
357  African buffalo-specific sequence, which was identified after aligning the African buffalo
358  genome to eight existing high quality bovid genome assemblies (cattle, water buffalo, yak
359  and goat [29-34]). *S. caffer* sequences that that did not match any regions in the other
360  assemblies were defined as African buffalo-specific sequence. These sequences were
361  validated by assessing coverage of these African buffalo-specific sequences in randomly
362  selected short read data from the population data, based on the expectation that if these were
363  genuine African buffalo-specific sequence there would be coverage detected in multiple
364  samples, and this was indeed the case. While 57.1% of these African buffalo-specific
365  sequences are repeats, there are 583 genes, 71 pseudogenes, and 194 ncRNAs that are
366  entirely within the identified regions. These were enriched for genes associated with the host
367  defence, and the genes within these regions would clearly be of interest in further studies to
368  identify traits that may be relevant to these African buffalo-specific sequences.

369  *Population genomic structure: taxonomic insights*

370  The admixture analysis suggests that the *S. caffer* population splits into three high-level
371  lineages, with a further nine subgroupings apparent that correlate with geographical location.
372  There is little support in our data for the current classification of the four IUCN recognised
373  subspecies; *S. c. caffer* (Eastern and Southern African savannah), *S. c. brachyceros* (Western
374  African savannah), *S. c. aequinoctialis* (Central African savannah) and *S. c. nanus* (Western
375  and Central African forest), with *S. c. brachyceros* and *S. c. aequinoctialis* sometimes being
376  lumped and treated as a single subspecies, viz. *S. c. brachyceros*. Historically these
377  classifications have been based on a combination of geographical distribution, habitat
378  preferences and morphological features. *Syncerus c. nanus*, the forest buffalo, is the most
379  divergent morphologically, being on average much smaller, predominantly rufous in colour
380  as opposed to black, and with a different horn shape. From this perspective it is perhaps
381  surprising that we could not detect substantial genetic divergence from the Western/Central
382  African savannah buffalo. However, this finding agrees with previous genetic analyses using
383  mitochondrial D-loop sequence markers, which similarly indicated a lack of support for
384  differentiation between Western/Central African 'subspecies' [5]. However, the limited
385  number of samples assigned to *S. c. nanus* did not enable balanced analyses, with in general a
386  smaller number of populations sampled for the Western and Central African regions

8

387 compared to Eastern and Southern Africa. This may have resulted in some bias in our
388 population analyses. However, we did attempt to mitigate this bias to some extent by
389 reducing populations to 10 samples per population where relevant and possible.
390 Notwithstanding this fact, the present database provides genome-level and -wide resolution
391 on variation (based upon 23,454,419 identified variants relative to the assembled reference
392 genome); a much more robust basis for identifying genetic differentiation than previous
393 methods used to identify genetic substructuring in this species. These insights have parallels
394 with previous genome/multilocus genetic data studies on African ungulates with similar pan-
395 Sub-Saharan distribution, the giraffe (*Giraffa camelopardis*) and zebra (*Equus quagga*),
396 which indicated a lack of correlation of genetic data with morphology-based speciation, in
397 those cases resulting in the identification of cryptic speciation [45, 46].

398 Admixture and EEMS analyses indicate that the population genomic structure is shaped by
399 geographical barriers, which limit where migration and therefore where lineage and
400 population mixing can happen. This is evidenced by Ugandan buffalo demonstrating ancestry
401 from both Eastern and Western African populations, and there being some signal of East
402 African ancestry in Central African buffalo (*S. c. aequinoctialis*, *S. c. nanus* and intermediate;
403 Figure 3A and Supplementary Figure 4). Both admixture and EEMS data indicate that
404 Uganda is likely to act as an interface zone between these lineages, although further sampling
405 in relevant populations (for example, known buffalo populations in Eastern CAR and DRC,
406 South Sudan and Western Ethiopia) would help resolve the extent of genetic flow. EEMS
407 analyses suggests that the divergence between the East and West lineages was most likely
408 driven by geography, with the Congo Basin and River effectively creating a barrier to North-
409 South gene flow in the West of the continent, and Uganda being the pinch point at which
410 Central African savannah and forest populations can intersect with Eastern African savannah
411 buffalo.

412 The driving forces shaping the differentiation between Northern and Southern populations of
413 *S. c. caffer* (i.e. between the Kenyan, Ugandan and Tanzanian cluster and the Mozambique,
414 Botswana, Zimbabwe and South Africa cluster) is less clear from our analyses. A potential
415 role of the Great Rift Valley acting as historical barrier to gene flow has been suggested
416 within other large savannah mammals [47-49]. However, all Tanzanian samples included in
417 the present study originated from the North of the country (the closest population in the
418 Southern cluster being Niassa Special Reserve in Mozambique –approximately 1,000 km
419 from the Northern Tanzanian parks); additional samples from Central and Southern Tanzania
420 where substantial buffalo populations exist (e.g. in Ruaha and Nyerere NPs) could potentially
421 identify animals that are genetically intermediate between the 'Northern' and 'Southern'
422 clusters, and reveal that there is a steady cline of differentiation within *S. c. caffer* from North
423 to South, as supported by the isolation-by-distance analysis. The data are consistent with the
424 findings of a previous genomic study of *S. c. caffer* across its range, which also concluded
425 that there was a primary split between northern and southern *S. c. caffer* populations
426 approximately 50,000 years ago, followed by gene flow [10].

427 *Effective population sizes*

428 Although effective population size estimates are difficult to estimate accurately and can be
429 confounded by population structure, the effective population size data interestingly suggests a
430 coincident drop in $N_e$ with the rise in human $N_e$ (obtained through the 1,000 Genomes data
431 [50]). This is observed in similar analyses applied to both other individual African ungulates
432 (giraffe) [51] and collated global ruminant data [26]. In the case of African buffalo, previous
433 studies based on both microsatellite and mitochondrial DNA data have suggested an

9

434    expansion approximately 80,000 years ago coincident with the spread of grassland habitat,
435    which was followed by a significant decline ~3-7,000 years ago, probably resulting from an
436    overall increase in arid areas across Africa that are inhospitable to African buffalo [7, 52, 53]
437    – our findings are also consistent with these conclusions.  For the African buffalo, it was
438    anticipated that the greater resolution provided by genomic data may detect a drop in $N_e$
439    observed as a result of the rinderpest virus epidemic of the 1890s [54], which anecdotally
440    caused very high mortality of the buffalo populations through Eastern and Southern Africa in
441    particular [55, 56]. However, given the relatively recent timing of the rinderpest epidemic
442    and the fact that the $N_e$ was reducing across the relevant timeframe in our analysis, from the
443    genome data we are not able to infer the impact of rinderpest upon population sizes. Other
444    analyses using lower resolution genetic markers [53, 57, 58] were also not able to detect a
445    drop in $N_e$ that correlated with the timing of the rinderpest epidemic, although a recent
446    genomic study using samples from *S. c caffer* did identify a very significant drop in $N_e$ over
447    the past 500 years, which could plausibly be explained by rinderpest [10] – notably the
448    decline was particularly steep in samples from Hluhluwe-Umfolozi. While we did not have
449    sufficient numbers in each population to robustly test $N_e$, for the closely related groupings of
450    Uganda, Tanzania/Kenya and Zimbabwe/Botswana $N_e$ estimates were approximately of
451    1300, 2,000 and 3,000 individuals, respectively. In these clusters at least, there is limited
452    evidence for inbreeding depression, in agreement with previous studies [6]. However, the *S.
453    c. nanus* and South African *S. c. caffer* Hluhluwe-Umfolozi samples showed high levels of
454    homozygosity, meaning that further population-specific work is required in order to assess
455    inbreeding risk. The *S. c. caffer* Hluhluwe-Umfolozi population is known to derive from very
456    small number of founder animals, and our finding is in agreement with previous data that has
457    indicated high inbreeding coefficients and low genome-wide heterozygosity levels in this
458    population [9, 10].

