

1 Soft sweeps predominate recent positive selection in bonobos (*Pan paniscus*) and chimpanzees
2 (*Pan troglodytes*)

3

4 Colin M. Brand¹, Frances J. White¹, Nelson Ting^{1,2}, Timothy H. Webster³

5

6 ¹Department of Anthropology, University of Oregon, Eugene, OR

7 ²Institute of Ecology and Evolution, University of Oregon, Eugene, OR

8 ³Department of Anthropology, University of Utah, Salt Lake City, UT

9

10 Corresponding Author: Colin M. Brand

11

12 University of Oregon

13 Department of Anthropology

14 1218 University of Oregon

15 Eugene, OR 97403

16 cbrand2@uoregon.edu

17 **Abstract**

18 Two modes of positive selection have been recognized: 1) hard sweeps that result in the
19 rapid fixation of a beneficial allele typically from a *de novo* mutation and 2) soft sweeps that are
20 characterized by intermediate frequencies of at least two haplotypes that stem from standing
21 genetic variation or recurrent *de novo* mutations. While many populations exhibit both hard and
22 soft sweeps throughout the genome, there is increasing evidence that soft sweeps, rather than
23 hard sweeps, are the predominant mode of adaptation in many species, including humans. Here,
24 we use a supervised machine learning approach to assess the extent of hard and soft sweeps in
25 the closest living relatives of humans: bonobos and chimpanzees (genus *Pan*). We trained
26 convolutional neural network classifiers using simulated data and applied these classifiers to
27 population genomic data for 71 individuals representing all five extant *Pan* lineages, of which
28 we successfully analyzed 60 individuals from four lineages. We found that recent adaptation in
29 *Pan* is largely the result of soft sweeps, ranging from 73.1 to 97.7% of all identified sweeps.
30 While few hard sweeps were shared among lineages, we found that between 19 and 267 soft
31 sweep windows were shared by at least two lineages. We also identify novel candidate genes
32 subject to recent positive selection. This study emphasizes the importance of shifts in the
33 physical and social environment, rather than novel mutation, in shaping recent adaptations in
34 bonobos and chimpanzees.

35

36 **Keywords:** adaptation, convolutional neural network, diploS/HIC, selective sweep, supervised
37 machine learning

38 **Introduction**

39 The identification of adaptative traits and their genetic basis is one of the central goals of
40 evolutionary biology. Two approaches, top-down and bottom-up, have been used to accomplish
41 this goal; the latter of which leverages population-level data to recognize the genomic signatures
42 of positive selection (Barrett and Hoekstra 2011). At the genomic level, the process of adaptation
43 results in a window of reduced variation that erodes over time. As these signatures do not persist,
44 they can only be used to infer selection over a particular time scale in a population. In most
45 species, this time frame is restricted to a few thousand generations, roughly ~ 200,000 years in
46 humans (Oleksyk et al. 2010). The classic model for positive selection for a given locus proposes
47 that a single, novel mutation, that confers a fitness advantage (i.e., a beneficial allele) will rapidly
48 spread in a population and eventually reach fixation (Maynard Smith and Haigh 1974). Neutral
49 polymorphism adjacent to the novel allele will ‘hitchhike’, resulting in a distinct pattern of
50 reduced genomic diversity at the locus and surrounding sites. The term ‘hard sweep’ has been
51 used to identify this pattern and process.

52 ‘Soft sweeps’ describe the presence of two or more haplotypes that occur at intermediate
53 frequencies (Hermisson and Pennings 2005). Thus, the signature of a soft sweep is intermediate
54 to those of neutral or ‘background’ genomic variation and the signature of a hard sweep. This
55 pattern can result from recurrent *de novo* mutations following positive selection. Alternatively,
56 soft sweeps can also result from positive selection on standing genetic variation where alleles
57 were already present in a population before selection. This variation may be the result of
58 independent mutations (multiple origin soft sweep) or when an adaptive allele arose before
59 selection, but multiple copies have subsequently swept through the population (single origin soft
60 sweep). Soft sweeps are often incorrectly viewed synonymously with standing genetic variation;

61 hard sweeps can emerge from standing genetic variation if a single copy of the beneficial allele
62 was the ancestor of all beneficial alleles in a sample (Hermisson and Pennings 2017).

63 Hard and soft sweeps are locus-specific and, thus, not mutually exclusive across a
64 genome. Unsurprisingly, soft sweeps are also much more difficult to recognize than hard sweeps
65 because their genomic patterns are intermediate. Additionally, the identification of selective
66 sweeps, hard or soft, is further complicated by the possibility that neutral loci linked to either soft
67 or hard sweeps may produce a false signature similar to that of a sweep (Schrider et al. 2015;
68 Kern and Schrider 2018).

69 With these challenges in mind, a considerable amount of work has been dedicated to both
70 developing robust methods to identify selective sweeps and also understanding the evolutionary
71 parameters that determine hard or soft sweeps. Mutation-limited scenarios are expected to
72 exclusively produce hard sweeps because beneficial alleles rarely occur (Hermisson and
73 Pennings 2017). Thus, the most important parameter for estimating the likelihood of hard vs soft
74 sweeps is the population-scaled mutation rate: $\theta = 4N_e\mu$, where N_e is the effective population size
75 and μ is the mutation rate. However, this single parameter can vary widely depending on the
76 advantage of the beneficial allele, the effective population size, the size of the mutational target,
77 and the timescale for adaptation (Messer and Petrov 2013; Hermisson and Pennings 2017).

78 While it has become clear that most populations will likely exhibit a mosaic of hard and soft
79 sweeps (Hermisson and Pennings 2017), additional data on sweep type frequencies in various
80 species are sorely needed to better tease apart which parameters may determine each of those
81 frequencies.

82 Both species of the *Pan* genus represent important evolutionary models due to their
83 phylogenetic proximity to humans. *Homo* and *Pan* diverged ~ 5 to 7 Ma (Sarich and Wilson

84 1967; Bradley 2008; Scally et al. 2012; Besenbacher et al. 2019) and the most recent estimates
85 for the divergence of bonobos and chimpanzees range between 1 and 2 Ma (Prüfer et al. 2012; de
86 Manuel et al. 2016). Four extant chimpanzee subspecies evolved from a chimpanzee common
87 ancestor that split ~ 600 Ka with both subsequent lineages further splitting: one ~ 250 Ka and the
88 other ~ 160 Ka (de Manuel et al. 2016). These two species exhibit stark differences in aspects of
89 their morphology, physiology, behavior, and ecology (Susman 1984; Goodall 1986; Wrangham
90 1986; Kano 1992; White 1996; Furuichi 2011; Nishida 2011; Stumpf 2011; Behringer et al.
91 2014; Turley and Frost 2014; Wilson et al. 2014). Many of these distinguishing traits are inferred
92 to have occurred shortly after divergence, while much less is known about recent evolutionary
93 processes in these lineages.

94 Understanding recent positive selection in *Pan* is intriguing because of the dynamic
95 physical and social environments in which they evolved. Climatic variation across Africa is well-
96 documented for the Pleistocene and has been proposed to drive the evolution of *Homo* (Potts
97 1998; Antón et al. 2014), and such variation probably impacted other taxa during this time
98 period, including the genus *Pan*. Chimpanzee populations living in more stable environments
99 that were closer to Pleistocene refugia were recently described to exhibit less behavioral
100 diversity than chimpanzees living in more seasonal habitats that are more distant to forest refugia
101 (Kalan et al. 2020). While the formation of these refugia may have resulted in periods of habitat
102 stability for some bonobo and chimpanzee populations during glacial periods (Takemoto et al.
103 2017; Barratt et al. 2020), climatic fluctuations throughout the Pleistocene likely affected both
104 the physical environment—via changes in habitat structure and type—and the social
105 environment—via changes in the frequency of dispersal and intergroup encounters. Further,
106 evidence of admixture within extant and between extant and extinct members of the *Pan* genus

107 adds even more variation to the social environments in which these apes evolved (Hey 2010;
108 Wegmann and Excoffier 2010; de Manuel et al. 2016; Kuhlwilm et al. 2019). A dynamic
109 environment may result in selection for multiple existing alleles, resulting in a greater frequency
110 of soft sweeps than in a more stable environment where one would expect a greater frequency of
111 hard sweeps.

112 In this study, we apply a recently developed supervised machine-learning approach to
113 population-level genomic data for bonobos (*Pan paniscus*) and chimpanzees (*Pan troglodytes*) to
114 assess the extent of different sweep types in these species. While a few studies have examined
115 recent positive selection in bonobos and chimpanzees (e.g., Cagan et al. 2016; Han et al. 2019;
116 Schmidt et al. 2019; Nye et al. 2020), the role of hard and soft sweeps in shaping their
117 adaptations is currently unknown. We sought to categorize genomic regions as subject to recent
118 hard or soft sweeps, as linked to recent hard or soft selective sweeps, or as evolving neutrally.
119 Data from simulations have predicted that hard sweeps would be common in humans because of
120 our low mutation rate (Hermisson and Pennings 2017). Under this “mutation limitation
121 hypothesis” and given the similarity in mutation rate between *Homo* and *Pan*, one could predict
122 that bonobos and chimpanzees should also exhibit a high degree of hard sweeps. However, hard
123 sweeps appear quite rare in recent human evolution (Hernandez et al. 2011; Schrider and Kern
124 2017) and adaptation in humans may not be mutation-limited. This could be explained by several
125 non-mutually exclusive alternatives including demographic effects. Larger populations can have
126 more standing variation for selection to act on (Hermisson and Pennings 2005) which may result
127 in more soft sweeps whereas bottlenecks can result in drift and thus potentially more hard
128 sweeps if intermediate frequency haplotypes are lost. For example, humans have experienced
129 recent demographic changes (e.g., Schiffels and Durbin 2014), including a bottleneck upon

130 leaving Africa (e.g., Henn et al. 2012). Indeed, Schrider and Kern (2017) found that hard sweeps
131 were more frequent in non-African than African populations. Chimpanzees and bonobos have
132 also experienced recent demographic changes, including in effective population size, within the
133 time frame (< 200 Ka) for selective sweeps, based on PSMC analyses (Prado-Martinez et al.
134 2013; de Manuel et al. 2016). We therefore predicted that we would observe a higher frequency
135 of soft sweeps in *Pan*, but that lineage-specific population histories might affect the degree to
136 which soft sweeps dominate.

137

138 **Methods**

139 *Genomic Data*

140 We retrieved raw short read data on bonobos and all four chimpanzee subspecies from
141 the Great Ape Genome Project (GAGP) (Prado-Martinez et al. 2013). This dataset contained
142 high coverage genomes (Figures S1, S2) from 13 bonobos (*P. paniscus*), 18 central chimpanzees
143 (*P. troglodytes troglodytes*), 19 eastern chimpanzees (*P. t. schweinfurthii*), 10 Nigeria-Cameroon
144 chimpanzees (*P. t. ellioti*), and 11 western chimpanzees (*P. t. verus*) (File S1).

145

146 *Read Mapping and Variant Calling*

147 Initial quality assessments in fastqc (Andrews 2010) and multiqc (Ewels et al. 2016)
148 indicated a number of quality issues, including failed runs, problematic tiles, and substantial
149 variation in base quality. We removed adapters and trimmed all reads for quality with BBduk
150 (<https://sourceforge.net/projects/bbmap/>). For trimming, we used the parameters “ktrim=r k=21
151 mink=11 hdist=2 qtrim=rl trimq=15 minlen=50 maq=20” for all reads and added “tpo and tpe”
152 for paired reads.

153 We used XYalign (Webster et al. 2019) to create versions of the chimpanzee reference
154 genome, panTro6 (Kronenberg et al. 2018), for male- and female-specific mapping. Specifically,
155 the version of the reference for female mapping has the Y chromosome completely masked, as
156 its presence can lead to mismapping (Webster et al. 2019). We then mapped reads with BWA
157 MEM (Li 2013) and used SAMtools (Li et al. 2009) to fix mate pairs, sort BAM files, merge
158 BAM files per individual, and index BAM files. We use Picard (Broad Institute 2018) to mark
159 duplicates with default parameters, before calculating BAM statistics with SAMtools. We next
160 measured depth of coverage with mosdepth (Pedersen and Quinlan 2018), removing duplicates
161 and reads with a mapping quality less than 30 for calculations. Visualizations for coverage and
162 demography (see Generation of Simulated Chromosomes below) were created in R, version 3.5.2
163 (R Core Team 2020), using ‘ggplot2’ (Wickham 2016).

164 We used GATK4 (Poplin et al. 2018) for joint variant calling across all samples. We used
165 default settings for all steps—HaplotypeCaller, CombineGVCFs, and GenotypeGVCFs—with
166 three exceptions. First, we turned off physical phasing for computational efficiency and
167 downstream VCF compatibility with filtering tools. Second, because multiple samples in this
168 dataset suffer from contamination from other samples both within and across taxa (Prado-
169 Martinez et al. 2013), we employed a contamination filter to randomly remove 10% of reads
170 during variant calling. This should have the effect of reducing confidence in contaminant alleles.
171 Finally, we output non-variant sites to allow equivalent filtering of all sites in the genome and
172 more accurate assessments of callability.

173 The above quality control, assembly, and variant calling steps are all contained in an
174 automated Snakemake (Köster and Rahmann 2012) available on Github
175 (https://github.com/thw17/Pan_reassembly). The repository also contains a Conda environment

176 with all software versions and origins, most of which are available through Bioconda (Grüning et
177 al. 2018).

178

179 *Variant Filtration and Genome Accessibility*

180 We considered only autosomes for this analysis as the X and Y chromosome violate
181 many of the assumptions for the following methods (Webster and Wilson Sayres 2016). We also
182 excluded unlocalized scaffolds (N = 4), unplaced contigs (N = 4,316), and the mitochondrial
183 genome from any downstream analyses. Additional filtration steps were completed using
184 bcftools (Li 2011); command line inputs are provided in parentheses. Given our focus on
185 selective sweeps, we only included single nucleotide variants (SNVs) (“-v snps”) that were
186 biallelic (“-m2 -M2”). On a per sample basis within each site, we marked genotypes where
187 sample read depth was less than 10 and/or genotype quality was less than 30 as uncalled (“-S . -i
188 FMT/DP \geq 10 && FMT/GT \geq 30”). To ensure that missing data did not bias our results, we
189 further excluded any sites where less than ~ 80% of individuals (N = 56) were confidently
190 genotyped (“AN \geq 112”). We also removed any positions that were monomorphic for either the
191 reference or alternate allele (“AC > 0 && AC \neq AN”). These filtrations steps yielded 41,869,892
192 SNVs for our downstream analyses (Table S1).

193 We considered sites in our sample with low to no coverage to be ‘inaccessible’ in the
194 reference genome. Using the output of mosdepth (see Read Mapping and Variant Calling above),
195 we identified and filtered sites exhibiting low coverage as defined above. We used the
196 ‘maskfasta’ function in bedtools (Quinlan and Hall 2010) to mark these sites (N) in the pantro6
197 FASTA, featuring only the autosomes, for use in downstream analyses. This resulted in 86.3% of
198 the assembled autosomes as accessible (File S2).

199

200 *Generation of Simulated Chromosomes*

201 We used the software ‘discoal’ to generate simulated chromosomes on which we trained
202 a classifier per lineage (Kern and Schrider 2016). We generated a matching number of simulated
203 haploid chromosomes for the sample size of each *Pan* lineage (i.e., 26 chromosomes for 13 *P.*
204 *paniscus*, 20 chromosomes for 10 *P. t. ellioti*, etc.). Simulated chromosomes were set to 1.1 Mb
205 in length and divided into 0.1 Mb subwindows for a total of 11 subwindows. These simulations
206 included a population-scaled mutation rate ($4N\mu L$), where N is the effective population size, μ is
207 the per base pair per generation mutation rate, and L is the length of the simulated chromosome.
208 We used the median of the previously reported effective population size range per lineage
209 (Prado-Martinez et al. 2013). As estimates of genome-wide mutation rates vary considerably and
210 are complicated in that mutation rates vary across individual genomes, we based our parameter
211 on a mutation rate of 1.6×10^{-8} , which falls between estimates from genome-wide data and
212 phylogenetic estimates (Narasimhan et al. 2017). We introduced some variation in this rate by
213 setting a lower and upper-bound to 1.5 and 1.7×10^{-8} and sampled a new mutation rate per
214 simulation drawing from this uniform prior. All simulations also included a population-scaled
215 recombination rate ($4NrL$), where r is the recombination rate per base pair per generation, again
216 calculated from the median effective population size for each lineage from Prado-Martinez et al.
217 (2013) and a recombination rate drawn from a uniform prior of $1.1 - 1.3 \times 10^{-8}$, based on the
218 mean genome-wide rate (1.2×10^{-8}) reported for bonobos, chimpanzees, and gorillas (Stevison et
219 al. 2015). We note that while some of the estimated recombination rates in bonobos and
220 chimpanzees are beyond the uniform distribution used in our simulations, many of these values
221 are the high rates present in the telomeres, regions that generally exhibit lower or no coverage

222 and thus will be largely if not entirely masked from this analysis (see Variant Filtration and
223 Genome Accessibility above). We also included a demographic string reflecting approximate
224 changes in population size for each lineage between ~ 0.05 and 2 Ma. Changes in population size
225 were set in units of $4N_0$ generations, N_0 was set to the approximate median effective population
226 size from (Prado-Martinez et al. 2013) and we used a generation time of 25 years (Langergraber
227 et al. 2012). Population size changes for this time period were drawn from a previous PSMC
228 analysis (de Manuel et al. 2016) (Figure S3). While this is only one study from which to draw
229 demographic information and reconstructions of *Pan* demography vary widely across studies, the
230 downstream program used to classify genomic windows, diploS/HIC, is robust to demographic
231 misspecification (Kern and Schrider 2018). We generated 2×10^3 simulations using these
232 parameters as a set of simulations under neutral evolution per lineage.