459    While we have very limited numbers of *S. c. nanus* samples, the finding of high levels of
460    homozygosity may be explained by the very different features of forest buffalo behaviour, in
461    that relative to savannah buffalo forest buffalo have smaller home ranges, shorter daily
462    movements, negligible seasonal movements and live in significantly smaller group sizes [2].
463    This is linked to the forest habitat likely generally acting as a greater barrier to gene flow than
464    savannah environments, limiting migration/dispersal and resulting in comparatively small
465    and isolated populations [5]. Genetic diversity metrics such as heterozygosity/homozygosity
466    and effective population size will clearly be an important feature for future studies,
467    particularly where there are increasingly fragmented and isolated populations, as is the case
468    for the West African Savannah buffalo.

469    *Selective Sweeps*

470    The selective sweep analyses identified tyrosine-protein phosphatase non-receptor type
471    substrate 1-like (SIRPA-like) as being under selection, independently detected using two
472    distinct and complementary methodologies ($P_R$ and XP-EHH), and across several population
473    groupings (Ugandan, Tanzanian/Kenyan, South African *S. c. caffer*, intermediate and *S. c
474    aequinoctialis* populations). The same locus was identified in selective sweep analyses of the
475    Asian buffalo *Bubalus bubalis* [42], and expression analysis in this species identified
476    upregulated gene expression in a macrophage-specific cluster. Interestingly SIRPA has been
477    associated with *Theileria annulata* infection in cattle [43], and its gene expression has been
478    shown in independent studies to be significantly upregulated in host cells following infection
479    and the cellular transformation associated with *T. annulata* infection [59, 60]. While SIRPA
480    will clearly be involved in the immune response to other pathogens, it is notable that *B.
481    bubalis* is the primary host of *T. annulata* (the tick-borne causative agent of tropical
482    theileriosis across North Africa and Asia). *Syncerus caffer* is similarly the primary host for

10

483    the related parasite *Theileria parva* (and the related *Theileria* sp. buffalo [61]), and it is
484    therefore plausible to link the described function of this gene with the long co-existence and
485    co-evolution of *S. caffer* with *T. parva.* Although only the Ugandan, Tanzanian/Kenyan and
486    South African *S. c. caffer* populations are within the current distribution of the tick vector
487    (*Rhipicephalus appendiculatus*) of *T. parva*, the historical range and selection of *T. parva*
488    cannot likely be inferred by the current vector distribution. Several other genes detected in
489    the selective sweep analysis have been implicated in the host response to apicomplexan
490    protozoa (which includes *Theileria* species), which lends credence to the hypothesis that the
491    ancient co-evolution and selection pressure exerted by *T. parva* in *S. caffer* may have played
492    a role in shaping the patterns of diversity in relevant regions of the current *S. caffer* genome.
493    The long relationship between *T. parva* and *S. caffer* is reflected in the limited pathology
494    caused by infection of *T. parva* in *S. caffer*, which is in stark contrast to the severe and often
495    fatal disease caused by *T. parva* infection in other hosts such as domestic cattle [18, 62]. The
496    latter have only co-existed with *T. parva* for 5,000-10,000 years [63]. This finding may
497    provide a route to identifying genes and pathways important in controlling disease during
498    infections by *Theileria* species, that can, for example, be translated to mitigating the effect of
499    these pathogens upon cattle or Asian buffalo owned by resource-poor farmers.

500    *Conclusion*

501    For the first time we have analysed genome-level data from all extant recognised African
502    buffalo subspecies, covering the majority of the remaining geographical distribution of the
503    species. Our findings demonstrate that the African buffalo is composed of three main
504    lineages, that further subdivide based on geographical location.  While current subspecies
505    nomenclature is likely to still have utility in terms of Management or Conservation Units,
506    more samples and data, particularly from *S. c. nanus*, *S. c. brachyceros* and *S. c.*
507    *aequinoctialis*, would help resolve the status of taxonomic units across the population range
508    of African buffalo. The data also demonstrated that genetic connectivity between populations
509    has historically been constrained by geographical barriers that have shaped the modern
510    population structure (particularly the Congo basin), and that human influence has been for
511    ~10,000 years and remains a main pressure on effective population size and population
512    fragmentation. While most populations do not show signs of inbreeding, particular
513    populations do, and this has implications for conservation and management of the species.
514    Finally, through analyses of selective sweeps, we identified infectious diseases as a likely
515    substantial contributor to historical selection, and hypothesise that protozoan pathogens for
516    which the buffalo has been primary host for millennia may be responsible for driving some of
517    this selection.
518
519    **Materials & Methods**

520    *Sample collection*

521    DNA samples were obtained through (1) active sampling of animals for this project; this was
522    done in collaboration with the Kenya Wildlife Service at the Ol Pejeta Conservancy, Kenya,
523    or (2) secondary use of DNA samples previously collected; this included samples previously
524    collected and published from Tanzania [14], Uganda [64], and Mozambique, Botswana,
525    Zimbabwe, South Africa, Niger, Burkina Faso, Gabon Central African Republic and Chad [5,
526    6, 8]. For sample collection in Kenya, buffalo were darted and sedated by qualified veterinary
527    personnel from KWS, and 10 ml blood collected into Paxgene Blood DNA tubes from
528    peripheral venous sampling. DNA was extracted from the Paxgene Blood DNA tubes using
529    the Paxgene Blood DNA kit (Qiagen) according to the manufacturer's instructions. Tissue
530    pieces (OPB4) were snap frozen in liquid nitrogen in the field. Tissue pieces were

531 homogenised using mortar and pestle over liquid nitrogen. The powder was resuspended in
532 Trireagent (Sigma-Aldrich) and RNA was isolated using the RNeasy kit (Qiagen) according
533 to the manufacturer's instructions.

534 Relevant research approvals were obtained in all instances; for the active sampling within this
535 study, approval was obtained from the Kenya Wildlife Services (permit number
536 KWS/BRM/5001). For secondary use of DNA samples previously collected, relevant permits
537 are Tanzania Wildlife Research Institute and Tanzania Commission for Science and
538 Technology (permit number 2021-262-NA-2021-066) [14] and Uganda Wildlife Authority
539 (permit number COD/96/05) [64], or details are provided in [5, 6, 8].

### Genome sequencing

541 For the reference genome, a buffalo sample from Ol Pejeta in Kenya (OPB4) was sequenced
542 using a combination of Illumina HiSeq (Dovetail Genomics & Edinburgh Genomics) and
543 Pacific BioSciences approaches (Dovetail Genomics & Edinburgh Genomics) to a final
544 sequencing coverage of 75x (Illumina) and 60x (PacBio). The same sample was also
545 sequenced using Illumina Hi-C (Dovetail Genomics) in order to facilitate scaffolding. For
546 the population samples, approximately 2.5 µg of total DNA from 196 animals sampled across
547 Africa (Kenya, Uganda, Tanzania, Mozambique, Botswana, Zimbabwe, South Africa, Niger,
548 Burkina Faso, Gabon, Central African Republic and Chad; Table 1, Supplementary Table 2)
549 was subjected to whole-genome sequencing by Illumina HiSeq; this was performed at a
550 coverage of 30x for 50 samples from Tanzania, with the remaining samples being sequenced
551 at 15x.