233 Hard and soft selective sweeps were simulated with all of the aforementioned parameters
234 and using a uniform prior of population-scaled selection coefficients ($\alpha = 2Ns$) derived from each
235 lineage's median effective population size (Prado-Martinez et al. 2013) and moderately weak to
236 moderately strong selection coefficients between 0.02 and 0.05. Sweeps also included a
237 parameter (τ) for the time to fixation of the beneficial allele over a uniform range in units of $4N$
238 generations. This value ranged from 0 to 0.001 for all lineages. Linked-hard and linked-soft
239 sweeps were generated by placing the selected site at the center of each of the 10 subwindows
240 flanking the center (6th) subwindow. Additionally, we included a uniform prior on the frequency
241 at which a mutation is segregating at the time it becomes beneficial for soft and linked-soft
242 sweeps, setting this range from 0 to 0.2. We generated 1×10^3 simulations per subwindow for
243 linked-hard and linked-soft sweeps ($N = 10$) and 2×10^3 simulations for hard and soft sweeps.

244 This resulted in a total of 2×10^3 hard, 1×10^4 hard-linked, 2×10^3 soft, and 1×10^4 soft-linked
245 simulated sweeps. Parameters for these simulations are presented in File S3.

246

247 *Calculation of Simulation Feature Vectors and Classifier Training*

248 We calculated feature vectors from these simulated chromosomes using the ‘fvecSim’
249 function in the program diploS/HIC (Kern and Schrider 2018). Briefly, diploS/HIC calculates 12
250 summary statistics for all 11 subwindows: π , Watterson’s θ , Tajima’s D , the variance, skew, and
251 kurtosis of genotype distance (g_{kl}), the number of multilocus genotypes, J_1 , J_{12} , J_2/J_1 , unphased
252 Z_{ns} , and the maximum value of unphased ω . Collectively, these summary statistics capture
253 information about the site frequency spectrum (SFS), haplotype structure, and linkage
254 disequilibrium (LD). diploS/HIC uses a convolutional neural network (CNN) to capture essential
255 aspects of a feature (the feature vector) by sliding a receptive field over the image to compute dot
256 product between the original filter and the convolutional filter. In diploS/HIC, the CNN uses
257 three branches of a CNN, of which each has two dimensional convolutional layers with ReLu
258 activations followed by max pooling. This is followed by a dropout layer to control for model
259 overfitting. Outputs from all three units are fed into two fully connected dense layers, which also
260 use dropout layers, before arriving at a softmax activation that outputs the probability for each
261 categorical class (hard, hard-linked, neutral, soft-linked, or soft). Complete details for this
262 procedure can be found in Kern and Schrider (2018).

263 When calculating feature vectors for the simulated chromosomes, we used the optional
264 arguments for the ‘fvecSim’ function to mask each simulation with 110,000 bp segment
265 randomly drawn from our masked FASTA where > 0.25 of SNVs in a subwindow were
266 accessible (i.e., not marked by Ns). This enabled us to train our classifiers on simulated data

267 featuring the same patterns of inaccessible genomic regions that the classifier would encounter in
268 the empirical data.

269 We created a balanced set with equal representation (2×10^3) of all five classes via
270 sampling without replacement in which to train the classifier using diploS/HIC's
271 'makeTrainingSets' function. These were divided into 8,000 training examples, 1,000 validation
272 examples, and 1,000 testing examples to test the accuracy of the classifier via the 'train' function
273 in diploS/HIC. We built ten classifiers per lineage and selected the one with the highest accuracy
274 to apply to the empirical data (File S4).

275 A second, independent set of simulated chromosomes was generated per lineage using
276 the same parameters. After calculating feature vectors and creating a balanced training set, we
277 used diploS/HIC's 'predict' function to assess the true positive rate, false positive rate, and
278 accuracy of each classifier (Tables S2 - S5).

279

280 *Empirical Data Feature Vectors and Prediction*

281 Upon achieving > 0.8 accuracy, each trained classifier was applied to its respective *Pan*
282 lineage. Each autosome was analyzed separately and feature vectors calculated using
283 diploS/HIC's 'fvecVcf' function. We supplied this function with the masked FASTA for that
284 chromosome and discarded windows where any subwindow had < 0.25 unmasked sites
285 following Schrider and Kern (2017) (File S5). This step reduces the potential effect of the
286 number of SNVs in a given window on sweep classification. Finally, the trained classifier was
287 applied to the feature vector files using the 'predict' function.

288

289 *Sweep Identification, Potential Target Genes, and Gene Ontology*

290 As diploS/HIC outputs the probability for each sweep class, we first report the class
291 inferred to be the most likely. However, as the difference between the most likely class and the
292 next most likely may be small, we further report windows where the sweep class probability is $>$
293 0.5 , > 0.75 , and > 0.9 (File S6). We also examined our data for spatial patterns. Windows
294 classified as immediately abutting other windows with the same sweep type for hard and soft
295 sweeps were considered to be a single sweep. Unique sweep windows and those shared between
296 two or more lineages were visualized using UpSet plots (Lex et al. 2014) in R (R Core Team
297 2020).

298 We examined what genes lie in the windows identified as being subject to a recent
299 selective sweep by extracting the genomic coordinates of all autosomal coding regions for the
300 longest transcript per gene ($N = 20,119$ genes) in the panTro6 genome via the panTro6 gff
301 (retrieved from: https://www.ncbi.nlm.nih.gov/genome/202?genome_assembly_id=380228). We
302 used the bedtools ‘intersect’ function (Quinlan and Hall 2010) to identify overlap between
303 coding regions and candidate sweep windows after converting both CDS and sweep window
304 coordinates to 0-start, half-open format. As some coding sequences may have been masked (see
305 Variant Filtration and Genome Accessibility above), we extracted FASTAs for each coding
306 sequence using bedtools ‘getfasta’ function (Quinlan and Hall 2010) and used a custom R script
307 to calculate the percent of each gene that was masked. Overall, 66.2% of all coding sequence
308 was unmasked. We excluded listing genes for candidate sweep regions if $> 50\%$ of the total
309 coding sequence per gene was masked. Thus, we considered 13,228 genes as potential targets for
310 selective sweeps (File S7).

311 We investigated the enrichment of particular pathways by performing a gene ontology
312 analysis using the Functional Annotation Tool in DAVID (Huang et al. 2008; Huang et al. 2009).

313 We used the custom background described above (genes whose total coding sequence was >
314 50% unmasked) rather than all pantro6 genes to ensure our analysis was not underpowered.
315 DAVID does not allow for official gene symbols to be used in a background list, so we
316 converted gene symbols to Entrez gene IDs. As not all gene symbols have a corresponding
317 Entrez gene ID, we removed genes for which there was no Entrez gene ID (N = 98 in
318 background list). We collated genes for both hard and soft sweeps into a single input per lineage.
319 We evaluated statistical significance for biological process gene ontology terms via p-values
320 adjusted using the Benjamini-Hochberg method (Benjamini and Hochberg 1995).

321 Scripts for all data analyses are available on Github
322 (https://github.com/brandcm/Pan_Selective_Sweeps).

323

324 **Results**

325 We generated four classifiers that reached an acceptable level of accuracy for bonobos
326 (*P. paniscus*), central chimpanzees (*P. t. troglodytes*), eastern chimpanzees (*P. t. schweinfurthii*),
327 and Nigeria-Cameroon (*P. t. ellioti*) chimpanzees. These classifiers ranged in accuracy from
328 85.6% (Nigeria-Cameroonian chimpanzees) to 93.9% (central chimpanzees) (File S4). We could
329 not produce a sufficiently accurate classifier using realistic parameters for western chimpanzees
330 (*P. t. verus*); therefore, they were excluded from downstream analyses. Following Kern and
331 Schrider (2018), we calculated false positive rates by testing our classifiers on a second,
332 independent set of simulated chromosomes per lineage. We used a binary classification,
333 considering the identification of either sweep type as a positive and identification of a linked or
334 neutral region to be negative. Our trained classifiers had considerable statistical power (1 - false
335 positives) ranging from 96.6 to 99.2% and a low false positive rate (false positives / false

336 positives + true negatives) that ranged from 1.4 to 4.3% across all four classifiers (Tables S2 -
337 S5). When considered separately—i.e., true positives only included one sweep type (hard or soft)
338 rather than both—we had greater power to detect hard sweeps than soft sweeps, averaging 99%
339 and 96.9% across lineages, respectively (Tables S2 - S5). Accuracy (true positives + true
340 negatives / total) for identifying sweep regions vs non-sweep regions ranged from 94.1 to 98.3%
341 while a second estimate (in addition to the first accuracy estimate that resulted from the
342 construction of the classifiers) of class-specific accuracy ranged from 81.6 to 92.1% (Tables S2 -
343 S5).

344 We classified ~ 91.6% of the assembled autosomes in each lineage (Table 1, File S8),
345 even after masking for inaccessible regions and excluding windows with few SNVs. We found
346 that soft sweeps were abundant in all four lineages, accounting for > 73% of all individual
347 sweeps, whereas hard sweeps were relatively rare (Table 1, File S8). This pattern held true even
348 when more stringent posterior probabilities were applied to consider a region a sweep and at
349 least 30% of hard sweep windows and 76% of soft sweep windows were called with 50% or
350 greater posterior probability (File S6). Genomic regions linked to sweeps were also quite
351 pervasive in all four lineages (Table 1); particularly among eastern chimpanzees, where roughly
352 86% of the genome was classified as linked to selective sweeps.

353 We examined overlap in windows classified as either a hard or soft sweep across
354 lineages, which may reflect either ancestral or parallel adaptation. Most hard sweep windows
355 were unique to each lineage; however, we did find some shared windows across lineages (Figure
356 1). Central and Nigeria chimpanzees shared the highest number of sweep windows (N = 33) but
357 when weighted by the total possible number of windows, the highest overlap for hard sweeps
358 was between eastern and Nigeria chimpanzees (7/32 or ~ 0.21). No hard sweeps windows were

359 shared across all lineages. Like hard sweeps, most soft sweep windows were also unique to each
360 lineage (Figure 2). Among pairs of lineages there was remarkable consistency in the number of
361 shared windows (N = 111-147), even when the total possible number of shared windows is
362 considered. One exception is eastern and central chimpanzees who shared nearly twice the
363 number of soft sweep windows (N = 267). The highest number of shared soft sweep windows
364 between three lineages occurred in the three chimpanzee subspecies (N = 80). Only 19 windows
365 were shared across all four lineages.

366 After excluding genes that were > 50% masked, we identified 1,671 candidate genes in
367 bonobo hard and soft sweeps, 1,761 genes in central chimpanzee sweeps, 1,372 genes in eastern
368 chimpanzee sweeps, and 1,844 genes in Nigeria-Cameroonian chimpanzee sweeps (File S9).
369 After correcting for multiple testing, across all lineages, we identified only two significantly
370 enriched pathways in central chimpanzees: nervous system development and central nervous
371 system development (File S10).

372

373 **Discussion**

374 Our study contributes to the emerging picture of recent evolution in *Pan* and adaptation
375 more broadly. Contrary to the predictions of a mutation-limitation hypothesis, yet concordant
376 with recent results for humans (e.g., Hernandez et al. 2011; Schrider and Kern 2017), we find
377 soft sweeps to overwhelmingly predominate regions of the genome experiencing selective
378 sweeps in both bonobos and the three chimpanzee subspecies we could analyze. These results
379 confirm the prediction from Schmidt et al. (2019) who speculated that soft sweeps played a
380 major role in the evolution of eastern and central chimpanzees. Those authors also posit that hard
381 sweeps should be more frequent in western chimpanzees relative to other subspecies because of

382 their low effective population size. While western chimpanzees are estimated to have the lowest
383 effective population size, it is estimated to be only slightly lower than that of bonobos for which
384 we found a high number (95.1%) of soft sweeps (e.g., Prado-Martinez et al. 2013; de Manuel et
385 al. 2016). It is curious that Nigeria-Cameroon chimpanzees exhibit the most hard sweeps in this
386 analysis. While this could be the result of a multitude of factors, a notable possibility is that this
387 lineage has experienced the most stable effective population size in recent evolutionary time as
388 estimated by PSMC, compared to bonobos, eastern chimpanzees, and central chimpanzees
389 (Prado-Martinez et al. 2013; de Manuel et al. 2016).

390 Our analysis of shared hard and soft sweeps found that most sweeps of both types were
391 unique to each lineage. However, there was a high number of hard sweep windows shared
392 between central and Nigeria-Cameroon chimpanzees as well as between eastern and Nigeria-
393 Cameroon chimpanzees when the total possible number of shared sweeps was considered.
394 Further, there were nearly twice the number of shared soft sweep windows shared between
395 eastern and central chimpanzees. These results are similar to other recent findings (Nye et al.
396 2020). It is impossible to discern whether or not the overlap in hard sweeps between central and
397 Nigeria-Cameroon chimpanzees and the overlap in soft sweeps for eastern and central
398 chimpanzees is the result of shared ancestry and/or similar environmental conditions because
399 both pairs of lineages share a geographic boundary: the Ubangi river for eastern and central
400 chimpanzees and Sanaga river for central and Nigeria-Cameroon chimpanzees. The overlap in
401 hard sweeps between eastern and Nigeria-Cameroon chimpanzees is more puzzling because they
402 are not sister taxa and share a common ancestor ~ 600 Ka. Therefore, parallel adaptation via
403 similar physical and/or social environments may serve as a more likely hypothesis. While the
404 lowest in overall frequency, we also identified a number of soft sweep windows that were shared

405 across three lineages as well as 19 windows that occurred in all four. Future work should further
406 investigate these shared sweep windows.

407 As mentioned above, soft sweeps are not exclusively the result of selection on standing
408 genetic variation (Pennings and Hermisson 2006a; Pennings and Hermisson 2006b). However,
409 given the mutation rates estimated for bonobos and chimpanzees, it appears unlikely that
410 recurrent *de novo* mutations explain the majority of these soft sweeps. We did not explicitly
411 model for different types of soft sweeps in our analysis. However, while soft sweeps from
412 standing genetic variation and *de novo* mutations may exhibit similar genomic signatures, the
413 hypothesis that these processes result in similar genomic signatures must be tested before any
414 additional conclusions are drawn. Nonetheless, our results reveal a major role of standing genetic
415 variation, and thus changes in the physical and social environment, in driving recent adaptations
416 in *Pan*.

417 A few recent studies have considered the impact of effective population size on adaptive
418 evolution in the great apes (Cagan et al. 2016; Nam et al. 2017). Theory predicts that the rate of
419 adaptive evolution should be positively correlated with effective population size when $N_e s$ is \gg
420 1 (Gossmann et al. 2012). Both Cagan et al. (2016) and Nam et al. (2017) found a positive
421 association between effective population size and the rate of adaptive evolution, measured by
422 proportion of adaptive substitutions and the number of selective sweeps, respectively. However,
423 we observed no clear linear relationship between the number of sweeps (hard, soft, or both)
424 estimated from this analysis and the estimated effective population sizes for these four lineages
425 (see File S3 for population sizes). This descriptive result should be considered cautiously
426 because of the limited number of lineages analyzed here and the potential confounding effect of
427 phylogeny. It is possible that this relationship may not be driven by the number of sweeps, but

428 rather the strength of sweeps a population experiences (Nam et al. 2017). Estimates of selection
429 strength are generally lacking for the great apes so this relationship remains a question for further
430 study.

431 In addition to characterizing broad patterns in the genomic landscape for bonobos and
432 chimpanzees, the results of this study also highlight thousands of candidate regions and genes for
433 further analysis. We also find additional support for previous selection candidates. For example,
434 disease has been long thought to shape evolution in primates (Nakajima et al. 2008; van der Lee
435 et al. 2017). The potential for disease transmission between non-human primates and humans has
436 also prompted much research, particularly focusing on the genomic underpinnings of host
437 responses to lentiviruses, which include HIV and SIV (Gao et al. 1999; Van Heuverswyn et al.
438 2006; Compton et al. 2013; Nakano et al. 2020). Cagan and colleagues (2016) found evidence of
439 recent positive selection within *IDO2*, a T-cell regulatory gene, among all four-chimpanzee
440 subspecies and bonobos. We identified a candidate soft sweep region for eastern chimpanzees
441 that overlaps this gene. However, this window had one of the lowest posterior probabilities in
442 this lineage (49.7%) and there was a nearly equally high probability that this window was linked
443 to a soft sweep (43.8%). Clearly, additional work is needed to understand the potential role of
444 *IDO2* in *Pan* evolution. Schmidt et al. (2019) recently described three chemokine receptor
445 genes—*CCR3*, *CCR9*, and *CXCR6*—had a significant number of highly differentiated SNVs in
446 central chimpanzees. We could evaluate all three of these genes in our analysis but only one fell
447 within a candidate sweep window: *CXCR6*. The window containing this gene was confidently
448 called as a soft sweep with a posterior probability of 85.5%. It is not known as to whether or not
449 *SIV_{cpz}* uses *CXCR6* to enter chimpanzee host cells (Wetzel et al. 2018). However, multiple lines
450 of evidence for selection either at this locus or within the window overlapping this gene prompt a

451 closer examination of this genomic region. Finally, *TRIM5* fell within a hard sweep window in
452 central chimpanzees. *TRIM5* is a well-known retrovirus restriction factor that appears subject to
453 ancient, multi-episodic positive selection in primates (Sawyer et al. 2005).

454 Recent attention has focused on admixture between lineages in the genus *Pan* and the
455 potential adaptiveness of introgressed genomic elements. de Manuel and colleagues (2016)
456 identified 221 genes that fell within putatively introgressed elements in central chimpanzees
457 from admixture with bonobos. Some of this admixture is estimated to occur < 200 Ka, thus
458 within the timeframe that the present analysis can detect selective sweeps. While we could not
459 evaluate six of these 221 genes, five fell within candidate sweep regions in central chimpanzees
460 from our study: *CDK8*, *EIF4E3*, *GRID2*, *PTPRM*, and *TRIM5*. As described above, *TRIM5* was
461 unique to central chimpanzees. We found *CDK8* in sweep windows for bonobos, eastern
462 chimpanzees, and Nigeria-Cameroon chimpanzees. In humans, *CDK8* mutations have been
463 associated with multiple phenotypic effects including hypotonia, behavioral disorders, and facial
464 dysmorphism (Calpena et al. 2019). We also identified *EIF4E3* in candidate sweeps for bonobos
465 whereas *GRID2* and *PTPRM* were found in eastern chimpanzees. *EIF4E3* is a translation
466 initiation factor (Osborne et al. 2013) while *PTPRM* is a member of the protein phosphatase
467 family (PTP) and has multiple functions including cell proliferation and differentiation (Sun et
468 al. 2012). *GRID2* generates ionotropic glutamate receptors and mutations have been associated
469 with abnormalities of the cerebellum (Lalouette et al. 1998).