### Genome assembly

553 A primary assembly of the single molecule PacBio sequencing from OPB4 (mean read
554 lengths > 10Kb) was generated using FALCON and consisted of 7,269 contigs and an N50 of
555 1.9 Mb. This primary assembly was scaffolded using the Hi-C libraries and the HiRise
556 software by Dovetail. The resulting scaffold-level assembly was further improved via gap
557 filling and polishing steps performed with PBJelly [27] and Pilon [28] respectively, as
558 described below. Gap filling: 7,085 gaps (both inter- and intra-scaffolds) were identified in
559 the scaffold-level assembly. A total of 78 inter-scaffold gaps were partially filled (i.e.
560 extended on one side) using PBJelly, with 476,665 bases added in total, while none of the
561 identified gaps were fully closed. This observation confirmed the high quality of the primary
562 assembly achieved from PacBio reads including a post-processing step using Arrow (part of
563 the GenomicConsensus package from PacBio). Polishing: An additional 75x Illumina short
564 read sequencing (101bp paired-end reads) of DNA from the same individual used to build the
565 reference genome assembly (OPB4), was used to polish the *de novo* scaffold-level reference
566 genome assembly. Polishing allows the correction of artefacts due to sequencing errors in
567 assemblies, using the pile up of short reads that are associated with low sequencing error
568 (~1%). This process was performed multiple times and improvement upon quality metrics
569 (i.e. reduced numbers of ambiguous bases, corrected SNPs, resolved small indels, closed
570 gaps) were assessed after each round of Pilon (see Supplementary Data 1). The rate of
571 improvements reached a plateau between the third (P3) and the fourth (P4) rounds of Pilon,
572 and therefore the resulting P4 polished assembly was considered optimal and used for
573 downstream analysis. Given the reference genome should not contain any homozygote
574 alternate variant calls relative to the short read data from the same sample, we compared how
575 the number of these changed following polishing. The Illumina short reads, sequenced from
576 the same animal as that used to generate the reference genome assembly (OPB4), were
577 mapped with bwa-mem (BWA v0.7.17) against the polished genome assemblies (P2 to P4).

578 The percentages of mapped reads were extremely high (>99%) and comparable across the P2,
579 P3 and P4 assemblies.

580 *Assembly statistics*

581 To directly compare the quality of the genome assembly at each step during the assembly
582 process, and to highlight improvements, QUAST (v 5.0.2) [65] was used to produce genome
583 assembly metrics for each iteration of the genome assembly, pre and post gap filling with
584 PBJelly, and for each successive round of polishing with Pilon (Supplementary Data 1).
585 QUAST further compares a given genome assembly to a reference genome, and for this the
586 genome assembly for the water buffalo *Bubalus bubalis* (GCF_003121395.1) [29] was
587 provided, to produce genome alignment metrics and details of suspected misassemblies
588 (Supplementary Data 1). A custom Python script
589 (https://raw.githubusercontent.com/evotools/CattleGraphGenomePaper/master/Assembly/AB
590 S.py) was used to calculate scaffold metrics, N, L, NG, LG and GC content for a given
591 proportion of the scaffold-level P4 genome assembly, in 5% increments (5-100,
592 Supplementary Data 1). The scaffold-level P4 genome assembly contains a total of 3,351
593 scaffolds, of which 1,381 scaffolds are greater than 10 kb. Quality values (QV) representative
594 of the single-base accuracy were computed using Merqury (v1.1)[66]with the K-mer counts
595 generated by Meryl (v1.2; https://github.com/marbl/meryl). For downstream analysis we
596 selected 1,381 contigs with a length of 10 kb or greater, representing 99.68% (2.653 Gb) of
597 the total length of the assembled genome. This subset of contigs were used for downstream
598 analyses.

599 *Detection of novel genomic sequences*

600 Following completion of the assembly, we identified the novel sequences in the genome in
601 comparison with other ruminant species. We selected a set of nine genome assemblies for
602 five species, and calculated the distances among them using mash v2.2 [67], using a K-mer
603 size of 32. We used the following genome assemblies to generate the alignment graph:
604 *Syncerus caffer* (accession number GCA_902825105.1), *Bubalus bubalis* Mediterranean
605 (GCF_003121395.1)[29], *Capra hircus* San Clemente (GCF_001704415.1)[34], *Bos
606 grunniens* (GCA_005887515.2)[30], *Bos taurus indicus* Brahman (GCF_003369695.1), *Bos
607 taurus taurus* Angus (GCA_003369685.2)[31], *Bos taurus taurus* Hereford
608 (GCF_002263795.1)[33], *Bos taurus taurus* N'Dama (GCA_905123515) and *Bos taurus
609 indicus* Ankole (GCA_905123885)[32]. We then generated a phylogenetic tree using the
610 neighbour-joining algorithm included in the neighbour software from Phylip (v3.698)[68]
611 which was used to create the following guide tree for CACTUS [69]:

612 ((angus:0.00187,hereford:0.00115)Anc1:0.0004,(ankole:0.00317,((yak:0.00671,((abuffalo:0.
613 01228,wbuffalo:0.0095)Anc6:0.00438,goat:0.04443)Anc5:0.01195)Anc4:0.00254,brahman:0
614 .00256)Anc3:0.00023)Anc2:0.0004,ndama:0.00195)Anc0;

615 The HAL archive of multiple whole genome alignments (mWGA) was generated using the
616 software CACTUS [69], and then converted to PackedGraph format using the hal2vg
617 software (v.2.1)[70] with the African buffalo genome as reference. We then used the nf-
618 GraphSeq workflow
619 (https://github.com/evotools/CattleGraphGenomePaper/tree/master/detectSequences/nf-
620 GraphSeq) described in Talenti et al. [32] based on libbdsg [71] to identify the nodes (i.e. the
621 fragment of genome) that are found exclusively in the backbone of the graph (i.e. African
622 buffalo genome), excluding all intervals overlapping a gap. We combined all interval regions
623 less than 5bp apart using BEDTools (v.2.30.0) [72]. We then annotated the regions by length
624 (short if < 10bp, intermediate if < 60bp and large if > 60bp), position (labelled telomeric if

13

625 &lt;10Kb from the end of a scaffold larger than 5Mb, flanking a gap if &lt;1Kb from an N-mer),
626 type of sequence (novel if > 95% of the bases in the region are not found in any other
627 genome, or haplotype if < 95% of the bases were found only in the African Buffalo) and
628 proportion of masked bases. We filtered out regions if they 1) were not classified as long, 2)
629 contained less than 50% novel bases and 3) were not telomeric or were not flanking a gap.

630 To validate that these regions corresponded to buffalo sequence, and did not derive, for
631 example, from contamination, 46 of the population WGS samples were randomly selected
632 and their coverage examined at these regions, with the assumption that if these regions
633 corresponded to contamination in our reference sample, they would not have aligned reads
634 from multiple buffalo samples. Mean read depth was calculated for each of the 74,659 novel
635 regions within the reference genome, for the 46 population samples, using Mosdepth (v0.3.4)
636 [73]. The distribution of average coverage values across the population samples, for each
637 novel region, is shown in Supplementary Figure 1. There are only 1494 novel regions with a
638 mean read depth $< 1$ and 419 regions with no reads mapped across these 46 samples,
639 suggesting that these putative African buffalo-specific regions do not derive from an artefact
640 such as contamination.

641 We characterized the content of the novel regions by 1) performing a motif analysis using
642 HOMER (v4.11.1)[74], and 2) by detecting the novel features. To identify these features, we
643 used the annotation generated by Ensembl and available in the rapid release database
644 (http://www.ensembl.info/2020/06/25/ensembl-rapid-release/; accession
645 GCA_902825105.1). We identified all gene features overlapping a novel sequence using
646 bedtools intersect (v2.30.0) [72], and identified only these fully overlapping a novel region
647 still with bedtools intersect with the -f 1.0 option (100% of overlap between the feature and
648 the novel region).

649 Once we identified these fully new gene features, we extracted the GO term and KEGG
650 pathways present in the annotation itself in embl format. To do so, we first converted the file
651 in GenBank format, and then extracted for each gene the transcript IDs, protein IDs and
652 biological terms. For these terms, we performed an enrichment analysis in R using a binomial
653 test with the genes not in novel regions as background.

654 *Reference genome annotation*

655 Genome annotation was undertaken at EMBL-EBI by Ensembl, primarily using RNA-seq
656 and full-length isoform sequencing (Iso-Seq) data generated from the animal for which the
657 genome was assembled. A TruSeq stranded total RNA-seq library with one round of Ribo-
658 Zero Gold kit (Illumina) was prepared from one pooled library consisting of RNA samples
659 from eight tissues (heart, prescapular and inguinal lymph nodes, testis, liver, kidney, lung and
660 spleen) collected from the animal for which the genome was assembled. RNA-seq was
661 performed at Edinburgh Genomics on an S2 lane of an Illumina NovaSeq 6000 platform
662 generating 100bp paired-end reads. Iso-Seq was performed at the Centre for Genomic
663 Research at the Univerity of Liverpool, using RNA samples from six different tissues
664 (prescapular lymph node, testis, liver, kidney, lung and spleen) collected from the same
665 animal. Full-length cDNA from total RNA was generated using TeloPrime full-length cDNA
666 amplification kit (v2) from Lexogen. A total of six barcoded TeloPrime libraries from six
667 RNA samples were multiplexed. Iso-seq was performed on the resulting multiplexed library
668 using six PacBio Sequel SMRT cells. The RNA-seq data were aligned to the reference
669 genome using STAR [75]. For loci where the structures derived from the transcriptomic data
670 appeared to be fragmented or absent, gap-filling using cross-species protein data was carried
671 out. For more information on the annotation process see Supplementary Information 1.