470 The gene ontology analysis produced only two statistically significant terms, nervous
471 system development and central nervous system development, for a single *Pan* lineage: central
472 chimpanzees. While cognitive and neurological differences are widely considered to differentiate
473 bonobos and chimpanzees (e.g., Rilling et al. 2012; Stimpson et al. 2016; Staes et al. 2019), we

474 are unaware of any studies that investigate variation among chimpanzee subspecies that may
475 explain enrichment for nervous system and central nervous system development related genes
476 specifically in central chimpanzees. We note that compared to other gene ontology analyses, our
477 level of enrichment is quite low. While we excluded a large number of genes from our analysis
478 due to poor coverage, our use of a custom background should increase, rather than decrease,
479 statistical power.

480 The results from our analysis should be interpreted with some caution. First, while our
481 classifiers achieved a high degree of accuracy, it is possible that some selective sweeps in each
482 lineage were not detected or regions were incorrectly identified as such (Tables S2 - S5). We
483 also note that we did not model small selection coefficients as we could not accurately classify
484 sweeps under weak selection. Overall, our classifiers were quite good at identifying hard and
485 linked-hard sweeps with both at approximately 95% accuracy across all lineages. Neutral and
486 linked-soft regions were the most difficult to recognize with neutral regions typically being
487 classed as soft-linked when they did not appear neutral. This suggests that the neutral portion of
488 the genome for each lineage is slightly underestimated here. Finally, some soft sweeps were
489 identified as hard sweeps in each of our classifiers, suggesting that some portion of identified
490 hard sweeps in each lineage are, in fact, soft sweeps. The low false positive rates demonstrate the
491 overall accuracy of the observed genomic patterns (i.e., the proportion of hard and soft sweeps)
492 for these taxa. However, this point underscores the need to conduct subsequent analyses of the
493 candidate regions and genes to confirm such the proposed mode of adaptation and investigate
494 any functional consequences of that adaptation. In the ‘era of -omics’, the generation of
495 candidate regions for any type of selection across populations and species appears to
496 overwhelmingly outpace the confirmation of such patterns. Avenues of research that investigate

497 these candidate genes in more detail are thus well poised to provide a deeper and more accurate
498 understanding of lineage-specific adaptations.

499 Second, background selection, the loss of a linked neutral site from purifying selection on
500 a deleterious allele, can potentially mimic patterns of selective sweeps and thus may impact the
501 results of this study (Charlesworth et al. 1993). We did not explicitly model background
502 selection in our analysis, however, evidence from simulations in various taxa demonstrate that
503 this pattern of selection does not substantially increase the rate of false positives in selective
504 sweep analyses (Schrider and Kern 2017; Schrider 2020: 20). Further, Nam et al. (2017)
505 considered the effect of background selection on genomic diversity in extant apes, including all
506 five *Pan* lineages, and note that background selection alone does not produce the observed
507 diversity reduction near genic regions in these lineages.

508 Further, sampling bias can reduce the accuracy of identifying selective sweeps. If
509 multiple haplotypes are present in a population but only individuals sharing one haplotype are
510 sampled, then the sweep would be classified as a hard sweep when it is a soft sweep. However,
511 this scenario would only underestimate the degree of recent adaptation from soft sweeps.
512 Therefore, if this sampling bias is present in this analysis, then soft sweeps may predominate
513 recent *Pan* evolution to an even larger degree than described here. Population structure adds
514 further complications to the classification of hard sweeps. Parallel adaptation produces multi-
515 origin soft sweeps at the global population level that would appear to be hard in local
516 populations, although even local samples may sometimes appear to be soft sweeps (Ralph and
517 Coop 2010). Thus, if samples stemmed from one or few local populations then global soft
518 sweeps may be misclassified as hard. A previous analysis estimated the geographic origin of
519 individuals used in this analysis (de Manuel et al. 2016). These authors found that individuals

520 from both eastern and central chimpanzee populations were sampled from multiple countries
521 across the geographic range for both subspecies. Therefore, any hard sweeps detected in these
522 populations are likely accurate at the subspecies level. Geographic origin could not be assessed
523 for any of the bonobos or all of the Nigeria-Cameroon chimpanzees used in this analysis (de
524 Manuel et al. 2016). As such, sampling or geographic bias may partially explain the high degree
525 of hard sweeps observed in Nigeria-Cameroon chimpanzees, if they were sampled from a smaller
526 geographic area than the other subspecies. We encourage future studies to consider this potential
527 bias when hard sweeps are encountered in existing data and during study design.

528 This analysis focuses on signatures of positive selection at single loci. However, there is
529 theoretical and empirical evidence that a number of adaptive traits have a complex, multilocus
530 architecture (Pritchard et al. 2010; Yang et al. 2017; Bergery et al. 2018). For these polygenic
531 traits, shifts in the physical or social environment might result in allele frequency changes at
532 many loci, of which, according to models, few to none of which would reach fixation (Pritchard
533 et al. 2010). This may, in part, explain why hard sweeps appear to be rare in humans and other
534 species if it represents a dominant mode of adaptation in these taxa. Unfortunately, at this point,
535 we lack the data and methods to investigate the extent of polygenic selection across the genome
536 in many non-model taxa such as *Pan*. It is also worthwhile to address that this analysis focused
537 on modelling very recent completed selective sweeps. Another future avenue of study is the
538 identification of incomplete or partial sweeps in bonobos and chimpanzees.

539 Finally, while our approach to identifying hard and soft sweeps is a logical first step,
540 future work should consider sweeps within subspecies to assess population-level (i.e., local),
541 rather than lineage-specific, adaptations. This is underscored by the extensive phenotypic
542 variation among chimpanzees, particularly that of behavioral variation, which includes key

543 characteristics that are often used to dichotomize bonobos and chimpanzees (Wilson et al. 2014).
544 Further investigation is also clearly warranted in bonobos, whose overall phenotypic variation is
545 likely underappreciated compared to chimpanzees (Hohmann and Fruth 2003; Sakamaki et al.
546 2016; Beaune et al. 2017; Wakefield et al. 2019).

547

548 **Conclusion**

549 This study highlights the importance of changes in physical and/or social environment via
550 soft selective sweeps in the recent evolution of our closest living relatives, chimpanzees and
551 bonobos. Our results also yield further support for the ubiquity of soft, rather than hard, sweeps
552 in adaptation. We contribute candidate regions and genes that may help identify unique
553 phenotypes in each *Pan* lineage. Our findings also prompt many new questions including the
554 estimation of selection strength coefficients and the degree of haplotypic diversity in candidate
555 sweep regions. While our study focuses on these lineages broadly, this point also underscores the
556 need for high-coverage genomic data collected using non-invasive methods at more local
557 geographies.

558

559 **Acknowledgements**

560 We thank Andy Kern for help with implementing this analysis. Hazel Byrne, Tina Lasisi,
561 Alan Rogers, Liz Tapanes, and Andrew Zamora provided valuable comments on this manuscript.
562 We also thank Elisabeth Goldman and Noah Simons for assistance with bioinformatics. We
563 gratefully acknowledge Brad Sherman (NIH) who provided assistance with our gene ontology
564 analysis. We thank Mark Allen, Mike Coleman, and Rob Yelle (University of Oregon Research
565 and Advanced Computing Services) for their help with use of UO's computing cluster—Talapas.

566 Finally, we thank the Center for High Performance Computing at the University of Utah for
567 resources and support.

568

569 **References**

570 Andrews S. 2010. FASTQC. A quality control tool for high throughput sequence data. Available
571 from: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>

572 Antón SC, Potts R, Aiello LC. 2014. Evolution of early *Homo*: An integrated biological
573 perspective. *Science* 345:1236828.

574 Barratt CD, Lester JD, Gratton P, Onstein RE, Kalan AK, McCarthy MS, Bocksberger G, White
575 LC, Vigilant L, Dieguez P, et al. 2020. Late Quaternary habitat suitability models for
576 chimpanzees (*Pan troglodytes*) since the Last Interglacial (120,000 BP). *bioRxiv*
577 [Internet]. Available from:
578 <http://biorxiv.org/content/early/2020/05/25/2020.05.15.066662>

579 Barrett RDH, Hoekstra HE. 2011. Molecular spandrels: tests of adaptation at the genetic level.
580 *Nature Reviews Genetics* 12:767–780.

581 Beaune D, Hohmann G, Serckx A, Sakamaki T, Narat V, Fruth B. 2017. How bonobo
582 communities deal with tannin rich fruits: Re-ingestion and other feeding processes.
583 *Behavioural Processes* 142:131–137.

584 Behringer V, Deschner T, Deimel C, Stevens JMG, Hohmann G. 2014. Age-related changes in
585 urinary testosterone levels suggest differences in puberty onset and divergent life history
586 strategies in bonobos and chimpanzees. *Hormones and Behavior* 66:525–533.

587 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful
588 approach to multiple testing. *J R Stat Soc Series B Stat Methodol* 57:289–300.

589 Bergey CM, Lopez M, Harrison GF, Patin E, Cohen JA, Quintana-Murci L, Barreiro LB, Perry
590 GH. 2018. Polygenic adaptation and convergent evolution on growth and cardiac genetic
591 pathways in African and Asian rainforest hunter-gatherers. *Proc Natl Acad Sci USA*
592 115:E11256.

593 Besenbacher S, Hvilsom C, Marques-Bonet T, Mailund T, Schierup MH. 2019. Direct estimation
594 of mutations in great apes reconciles phylogenetic dating. *Nature Ecology & Evolution*
595 3:286–292.

596 Bradley BJ. 2008. Reconstructing phylogenies and phenotypes: a molecular view of human
597 evolution. *Journal of Anatomy* 212:337–353.

598 Broad Institute. 2018. Picard Tools. Available from: <http://broadinstitute.github.io/picard/>

- 599 Cagan A, Theunert C, Laayouni H, Santpere G, Pybus M, Casals F, Prüfer K, Navarro A,
600 Marques-Bonet T, Bertranpetit J, et al. 2016. Natural selection in the great apes. *Mol Biol*
601 *Evol* 33:3268–3283.
- 602 Calpena E, Hervieu A, Kaserer T, Swagemakers SMA, Goos JAC, Popoola O, Ortiz-Ruiz MJ,
603 Barbaro-Dieber T, Bownass L, Brilstra EH, et al. 2019. De Novo Missense Substitutions
604 in the Gene Encoding CDK8, a Regulator of the Mediator Complex, Cause a Syndromic
605 Developmental Disorder. *The American Journal of Human Genetics* 104:709–720.
- 606 Charlesworth B, Morgan MT, Charlesworth D. 1993. The effect of deleterious mutations on
607 neutral molecular variation. *Genetics* 134:1289–1303.
- 608 Compton AA, Malik HS, Emerman M. 2013. Host gene evolution traces the evolutionary history
609 of ancient primate lentiviruses. *Philosophical Transactions of the Royal Society B:*
610 *Biological Sciences* 368:20120496.
- 611 Ewels P, Magnusson M, Lundin S, Käller M. 2016. MultiQC: summarize analysis results for
612 multiple tools and samples in a single report. *Bioinformatics* 32:3047–3048.
- 613 Furuichi T. 2011. Female contributions to the peaceful nature of bonobo society. *Ev Anth*
614 20:131–142.
- 615 Gao F, Bailes E, Robertson DL, Chen Y, Rodenburg CM, Michael SF, Cummins LB, Arthur LO,
616 Peeters M, Shaw GM, et al. 1999. Origin of HIV-1 in the chimpanzee *Pan troglodytes*
617 troglodytes. *Nature* 397:436–441.
- 618 Goodall J. 1986. The chimpanzees of Gombe: Patterns of behavior. Cambridge, MA: Belknap
619 Press
- 620 Gossmann TI, Keightley PD, Eyre-Walker A. 2012. The Effect of Variation in the Effective
621 Population Size on the Rate of Adaptive Molecular Evolution in Eukaryotes. *Genome*
622 *Biology and Evolution* 4:658–667.
- 623 Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, Valieris R, Köster J,
624 The Bioconda Team. 2018. Bioconda: sustainable and comprehensive software
625 distribution for the life sciences. *Nature Methods* 15:475–476.
- 626 Han S, Andrés AM, Marques-Bonet T, Kuhlwilm M. 2019. Genetic variation in *Pan* species is
627 shaped by demographic history and harbors lineage-specific functions. *Genome Biology*
628 *and Evolution* 11:1178–1191.
- 629 Henn BM, Cavalli-Sforza LL, Feldman MW. 2012. The great human expansion. *Proc Natl Acad*
630 *Sci USA* 109:17758.
- 631 Hermisson J, Pennings PS. 2005. Soft sweeps: Molecular population genetics of adaptation from
632 standing genetic variation. *Genetics* 169:2335–2352.

- 633 Hermisson J, Pennings PS. 2017. Soft sweeps and beyond: understanding the patterns and
634 probabilities of selection footprints under rapid adaptation. *Methods Ecol Evol* 8:700–
635 716.
- 636 Hernandez RD, Kelley JL, Elyashiv E, Melton SC, Auton A, McVean G, Project 1000 Genomes,
637 Sella G, Przeworski M. 2011. Classic selective sweeps were rare in recent human
638 evolution. *Science* 331:920–924.
- 639 Hey J. 2010. The divergence of chimpanzee species and subspecies as revealed in
640 multipopulation isolation-with-migration analyses. *Mol Biol Evol* 27:921–933.
- 641 Hohmann G, Fruth B. 2003. Culture in bonobos? Between \square species and within \square species variation
642 in behavior. *Curr Anthropol* 44:563–571.
- 643 Huang DW, Sherman BT, Lempicki RA. 2008. Bioinformatics enrichment tools: paths toward
644 the comprehensive functional analysis of large gene lists. *Nucleic Acids Research* 37:1–
645 13.
- 646 Huang DW, Sherman BT, Lempicki RA. 2009. Systematic and integrative analysis of large gene
647 lists using DAVID bioinformatics resources. *Nature Protocols* 4:44–57.
- 648 Kalan AK, Kulik L, Arandjelovic M, Boesch C, Haas F, Dieguez P, Barratt CD, Abwe EE,
649 Agbor A, Angedakin S, et al. 2020. Environmental variability supports chimpanzee
650 behavioural diversity. *Nature Communications* 11:4451.
- 651 Kano T. 1992. The last ape: Pygmy chimpanzee behavior and ecology. Stanford: Stanford
652 University Press
- 653 Kern AD, Schrider DR. 2016. Discoal: flexible coalescent simulations with selection.
654 *Bioinformatics* 32:3839–3841.
- 655 Kern AD, Schrider DR. 2018. diploS/HIC: An updated approach to classifying selective sweeps.
656 *G3: Genes, Genomes, Genetics* 8:1959–1970.
- 657 Köster J, Rahmann S. 2012. Snakemake—a scalable bioinformatics workflow engine.
658 *Bioinformatics* 28:2520–2522.
- 659 Kronenberg ZN, Fiddes IT, Gordon D, Murali S, Cantsilieris S, Meyerson OS, Underwood JG,
660 Nelson BJ, Chaisson MJP, Dougherty ML, et al. 2018. High-resolution comparative
661 analysis of great ape genomes. *Science* [Internet] 360. Available from:
662 <https://science.sciencemag.org/content/360/6393/ear6343>
- 663 Kuhlwilm M, Han S, Sousa VC, Excoffier L, Marques-Bonet T. 2019. Ancient admixture from
664 an extinct ape lineage into bonobos. *Nat Ecol Evol* 3:957–965.
- 665 Lalouette A, Guénet J-L, Vríz S. 1998. Hotfoot Mouse Mutations Affect the $\delta 2$ Glutamate
666 Receptor Gene and Are Allelic to Lurcher. *Genomics* 50:9–13.

- 667 Langergraber KE, Prüfer K, Rowney C, Boesch C, Crockford C, Fawcett K, Inoue E, Inoue-
668 Muruyama M, Mitani JC, Muller MN, et al. 2012. Generation times in wild chimpanzees
669 and gorillas suggest earlier divergence times in great ape and human evolution. *Proc Natl*
670 *Acad Sci USA* 109:15716.
- 671 Lex A, Gehlenborg N, Strobel H, Vuillemot R, Pfister H. 2014. UpSet: Visualization of
672 Intersecting Sets. *IEEE Transactions on Visualization and Computer Graphics* 20:1983–
673 1992.
- 674 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping
675 and population genetical parameter estimation from sequencing data. *Bioinformatics*
676 27:2987–2993.
- 677 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
678 *arXiv:1303.3997*.
- 679 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
680 1000 Genome Project Data Processing Subgroup. 2009. The sequence alignment/map
681 format and SAMtools. *Bioinformatics (Oxford, England)* 25:2078–2079.
- 682 de Manuel M, Kuhlwilm M, Frandsen P, Sousa VC, Desai T, Prado-Martinez J, Hernandez-
683 Rodriguez J, Dupanloup I, Lao O, Hallast P, et al. 2016. Chimpanzee genomic diversity
684 reveals ancient admixture with bonobos. *Science* 354:477–481.
- 685 Maynard Smith J, Haigh J. 1974. The hitch-hiking effect of a favourable gene. *Genetics*
686 *Research* 23:23–35.
- 687 Messer PW, Petrov DA. 2013. Population genomics of rapid adaptation by soft selective sweeps.
688 *Trends Ecol Evol* 28:659–669.
- 689 Nakajima T, Ohtani H, Satta Y, Uno Y, Akari H, Ishida T, Kimura A. 2008. Natural selection in
690 the TLR-related genes in the course of primate evolution. *Immunogenetics* 60:727–735.
- 691 Nakano Y, Yamamoto K, Ueda MT, Soper A, Konno Y, Kimura I, Uriu K, Kumata R, Aso H,
692 Misawa N, et al. 2020. A role for gorilla APOBEC3G in shaping lentivirus evolution
693 including transmission to humans. *PLOS Pathogens* 16:e1008812.
- 694 Nam K, Munch K, Mailund T, Nater A, Greminger MP, Krützen M, Marquès-Bonet T, Schierup
695 MH. 2017. Evidence that the rate of strong selective sweeps increases with population
696 size in the great apes. *PNAS* 114:1613–1618.
- 697 Narasimhan VM, Rahbari R, Scally A, Wuster A, Mason D, Xue Y, Wright J, Trembath RC,
698 Maher ER, Heel DA van, et al. 2017. Estimating the human mutation rate from
699 autozygous segments reveals population differences in human mutational processes. *Nat*
700 *Commun* 8:1–7.
- 701 Nishida T. 2011. Chimpanzees of the lakeshore: Natural history and culture at Mahale.
702 Cambridge: Cambridge University Press