14

672 *Detection of variants in WGS samples across Africa.*

673 For all 196 WGS samples from *S. caffer* across Africa (raw data is available at ENA via
674 accession numbers PRJEB59220 and ERP144275), reads were mapped with bwa-mem
675 (BWA v0.7.17) against the reference genome generated as above. The GATK (v4.0.11.0)
676 pipeline, following the best practices as outlined at https://gatk.broadinstitute.org/hc/en-
677 us/articles/360036194592-Getting-started-with-GATK4 was used with HaplotypeCaller to
678 identify variants (SNPs and Indels). The GATK best practice includes a Variant Quality
679 Score Recalibration (VQSR) step that compares all variant calls to those in a high quality set
680 to identify and flag potential false positives. Unlike in well-characterised species no gold-
681 standard set of variants is available for the African buffalo. We therefore used a consensus set
682 of 6,806,905 variants called from the Illumina data generated for the same sample as the
683 reference genome using three software tools (GATK, Arrow and Longshot [76]). Although
684 we do not expect this set to be free of false variant calls, we expect it to be enriched for true
685 positives and was therefore used in VQSR. Three VQSR tranches, 99, 99.9 and 100 (each
686 representing the proportion of gold-standard variants that are retained at each quality
687 threshold), were assessed. The variant set resulting from the 99.9 tranche was selected for
688 downstream analyses with a Ti/Tv ratio of 2.07 and >120M variants. The variant set was
689 further filtered for GQ (Phred-scaled Probability that the call is incorrect) values less than 30
690 and site missingness of 0.9 (at least 90% of the samples contain data at this site). PLINK
691 (v1.90) was used to calculate sample missingness, the proportion of variant sites missing
692 from each sample, and vcftools (v0.1.13) to calculate the relatedness of all individuals. For
693 downstream analyses, individual samples with a missingness greater than 0.15, additionally
694 individuals that were closer than fourth degree relatedness (relatedness value 0.0625), were
695 also removed, resulting in a variant dataset covering 163 individual animals. We checked for
696 any mapping biases due to use of an East African reference genome, by randomly sampling
697 three animals per country and comparing how read mapping rates differed by longitude
698 (Supplementary Figure 7). No obvious mapping bias was observed among the West African
699 samples when mapping to the reference genome obtained from an East African sample.

700 *Genomic diversity analyses*

701 The VCF file for the set of unrelated samples was first filtered through bcftools (v1.9;
702 https://samtools.github.io/bcftools/) to keep only unrelated individuals according to the KING
703 method implemented in vcftools [77, 78]. A cutoff of 0.0625 was applied to exclude 3$^{rd}$
704 degree relatives or closer. Furthermore only biallelic SNPs in large contigs (>10Kb) were
705 retained. Variants were further filtered using plink (v1.90b4) [79] to restrict to those with a
706 minor allele frequency >0.05. This dataset was then used to carry out analyses of migration
707 events and effective population size. ADMIXTURE and the identity-by-state phylogenetic
708 tree can benefit from having an even sample size for the different populations/samples
709 deriving from the same location that were tested [38]. Therefore, for these analyses we
710 identified a representative subsample for the populations with more than 15 animals. Sample
711 size reduction was carried out using the BITE R package [80] to select a representative set of
712 individuals for each population. The reduction process was performed on each population
713 separately. For each group we selected the variants with very high call rate (99%) and highly
714 polymorphic (--maf 0.3). The reduction step in BITE was performed considering only
715 individuals with 95% call rate and up to 10K markers to compute the kinship matrix (options
716 n.trials = 100000, ibs.marker=10000, n.k=2, ibs.thr = 0.95, id.cr=0.95). Principal component
717 analysis (PCA) was performed post reduction in sample numbers using plink v1.90b4.
718 Admixture analysis was performed using ADMIXBoots
719 (https://github.com/RenzoTale88/ADMIXBoots), a Nextflow workflow that performs
720 bootstrapped admixture (v1.3.0) [81], defining a consensus of the different K at different

15

721  iterations using CLUMPP [82] and generating plots in R. The workflow was run pruning for
722  variants in linkage (plink --indep-pairwise 5000 100 0.3), testing every K between 2 and 15,
723  and with 100 bootstraps of 100,000 markers each. A consensus of the different bootstraps
724  was called using CLUMPP in LargeKGreedy mode. Bar charts for each consensus K,
725  boxplots for the distribution of the CV errors and line plots of the H' scores of each K were
726  generated from the pipeline automatically. Bootstrapped identity-by-state (IBS) phylogenetic
727  tree was calculated using the nf-PhyloTree workflow (https://github.com/RenzoTale88/nf-
728  PhyloTree). This workflow uses plink to generate a matrix of identity-by-state distances
729  across individuals. These workflows then use a series of custom scripts to generate the
730  individual phylogenetic tree, call the consensus tree and generate the input compliant for
731  GraphLan [83]. In our case, we ran the workflow allowing for pruning variants in high
732  linkage disequilibrium using plink (--indep-pairwise 5000 100 0.3), then generated  100-
733  bootstrapped IBS-based distances using plink (--distance 1-ibs square flat-missing), allowing
734  repeated variants in each bootstrap and sampling a number of SNPs equalling the number of
735  pruned variants. A consensus tree was generated using a custom python script, converted to
736  phyloXML, annotated for colour and consistency nodes and finally plotted with GraphLan
737  [83]. For the isolation by distance analysis, pairwise $F_{ST}$ values between populations were
738  calculated using vcftools, and the Haversine formula was used to calculate the distances
739  between the centre points of population sampling sites.

740  *Estimated Effective Migration Surfaces (EEMS)*

741  The EEMS package developed by Petkova et al. [84] was used
742  (https://github.com/dipetkov/eems) to estimate effective migration surfaces. The
743  runeems_snps program was used to visualise spatial population structure in the African
744  buffalo populations and to identify the geographic barriers to migration preventing gene flow
745  across these populations. The runeems_snps program requires the following data as input
746  files: (1) a matrix of average pairwise genetic dissimilarities, (2) sample coordinates, and (3)
747  a list of habitat coordinates, here covering the natural distribution of African buffalo
748  populations on the African continent, and listed as a sequence of vertices organised as a
749  closed polygon. For the input files for EEMS analysis, a matrix of average pairwise genetic
750  dissimilarities was generated from the pruned set of SNP data, using the bed2diffs_v1
751  program within the EEMS package. The locations of all African buffalo animals, from which
752  DNA samples were collected for WGS and variant detection, were inputted as longitude and
753  latitude coordinates, indicating either specific sampling locations or the centre of specific
754  geographical regions (*e.g.* national parks) when no other information was available. The list
755  of habitat coordinates was generated based on the known past and present natural distribution
756  of the four subspecies of African buffalo populations (as described in [2]) and using the
757  https://www.latlong.net/ website to identify the latitude and longitude geocoding of point
758  locations on the African continent. EEMS analysis was run using the runeems_snps program
759  within the EEMS package based on the African buffalo pruned SNP data. Parameters used to
760  run EEMS analysis were set as follows: nIndiv = 163; nSites = 6000; nDemes = 400; diploid
761  = true; numMCMCIter = 4000000; numBurnIter = 1000000; numThinIter = 9999.
762  Description for all parameters used are defined in the EEMS instruction manual
763  (v.0.0.0.9000). Results of EEMS analysis were plotted using the rEEMSplot package in R to
764  generate contour plots of effective migration and effective diversity surfaces from EEMS
765  outputs. Additionally, posterior probability trace plots (pilogl) were used to check the MCMC
766  sampler had successfully converged using four million MCMC iterations. The effective
767  migration and diversity surfaces plots also include the addition of lakes and rivers depicted in
768  blue based on data extracted from the Natural Earth website
769  (https://www.naturalearthdata.com/download/50m/physical/).