- 703 Nye J, Mondal M, Bertranpetit J, Laayouni H. 2020. A fully integrated machine learning scan of
704 selection in the chimpanzee genome. *NAR Genomics and Bioinformatics* [Internet] 2.
705 Available from: <https://doi.org/10.1093/nargab/lqaa061>
- 706 Oleksyk TK, Smith MW, O'Brien SJ. 2010. Genome-wide scans for footprints of natural
707 selection. *Philosophical Transactions of the Royal Society B: Biological Sciences*
708 365:185–205.
- 709 Osborne MJ, Volpon L, Kornblatt JA, Culjkovic-Kraljacic B, Baguet A, Borden KLB. 2013.
710 eIF4E3 acts as a tumor suppressor by utilizing an atypical mode of methyl-7-guanosine
711 cap recognition. *Proc Natl Acad Sci USA* 110:3877.
- 712 Pedersen BS, Quinlan AR. 2018. Mosdepth: quick coverage calculation for genomes and
713 exomes. *Bioinformatics* 34:867–868.
- 714 Pennings PS, Hermisson J. 2006a. Soft sweeps II—Molecular population genetics of adaptation
715 from recurrent mutation or migration. *Mol Biol Evol* 23:1076–1084.
- 716 Pennings PS, Hermisson J. 2006b. Soft sweeps III: The signature of positive selection from
717 recurrent mutation. *PLOS Genetics* 2:e186.
- 718 Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Van der Auwera GA, Kling
719 DE, Gauthier LD, Levy-Moonshine A, Roazen D, et al. 2018. Scaling accurate genetic
720 variant discovery to tens of thousands of samples. *bioRxiv*:201178.
- 721 Potts R. 1998. Variability selection in hominid evolution. *Ev Anth* 7:81–96.
- 722 Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, Veeramah KR,
723 Woerner AE, O'Connor TD, Santpere G, et al. 2013. Great ape genetic diversity and
724 population history. *Nature* 499:471–475.
- 725 Pritchard JK, Pickrell JK, Coop G. 2010. The genetics of human adaptation: Hard sweeps, soft
726 sweeps, and polygenic adaptation. *Curr Biol* 20:R208–R215.
- 727 Prüfer K, Munch K, Hellmann I, Akagi K, Miller JR, Walenz B, Koren S, Sutton G, Kodira C,
728 Winer R, et al. 2012. The bonobo genome compared with the chimpanzee and human
729 genomes. *Nature* 486:527–531.
- 730 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
731 features. *Bioinformatics* 26:841–842.
- 732 R Core Team. 2020. R: A Language and Environment for Statistical Computing. Vienna,
733 Austria: R Foundation for Statistical Computing Available from: [https://www.R-](https://www.R-project.org/)
734 [project.org/](https://www.R-project.org/)
- 735 Ralph P, Coop G. 2010. Parallel adaptation: One or many waves of advance of an advantageous
736 allele? *Genetics* 186:647–668.

- 737 Rilling JK, Scholz J, Preuss TM, Glasser MF, Errangi BK, Behrens TE. 2012. Differences
738 between chimpanzees and bonobos in neural systems supporting social cognition. *Social*
739 *Cognitive and Affective Neuroscience* 7:369–379.
- 740 Sakamaki T, Maloueki U, Bakaa B, Bongoli L, Kasalevo P, Terada S, Furuichi T. 2016.
741 Mammals consumed by bonobos (*Pan paniscus*): new data from the Iyondji forest,
742 Tshuapa, Democratic Republic of the Congo. *Primates* 57:295–301.
- 743 Sarich VM, Wilson AC. 1967. Immunological time scale for hominid evolution. *Science*
744 158:1200.
- 745 Sawyer SL, Wu LI, Emerman M, Malik HS. 2005. Positive selection of primate *TRIM5a*
746 identifies a critical species-specific retroviral restriction domain. *Proc Natl Acad Sci U S*
747 *A* 102:2832.
- 748 Scally A, Dutheil JY, Hillier LW, Jordan GE, Goodhead I, Herrero J, Hobolth A, Lappalainen T,
749 Mailund T, Marques-Bonet T, et al. 2012. Insights into hominid evolution from the
750 gorilla genome sequence. *Nature* 483:169–175.
- 751 Schiffels S, Durbin R. 2014. Inferring human population size and separation history from
752 multiple genome sequences. *Nature Genetics* 46:919–925.
- 753 Schmidt JM, Manuel M de, Marques-Bonet T, Castellano S, Andrés AM. 2019. The impact of
754 genetic adaptation on chimpanzee subspecies differentiation. *PLOS Genetics*
755 15:e1008485.
- 756 Schrider DR. 2020. Background selection does not mimic the patterns of genetic diversity
757 produced by selective sweeps. *Genetics* 216:499.
- 758 Schrider DR, Kern AD. 2017. Soft sweeps are the dominant mode of adaptation in the human
759 genome. *Mol Biol Evol* 34:1863–1877.
- 760 Schrider DR, Mendes FK, Hahn MW, Kern AD. 2015. Soft shoulders ahead: Spurious signatures
761 of soft and partial selective sweeps result from linked hard sweeps. *Genetics* 200:267–
762 284.
- 763 Staes N, Smaers JB, Kunkle AE, Hopkins WD, Bradley BJ, Sherwood CC. 2019. Evolutionary
764 divergence of neuroanatomical organization and related genes in chimpanzees and
765 bonobos. *Cortex* 118:154–164.
- 766 Stevison LS, Woerner AE, Kidd JM, Kelley JL, Veeramah KR, McManus KF, Great Ape
767 Genome Project, Bustamante CD, Hammer MF, Wall JD. 2015. The time scale of
768 recombination rate evolution in great apes. *Mol Biol Evol* 33:928–945.
- 769 Stimpson CD, Barger N, Tagliabue JP, Gendron-Fitzpatrick A, Hof PR, Hopkins WD,
770 Sherwood CC. 2016. Differential serotonergic innervation of the amygdala in bonobos
771 and chimpanzees. *Social Cognitive and Affective Neuroscience* 11:413–422.

- 772 Stumpf RM. 2011. Chimpanzees and bonobos: Inter- and intraspecies diversity. In: Campbell CJ,
773 Fuentes A, MacKinnon KC, Bearder SK, Stumpf RM, editors. *Primates in perspective*.
774 New York: Oxford University Press. p. 340–356.
- 775 Sun P-H, Ye L, Mason MD, Jiang WG. 2012. Protein Tyrosine Phosphatase μ (PTP μ or
776 PTPRM), a Negative Regulator of Proliferation and Invasion of Breast Cancer Cells, Is
777 Associated with Disease Prognosis. *PLOS ONE* 7:e50183.
- 778 Susman RL ed. 1984. *The pygmy chimpanzee: Evolutionary biology and behavior*. New York:
779 Springer
- 780 Takemoto H, Kawamoto Y, Higuchi S, Makinose E, Hart JA, Hart TB, Sakamaki T, Tokuyama
781 N, Reinartz GE, Guislain P, et al. 2017. The mitochondrial ancestor of bonobos and the
782 origin of their major haplogroups. *PLOS ONE* 12:e0174851.
- 783 Turley K, Frost SR. 2014. The appositional articular morphology of the talo-crural joint: The
784 influence of substrate use on joint shape. *Anat Rec* 297:618–629.
- 785 Van Heuverswyn F, Li Y, Neel C, Bailes E, Keele BF, Liu W, Loul S, Butel C, Liegeois F,
786 Bienvenue Y, et al. 2006. SIV infection in wild gorillas. *Nature* 444:164–164.
- 787 van der Lee R, Wiel L, van Dam TJP, Huynen MA. 2017. Genome-scale detection of positive
788 selection in nine primates predicts human-virus evolutionary conflicts. *Nucleic Acids*
789 *Research* 45:10634–10648.
- 790 Wakefield ML, Hickmott AJ, Brand CM, Takaoka IY, Meador LM, Waller MT, White FJ. 2019.
791 New observations of meat eating and sharing in wild bonobos (*Pan paniscus*) at Iyema,
792 Lomako Forest Reserve, Democratic Republic of the Congo. *Fol Primatol* 90:179–189.
- 793 Webster TH, Couse M, Grande BM, Karlins E, Phung TN, Richmond PA, Whitford W, Wilson
794 MA. 2019. Identifying, understanding, and correcting technical artifacts on the sex
795 chromosomes in next-generation sequencing data. *Gigascience* [Internet] 8. Available
796 from: <https://academic.oup.com/gigascience/article/8/7/giz074/5530326>
- 797 Webster TH, Wilson Sayres MA. 2016. Genomic signatures of sex-biased demography: progress
798 and prospects. *Current Opinion in Genetics & Development* 41:62–71.
- 799 Wegmann D, Excoffier L. 2010. Bayesian inference of the demographic history of chimpanzees.
800 *Mol Biol Evol* 27:1425–1435.
- 801 Wetzl KS, Yi Y, Yadav A, Bauer AM, Bello EA, Romero DC, Bibollet-Ruche F, Hahn BH,
802 Paiardini M, Silvestri G, et al. 2018. Loss of CXCR6 coreceptor usage characterizes
803 pathogenic lentiviruses. *PLOS Pathogens* 14:e1007003.
- 804 White FJ. 1996. *Pan paniscus* 1973 to 1996: Twenty-three years of field research. *Ev Anth* 5:11–
805 17.