770 *Estimating effective population sizes and selective sweeps*

771 To calculate the XP-EHH scores the African buffalo genotype data was first phased using
772 Beagle 5.1 [85]. A recombination rate of 1cM/Mb was assumed and XP-EHH scores
773 calculated between each pair of populations using hapbin [86]. Peaks were called as
774 previously described [42]. Briefly XP-EHH scores were smoothed by averaging across 1000
775 SNP windows and putative selective sweep regions were those with an absolute XP-EHH >4,
776 with the start and end coordinates defined where the XP-EHH scores fell back below two.
777 The locations of XP-EHH peaks in the water buffalo and cattle genomes were obtained from
778 Dutta et al. [42] and the peaks for all three species mapped to the orthologous regions of the
779 water buffalo genome.
780 The effective population sizes over time of the three largest African buffalo populations were
781 calculated using Relate v1.1.6 [40] using the same phased haplotypes from Beagle. An
782 estimated generation time of 11 years for the African buffalo was used in this analysis [87].
783 Previously calculated estimated effective population sizes for human African populations
784 were obtained from Speidel et al [40]. The $P_R$ statistic was also calculated using Relate [40]
785 and the same Beagle haplotype files using an estimated mutation rate of $1.25 \times 10^{-8}$. Variants
786 with a P less than $5 \times 10^{-8}$ were retained. The circular Manhattan plot was created using the
787 CMplot R package [88]. The water buffalo genes were lifted over to the African buffalo
788 genome to identify which genes fell under putative selective sweep peaks.

789 *Data availability statement*

790 All raw data generated in this project is available through the following routes: for the buffalo
791 reference genome, PacBio, Illumina and Hi-C data is available via Genbank
792 (GCA_902825105.1) and ENA (PRJEB3658), and the assembly, annotation and associated
793 flat files can be accessed through
794 https://rapid.ensembl.org/Syncerus_caffer_GCA_902825105.1; transcriptomic data were
795 deposited to ENA (https://www.ebi.ac.uk/ena/browser/view/GCA_902825105.1) with
796 accession numbers PRJEB36587 and PRJEB36588 for RNA-seq and Iso-Seq, respectively,
797 and population genome data is available through ENA via accession number PRJEB59220.

798 **Acknowledgements**

816
817

818 **Figure Legends**

819

820 Figure 1. Genome assembly metrics. The BlobToolKit Snailplot shows N50 metrics and
821 BUSCO gene completeness. A) *Syncerus caffer de novo* assembly. The main plot represents
822 the full genome length of 2.65 Gb. The distribution of scaffold length is shown in dark grey
823 with the plot radius scaled to the longest chromosome present in the assembly (190 Mb,
824 shown in red). Dark and light orange sections represent N50 and N90 (69Mb and 8.8Mb),
825 respectively. The pale grey spiral shows the cumulative scaffold count on a log scale with
826 white scale lines showing successive orders of magnitude. The blue/pale-blue/white ring
827 graph shows the distribution of GC, AT and N percentages for the given range in the main
828 plot. A summary of complete, fragmented, duplicated and missing BUSCO genes in the
829 mammalia_odb9 set is shown in the top right. B) Chen et al published genome. BlobToolKit
830 Snailplot representing the *S. caffer* assembly as presented by Chen et al. [26] C) Glanzmann
831 et al published genome. BlobToolKit Snailplot representing the *S. caffer* assembly as
832 presented by Glanzmann et al. [25] The plot radius for both the Chen and Glanzmann
833 genomes has been scaled to the maximum contig length (190Mb) in the *S. caffer* genome
834 assembled here to enable comparison of metrics.

835 Figure 2. (A) The sampling locations of African buffalo samples sequenced in the current
836 study (circled letters), mapped on to the approximate current distribution of the four
837 subspecies. a: Singou and Pama Game Reserves (GR)/Arli National Park (NP) complex,
838 Burkina Faso (n=10 samples [before data filtering]); b: W NP, Niger (n=10); c: Zakouma NP,
839 Chad (n=13), d: Manovo-Gounda-St. Floris NP, Central African Republic (CAR; n=2); e:
840 Bamingui-Bangoran NP, CAR (n=2); f: Sangba, CAR (n=1); g: N'Gotto Forest Reserve,
841 CAR (n=2); h: Lekedi NP, Gabon (n=8); i: Murchison Falls NP, Uganda (n=13), j: Kidepo
842 NP, Uganda (n=20); k: Ol Pejeta Game Reserve, Kenya (n=12); l: Serengeti NP, Tanzania
843 (n=15); m: Ngorongoro Conservation Area, Tanzania (n=15); n: Tarangire NP, Tanzania
844 (n=10); o: Arusha NP, Tanzania (n=10); p: Niassa National Reserve (NR), Mozambique
845 (n=9); q: Marromeu NR, Mozambique (n=9); r: Chobe NP, Botswana (n=9); s: Okavango
846 Delta, Botswana (n=9); t: Gonarezhou NP/Crook's Corner, Zimbabwe (n=18), u: Hluhluwe-
847 Umfolozi NP, South Africa (n=8; for full sample data see Supplementary Table 2). (B)
848 Principal Component Analysis of population samples, with data for components 1 and 2
849 illustrated. Samples are coloured by country of origin, with different symbols indicating the
850 previously recognised subspecies.

851 Figure 3. Population genetic analyses based on genome sequences: A) Admixture analysis for
852 K = 3 showing the three major clusters of diversity (Western/Central Africa in red, Eastern
853 Africa in yellow and Southern Africa in blue). B) Neighbour-joining phylogenetic tree from
854 the Identity-by-State (IBS) distances showing the gradient from Southern- to Eastern- to
855 Western/Central-African populations (clockwise from the root). Symbols indicate confidence
856 value for each node (circle is less than 50%; square between 50% and 74%; star between
857 75% and 89%; hexagon above 90%). C) Isolation by distance (IBD) analysis of African
858 buffalo populations. The $F_{ST}$ values were calculated between all pairs of populations and
859 plotted against their geographic distance apart. Pairwise comparisons involving *S. c. nanus*
860 are indicated in blue, pairwise comparisons involving the Hluhluwe-Umfolozi population are
861 shown in green, the single pairwise comparison comparing *S. c. nanus* and Hluhluwe-
862 Umfolozi in purple, and all remaining pairwise comparisons in red. The predicted pairwise
863 $F_{ST}$ values outside of the observed distances are indicated by dashed lines. D) The proportion
864 of homozygous segments per sample (FROH) indicating that the Hluhluwe-Umfolozi
865 population has unusually high levels of homozygosity.

866 Figure 4. (A) The contour map shows the mean of two independent Estimating Effective
867 Migration Surfaces (EEMS) posterior migration rate estimates between 400 demes modelled
868 over the land surface of Sub-Saharan Africa. A value of 1 (blue) indicates a tenfold greater
869 migration rate over the average; –1 (orange) indicates tenfold lower migration than average.
870 The courses of the major river systems (Niger, Congo, Nile and Orange rivers), as well as
871 water bodies with a surface area greater than 5,000 km$^2$ are included to highlight their
872 potential relationships with migratory routes and barriers. Red diamonds indicate
873 geographical location of samples in the dataset. (B) Estimated effective population sizes of
874 African buffalo (solid lines) and human (dashed lines) populations over time. The countries
875 of sampling for each population are indicated in the legend along with the three letter 1,000
876 Genomes consortium population code for the human data. Only human populations from the
877 1,000 Genomes consortium dataset of recent African origin are shown.

878 Figure 5. (A) The coloured outermost track and legend indicates the SNP density across 41
879 large contigs. The next three tracks show the $P_R$ scores in the Uganda (centremost),
880 Zimbabwe/Botswana (middle) and Tanzania/Kenya (outer) populations. Red points indicate
881 SNPs with a P-value smaller than $5 \times 10^{-8}$. The peaks at LOC102396916 are highlighted (B)
882 Absolute XP-EHH scores across Contig 187 for the Hluhluwe-Umfolozi versus intermediate
883 populations indicating the peak also detected at the LOC102396916 locus.

884 Supplementary Data 1 Genome assembly statistics.

885 Supplementary Figure 1 Mean read depth across the putative African buffalo specific regions.
886 Mean read depth was calculated for each of the 74,659 novel regions within the reference
887 genome, for 46 of the population samples, using Mosdepth (v0.3.4) [73]. The distribution of
888 average coverage values across the population samples, for each novel region, is shown.
889 There are only 1494 novel regions with a mean read depth < 1 and 419 regions with no reads
890 mapped across these 46 samples.

891 Supplementary Table 1. Genes identified in the buffalo-specific sequence, with Ensembl
892 transcript, gene and protein identifiers, and GO terms where relevant. Note that the list is
893 greater than the 583 identifed genes, as some genes appear in the list more than once due to
894 having different transcripts.

895 Supplementary Table 2. Details of buffalo samples for which genome sequences were
896 generated and included in this study, including sample identification, subspecies, country of
897 origin, region of origin, whether sequences were retained in analysis following filtering steps
898 (0.0625 relatedness, 0.2 missingness), the population group the sample was assigned to, and a
899 latitude/longitude of a central point in the respective sampling area.

900 Supplementary Figure 2 Principal Component analysis pre- & post-reduction (i.e. following
901 sample removal post  filtering steps; 0.0625 relatedness, 0.2 missingness), with data for
902 components 1 and 2 illustrated, samples are coloured by population grouping.

903 Supplementary Figure 3. Admixture evaluation metrics ((A) cross-validation error, (B)
904 number of iterations to converge and (C) H') at different values of K calculated using 100
905 bootstraps of 100,000 variants each.

906 Supplementary Figure 4. Admixture analysis for K = 2-10.

907 Supplementary Figure 5. Relate-inferred inverse coalescence rates (effective population
908 sizes) for each of the larger sub-groups to themselves (dashed lines) and each other (solid
909 lines). For this comparison, due to the smaller sample sizes, all West African animals were
910 collated into one group.

19

911 Supplementary Table 3. Pairwise $F_{ST}$ values for the nine population groupings (*S. c
912 brachyceros*, *S. c. nanus*, *S. c. aequinoctialis*, intermediate (putative hybrids between *S. c.
913 nanus*, *S. c. aequinoctialis*), *S. c. caffer* Uganda, *S. c. caffe*r Kenya/Tanzania, *S. c. caffer*
914 Mozambique, *S. c. caffer* Zimbabwe/Botswana and *S. c. caffer* South Africa), and geographic
915 distance as measured to centred latitude/longitude measurement for each grouping.

916 Supplementary Table 4. Details of genes identified to be under selection in the African
917 buffalo, whether the gene has been previously identified to be in a selection peak in either the
918 cow or water buffalo, and whether the gene is related to immune response function. Genes
919 are grouped by (a) detected in both XPEHH and PR analyses of African buffalo (dark green),
920 (b) detected in either XPEHH or PR analyses of African buffalo, and in both metrics for
921 water buffalo or cow analyses (medium green), (c) detected in either XPEHH or PR analyses
922 of African buffalo, and in one of the metrics for water buffalo or cow analyses (light green),
923 or (d) none of the above (no colour).

924 Supplementary Figure 6. The absolute XP-EHH scores at the LOC102396916 locus on contig
925 187. The boundaries of the called peak are indicated by dashed vertical lines. The location
926 and direction of LOC102396916 is shown in blue below.

927 Supplementary Information 1 A methodological summary of the genome annotation process
928 undertaken at ENSEMBL.

929 Supplementary Figure 7. Mapping rate by longitude of three randomly selected samples per
930 country. No obvious mapping bias was observed among the West African samples when
931 mapping to the reference genome obtained from an East African sample.

932 **References**
933
934 1.    East, R., *African Antelope Database 1999*, I.S.A.S. Group, Editor. 1999: Gland,
935       Switzerland and Cambridge, UK.
936 2.    Cornelis, D., et al., *African buffalo Syncerus caffer (Sparrman, 1779)*, in *Ecology,
937       evolution and behaviour of wild cattle: implications for conservation.*, M. Melletti
938       and J. Burton, Editors. 2014, Cambridge University Press: Cambridge, UK.
939 3.    Cornelis, D., et al., *Conservation status of the African buffalo: a continent-wide
940       assessment.*, in *Ecology and Management of the African buffalo*, A. Caron, et al.,
941       Editors. 2023, Cambridge Univeristy Press: Cambridge, UK.
942 4.    Michaux, J., N. Smitz, and P. Van Hooft, *Taxonomic status of the African buffalo.*, in
943       *Ecology and Management of the African buffalo*, A. Caron, et al., Editors. 2023,
944       Cambridge University Press: Cambridge, UK.
945 5.    Smitz, N., et al., *Pan-African genetic structure in the African buffalo (Syncerus
946       caffer): investigating intraspecific divergence.* PLoS One, 2013. **8**(2): p. e56235.
947 6.    Smitz, N., et al., *Genetic structure of fragmented southern populations of African
948       Cape buffalo (Syncerus caffer caffer).* BMC Evol Biol, 2014. **14**: p. 203.
949 7.    Heller, R., A. Bruniche-Olsen, and H.R. Siegismund, *Cape buffalo mitogenomics
950       reveals a Holocene shift in the African human-megafauna dynamics.* Mol Ecol, 2012.
951       **21**(16): p. 3947-59.
952 8.    Smitz, N., et al., *Genome-wide single nucleotide polymorphism (SNP) identification
953       and characterization in a non-model organism, the African buffalo (Syncerus caffer),
954       using next generation sequencing.* Mammalian Biology, 2016. **81**: p. 595-603.
955 9.    de Jager, D., et al., *High diversity, inbreeding and a dynamic Pleistocene
956       demographic history revealed by African buffalo genomes.* Sci Rep, 2021. **11**(1): p.
957       4540.

10. Quinn, L., et al., *Colonialism in South Africa leaves a lasting legacy of reduced genetic diversity in Cape buffalo.* Mol Ecol, 2023. **32**(8): p. 1860-1874.

11. Pizzutto, C.S., H. Colbachini, and P.N. Jorge-Neto, *One Conservation: the integrated view of biodiversity conservation.* Anim Reprod, 2021. **18**(2): p. e20210024.

12. Hohenlohe, P.A., W.C. Funk, and O.P. Rajora, *Population genomics for wildlife conservation and management.* Mol Ecol, 2021. **30**(1): p. 62-82.

13. Auty, H., et al., *Cattle trypanosomosis: the diversity of trypanosomes and implications for disease epidemiology and control.* Rev Sci Tech, 2015. **34**(2): p. 587-98.

14. Casey-Bryars, M., et al., *Waves of endemic foot-and-mouth disease in eastern Africa suggest feasibility of proactive vaccination approaches.* Nat Ecol Evol, 2018. **2**(9): p. 1449-1457.

15. Bengis, R., et al., *Infections and parasites of the free-ranging African buffalo.*, in *Ecology and Management of the African buffalo*, A. Caron, et al., Editors. 2023, Cambridge University Press: Cambridge, Uk.

16. Dwinger, R.H., et al., *Susceptibility of buffaloes, cattle and goats to infection with different stocks of Trypanosoma vivax transmitted by Glossina morsitans centralis.* Res Vet Sci, 1986. **41**(3): p. 307-15.

17. Grootenhuis, J.G., et al., *Susceptibility of African buffalo and Boran cattle to Trypanosoma congolense transmitted by Glossina morsitans centralis.* Vet Parasitol, 1990. **35**(3): p. 219-31.

18. Morrison, W.I., J.D. Hemmink, and P.G. Toye, *Theileria parva: a parasite of African buffalo, which has adapted to infect and undergo transmission in cattle.* Int J Parasitol, 2020. **50**(5): p. 403-412.

19. Gifford-Gonzalez, D., *Animal disease challenges to the emergence of pastoralism in Sub-Saharan Africa.* . African Archaeological Review, 2000. **17**(3): p. 95-139.

20. Lankester, F. and A. Davis, *Pastoralism and wildlife: historical and current perspectives in the East African rangelands of Kenya and Tanzania.* Rev Sci Tech, 2016. **35**(2): p. 473-484.

21. Michel, A.L., *Implications of tuberculosis in African wildlife and livestock.* Ann N Y Acad Sci, 2002. **969**: p. 251-5.

22. Kock, R., et al., *Livestock and buffalo (Syncerus caffer) interfaces in Africa: ecology of disease transmission and implications for conservation and development.*, in *Ecology, Evolution and Behaviour of Wild Cattle*, M. Melletti and J. Burton, Editors. 2014, Cambridge University Press: Cambridge, UK. p. 431-425.

23. Caron, A., et al., *Relationship between burden of infection in ungulate populations and wildlife/livestock interfaces.* Epidemiol Infect, 2013. **141**(7): p. 1522-35.

24. Kock, R., et al., *African buffalo and colonial cattle: is "systems change" the best future for farming and nature in Africa?*, in *Ecology and Management of the African buffalo*, A. Caron, et al., Editors. 2023, Cambridge University Press: Cambridge, UK.

25. Glanzmann, B., et al., *The complete genome sequence of the African buffalo (Syncerus caffer).* BMC Genomics, 2016. **17**(1): p. 1001.

26. Chen, L., et al., *Large-scale ruminant genome sequencing provides insights into their evolution and distinct traits.* Science, 2019. **364**(6446).

27. English, A.C., et al., *Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology.* PLoS One, 2012. **7**(11): p. e47768.

28. Walker, B.J., et al., *Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement.* PLoS One, 2014. **9**(11): p. e112963.

29. Low, W.Y., et al., *Chromosome-level assembly of the water buffalo genome surpasses human and goat genomes in sequence contiguity.* Nat Commun, 2019. **10**(1): p. 260.

30. Zhang, S., et al., *Structural Variants Selected during Yak Domestication Inferred from Long-Read Whole-Genome Sequencing.* Mol Biol Evol, 2021. **38**(9): p. 3676-3680.

31. Koren, S., et al., *De novo assembly of haplotype-resolved genomes with trio binning.* Nat Biotechnol, 2018.

32. Talenti, A., et al., *A cattle graph genome incorporating global breed diversity.* Nat Commun, 2022. **13**(1): p. 910.

33. Rosen, B.D., et al., *De novo assembly of the cattle reference genome with single-molecule sequencing.* Gigascience, 2020. **9**(3).

34. Bickhart, D.M., et al., *Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome.* Nat Genet, 2017. **49**(4): p. 643-650.

35. Girgis, H.Z., *Red: an intelligent, rapid, accurate tool for detecting repeats de-novo on the genomic scale.* BMC Bioinformatics, 2015. **16**: p. 227.

36. Franza, B.R., Jr., et al., *The Fos complex and Fos-related antigens recognize sequence elements that contain AP-1 binding sites.* Science, 1988. **239**(4844): p. 1150-3.

37. Klopfenstein, D.V., et al., *GOATOOLS: A Python library for Gene Ontology analyses.* Sci Rep, 2018. **8**(1): p. 10872.

38. Meirmans, P.G., *Subsampling reveals that unbalanced sampling affects STRUCTURE results in a multi-species dataset.* Heredity (Edinb), 2019. **122**(3): p. 276-287.

39. Sabeti, P.C., et al., *Genome-wide detection and characterization of positive selection in human populations.* Nature, 2007. **449**(7164): p. 913-8.

40. Speidel, L., et al., *A method for genome-wide genealogy estimation for thousands of samples.* Nat Genet, 2019. **51**(9): p. 1321-1329.

41. Moon, J.M., et al., *Examination of Signatures of Recent Positive Selection on Genes Involved in Human Sialic Acid Biology.* G3 (Bethesda), 2018. **8**(4): p. 1315-1325.

42. Dutta, P., et al., *Whole genome analysis of water buffalo and global cattle breeds highlights convergent signatures of domestication.* Nat Commun, 2020. **11**(1): p. 4739.

43. Prajapati, B.M., et al., *Molecular markers for resistance against infectious diseases of economic importance.* Vet World, 2017. **10**(1): p. 112-120.

44. Young, R., et al., *A Gene Expression Atlas of the Domestic Water Buffalo (Bubalus bubalis).* Front Genet, 2019. **10**: p. 668.

45. Fennessy, J., et al., *Multi-locus Analyses Reveal Four Giraffe Species Instead of One.* Curr Biol, 2016. **26**(18): p. 2543-9.

46. Pedersen, C.T., et al., *A southern African origin and cryptic structure in the highly mobile plains zebra.* Nat Ecol Evol, 2018. **2**(3): p. 491-498.

47. Lohay, G.G., et al., *Genetic connectivity and population structure of African savanna elephants (Loxodonta africana) in Tanzania.* Ecol Evol, 2020. **10**(20): p. 11069-11089.

48. Bertola, L.D., et al., *Phylogeographic Patterns in Africa and High Resolution Delineation of Genetic Clades in the Lion (Panthera leo).* Sci Rep, 2016. **6**: p. 30807.

49. Smitz, N., et al., *A genome-wide data assessment of the African lion (Panthera leo) population genetic structure and diversity in Tanzania.* PLoS One, 2018. **13**(11): p. e0205395.

50. Genomes Project, C., et al., *A global reference for human genetic variation.* Nature, 2015. **526**(7571): p. 68-74.

51. Coimbra, R.T.F., et al., *Conservation Genomics of Two Threatened Subspecies of Northern Giraffe: The West African and the Kordofan Giraffe.* Genes (Basel), 2022. **13**(2).

52. Van Hooft, W.F., A.F. Groen, and H.H. Prins, *Phylogeography of the African buffalo based on mitochondrial and Y-chromosomal loci: Pleistocene origin and population expansion of the Cape buffalo subspecies.* Mol Ecol, 2002. **11**(2): p. 267-79.

53. Heller, R., et al., *Mid-Holocene decline in African buffalos inferred from Bayesian coalescent-based analyses of microsatellites and mitochondrial DNA.* Mol Ecol, 2008. **17**(22): p. 4845-58.

54. Mack, R., *The great African cattle plague epidemic of the 1890's.* Trop. Anim. Hlth. Prod., 1970. **2**: p. 210-219.

55. Plowright, W., *The effects of rinderpest and rinderpest control on wildlife in Africa. .* Symposia of the Zoological Society of London, 1982. **50**: p. 1-28.

56. Estes, R.D., *The Behaviour Guide to African Mammals.* 1991, Berkeley, USA: University of California Press.

57. Van Hooft, W.F., A.F. Groen, and H.H. Prins, *Microsatellite analysis of genetic diversity in African buffalo (Syncerus caffer) populations throughout Africa.* Mol Ecol, 2000. **9**(12): p. 2017-25.

58. Simonsen, B.T., H.R. Siegismund, and P. Arctander, *Population structure of African buffalo inferred from mtDNA sequences and microsatellite loci: high variation but low differentiation.* Mol Ecol, 1998. **7**(2): p. 225-37.

59. Stephens, S.A. and C.J. Howard, *Infection and transformation of dendritic cells from bovine afferent lymph by Theileria annulata.* Parasitology, 2002. **124**(Pt 5): p. 485-93.

60. Glass, E.J., S. Crutchley, and K. Jensen, *Living with the enemy or uninvited guests: functional genomics approaches to investigating host resistance or tolerance traits to a protozoan parasite, Theileria annulata, in cattle.* Vet Immunol Immunopathol, 2012. **148**(1-2): p. 178-89.

61. Bishop, R.P., et al., *The African buffalo parasite Theileria sp. (buffalo) can infect and immortalize cattle leukocytes and encodes divergent orthologues of Theileria parva antigen genes.* Int J Parasitol Parasites Wildl, 2015. **4**(3): p. 333-42.

62. Wragg, D., et al., *A locus conferring tolerance to Theileria infection in African cattle.* PLoS Genet, 2022. **18**(4): p. e1010099.

63. Decker, J.E., et al., *Worldwide patterns of ancestry, divergence, and admixture in domesticated cattle.* PLoS Genet, 2014. **10**(3): p. e1004254.

64. Obara, I., et al., *The Rhipicephalus appendiculatus tick vector of Theileria parva is absent from cape buffalo (Syncerus caffer) populations and associated ecosystems in northern Uganda.* Parasitol Res, 2020. **119**(7): p. 2363-2367.

65. Gurevich, A., et al., *QUAST: quality assessment tool for genome assemblies.* Bioinformatics, 2013. **29**(8): p. 1072-5.

66. Rhie, A., et al., *Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies.* Genome Biol, 2020. **21**(1): p. 245.

67. Ondov, B.D., et al., *Mash: fast genome and metagenome distance estimation using MinHash.* Genome Biol, 2016. **17**(1): p. 132.

68. Felsenstein, J., *PHYLIP (Phylogeny Inference Package)* 2009: Department of Genome Sciences, University of Washington, Seattle.

69. Armstrong, J., et al., *Progressive Cactus is a multiple-genome aligner for the thousand-genome era.* Nature, 2020. **587**(7833): p. 246-251.

70. Hickey, G., et al., *HAL: a hierarchical format for storing and analyzing multiple genome alignments.* Bioinformatics, 2013. **29**(10): p. 1341-2.

71. Eizenga, J.M., et al., *Efficient dynamic variation graphs.* Bioinformatics, 2020. **36**(21): p. 5139-5144.

72. Quinlan, A.R. and I.M. Hall, *BEDTools: a flexible suite of utilities for comparing genomic features.* Bioinformatics, 2010. **26**(6): p. 841-2.

1108    73.    Pedersen, B.S. and A.R. Quinlan, *Mosdepth: quick coverage calculation for genomes*
1109           *and exomes.* Bioinformatics, 2018. **34**(5): p. 867-868.
1110    74.    Heinz, S., et al., *Simple combinations of lineage-determining transcription factors*
1111           *prime cis-regulatory elements required for macrophage and B cell identities.* Mol
1112           Cell, 2010. **38**(4): p. 576-89.
1113    75.    Dobin, A., et al., *STAR: ultrafast universal RNA-seq aligner.* Bioinformatics, 2013.
1114           **29**(1): p. 15-21.
1115    76.    Edge, P. and V. Bansal, *Longshot enables accurate variant calling in diploid genomes*
1116           *from single-molecule long read sequencing.* Nat Commun, 2019. **10**(1): p. 4660.
1117    77.    Danecek, P., et al., *The variant call format and VCFtools.* Bioinformatics, 2011.
1118           **27**(15): p. 2156-8.
1119    78.    Manichaikul, A., et al., *Robust relationship inference in genome-wide association*
1120           *studies.* Bioinformatics, 2010. **26**(22): p. 2867-73.
1121    79.    Chang, C.C., et al., *Second-generation PLINK: rising to the challenge of larger and*
1122           *richer datasets.* Gigascience, 2015. **4**: p. 7.
1123    80.    Milanesi, M., et al., *BITE: an R package for biodiversity analyses.* BioRxiv, 2017.
1124    81.    Alexander, D.H., J. Novembre, and K. Lange, *Fast model-based estimation of*
1125           *ancestry in unrelated individuals.* Genome Res, 2009. **19**(9): p. 1655-64.
1126    82.    Jakobsson, M. and N.A. Rosenberg, *CLUMPP: a cluster matching and permutation*
1127           *program for dealing with label switching and multimodality in analysis of population*
1128           *structure.* Bioinformatics, 2007. **23**(14): p. 1801-6.
1129    83.    Asnicar, F., et al., *Compact graphical representation of phylogenetic data and*
1130           *metadata with GraPhlAn.* PeerJ, 2015. **3**: p. e1029.
1131    84.    Petkova, D., J. Novembre, and M. Stephens, *Visualizing spatial population structure*
1132           *with estimated effective migration surfaces.* Nat Genet, 2016. **48**(1): p. 94-100.
1133    85.    Browning, S.R. and B.L. Browning, *Rapid and accurate haplotype phasing and*
1134           *missing-data inference for whole-genome association studies by use of localized*
1135           *haplotype clustering.* Am J Hum Genet, 2007. **81**(5): p. 1084-97.
1136    86.    Maclean, C.A., N.P. Chue Hong, and J.G. Prendergast, *hapbin: An Efficient Program*
1137           *for Performing Haplotype-Based Scans for Positive Selection in Large Genomic*
1138           *Datasets.* Mol Biol Evol, 2015. **32**(11): p. 3027-9.
1139    87.    Pacifici, M., et al., *Database on generation length of mammals.* Nature Conservation,
1140           2013. **5**: p. 89-94.
1141    88.    Yin, L., *Package "CMplot".* 2019.
1142
1143

1144    Table 1. Sample number by country, subspecies and pre- and post-data filtering.

| Sample origin | Unfiltered | Filtered Missingness 0.20; Relatedness 0.0625 |
|---|---|---|
| Botswana | 17 | 15 |
| Burkina Faso | 9 | 7 |
| Central African Republic | 6 | 6 |
| Chad | 12 | 9 |
| Gabon | 7 | 2 |
| Kenya | 12 | 11 |
| Mozambique | 20 | 11 |
| Niger | 10 | 9 |
| South Africa | 8 | 6 |
| Tanzania | 50 | 48 |
| Uganda | 30 | 27 |
| Zimbabwe | 15 | 12 |
| Total | 196 | 163 |
| By subspecies | | |
| *S. c. caffer* | 152 | 130 |
| *S. c. brachyceros* | 19 | 16 |
| *S. c. aequinoctialis (S.c.a)* | 12 | 9 |
| Putative intermediate (*S.c.n/S.c.a*) | 6 | 6 |
| *S. c. nanus (S. c. n)* | 7 | 2 |
| Total | 196 | 163 |

1145
1146
1147
1148

**A**

**Scaffold statistics**
- Log10 scaffold count (total 3.4k)
- Scaffold length (total 2.7G)
- Longest scaffold (190M)
- N50 length (69M)
- N90 length (8.8M)

**BUSCO** mammalia_odb9 (4104)
- Complete (91.5%)
- Fragmented (5.0%)
- Duplicated (1.0%)
- Missing (3.5%)

**Scale**
- 2.7G
- 190M

Dataset: Syncerus caffer

**Composition**
- GC (41.8%)
- AT (58.2%)
- N (0.1%)

**B**

**Scaffold statistics**
- Log10 scaffold count (total 150k)
- Scaffold length (total 2.9G)
- Longest scaffold (11M)
- N50 length (2.3M)
- N90 length (270k)

**Scale**
- 2.9G
- 190M

Dataset: Chen Genome

**Composition**
- GC (41.8%)
- AT (58.2%)
- N (2.1%)

**C**

**Scaffold statistics**
- Log10 scaffold count (total 440k)
- Scaffold length (total 2.7G)
- Longest scaffold (17M)
- N50 length (2.4M)
- N90 length (420k)

**Scale**
- 2.7G
- 190M

Dataset: Glanzmann Genome

**Composition**
- GC (41.8%)
- AT (58.2%)
- N (2.8%)

A



B

**A**

**B**

log(m)

YRI
MSL
ASW
GWD
ESN
LWK
ACB

ZB
TK
UG

Estimated effective population size (Ne)

Years ago

Species — Buffalo
‑ ‑ ‑ Human

Population

Uganda (UG)
Tanzania/Kenya (TK)
Zimbabwe/Botswana (ZB)
Nigeria (YRI)
Kenya (LWK)

The Gambia (GWD)
Sierra Leone (MSL)
Nigeria (ESN)
USA (ASW)
Barbados (ACB)

**A**

**B**

PCA components 1 and 2 (full dataset)

PCA components 1 and 2 (representative subsample)

**A**



CV error − 100 bootstrap

**B**



Iterations to converge − 100 bootstrap

**C**



H' − 100 bootstrap

LOC102396916