- 806 Wickham H. 2016. *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag
807 Available from: <https://ggplot2.tidyverse.org>
- 808 Wilson ML, Boesch C, Fruth B, Furuichi T, Gilby IC, Hashimoto C, Hobaiter CL, Hohmann G,
809 Itoh N, Koops K, et al. 2014. Lethal aggression in *Pan* is better explained by adaptive
810 strategies than human impacts. *Nature* 513:414–417.
- 811 Wrangham RW. 1986. Ecology and social relationships in two species of chimpanzee. In:
812 Rubenstein DI, Wrangham RW, editors. *Ecological aspects of social evolution: Birds and*
813 *mammals*. Princeton, NJ: Princeton University Press. p. 352–378.
- 814 Yang J, Jin Z-B, Chen J, Huang X-F, Li X-M, Liang Y-B, Mao J-Y, Chen X, Zheng Z, Bakshi
815 A, et al. 2017. Genetic signatures of high-altitude adaptation in Tibetans. *Proc Natl Acad*
816 *Sci USA* 114:4189.
- 817

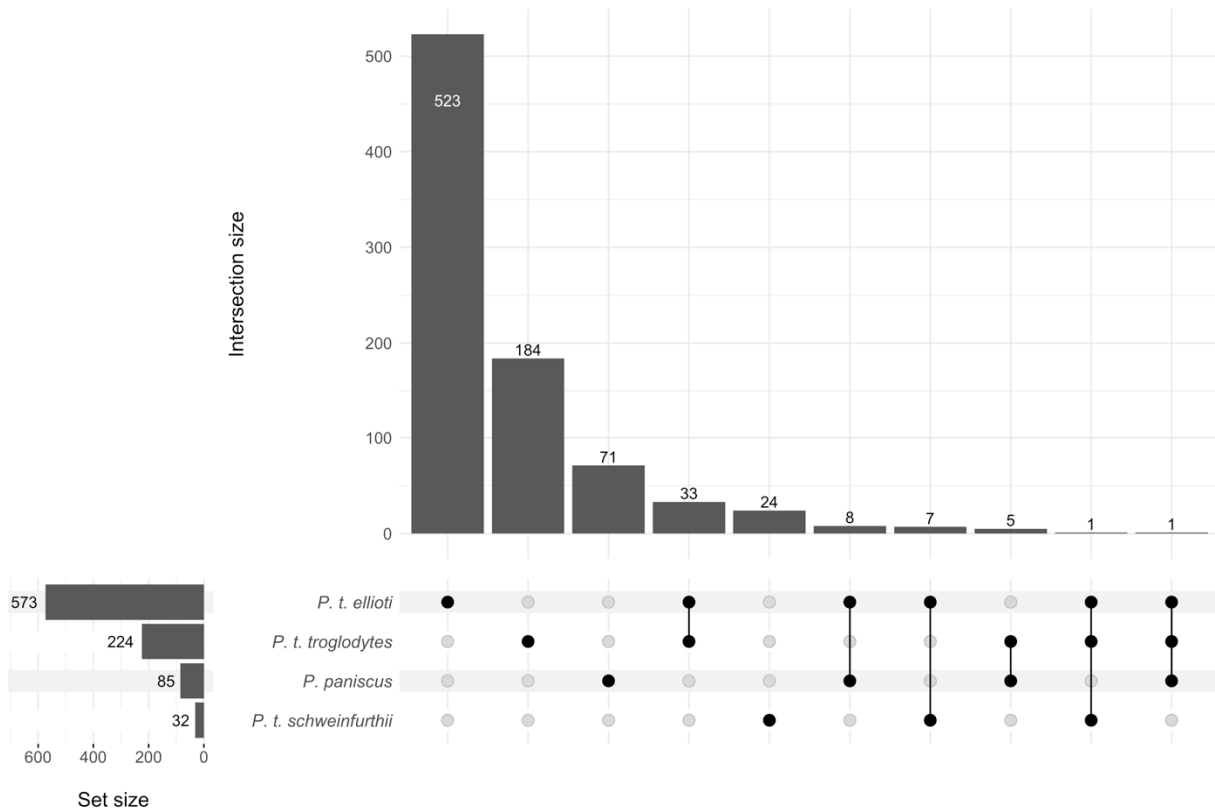
818 Table 1. Selective sweep summary per population.

Lineage	Number / Percent of Windows per Class Type						Number and Percent of Sweep Type		
	Hard	Linked- hard	Neutral	Linked- soft	Soft	Total	Hard	Soft	Total
<i>P. paniscus</i>	85 (0.4%)	1,576 (6.5%)	7,488 (30.8%)	13,168 (54.1%)	2,002 (8.2%)	24,319	81 (4.9%)	1,585 (95.1%)	1,666
<i>P. t. ellioti</i>	573 (2.4%)	6,358 (26.1%)	1,389 (5.7%)	14,498 (59.6%)	1,505 (6.2%)	24,323	488 (26.9%)	1,323 (73.1%)	1,811
<i>P. t. schweinfurthii</i>	32 (0.1%)	696 (2.9%)	1,835 (7.5%)	20,179 (83.0%)	1,581 (6.5%)	24,323	32 (2.3%)	1,376 (97.7%)	1,408
<i>P. t. troglodytes</i>	224 (0.9%)	1,746 (7.2%)	5,483 (22.5%)	15,121 (62.2%)	1,749 (7.2%)	24,323	184 (10.6%)	1,557 (89.4%)	1,741

819

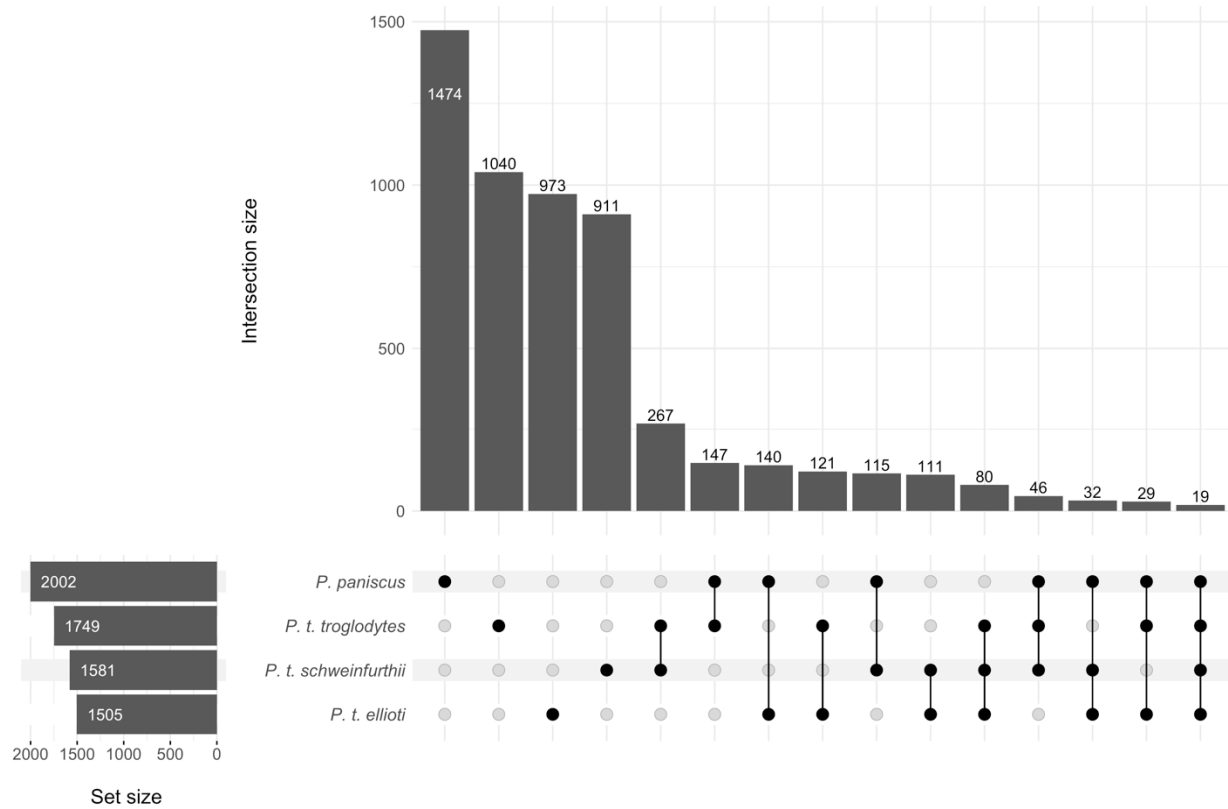
820

821 Figure 1. Unique and shared hard sweep windows. The frequency of windows shared by two or
822 more lineages should be considered relative to the total possible number of shared windows (i.e.,
823 the set size of the lineage with the smallest set size).



824

825 Figure 2. Unique and shared soft sweep windows. The frequency of windows shared by two or
826 more lineages should be considered relative to the total possible number of shared windows (i.e.,
827 the set size of the lineage with the smallest set size).



828

829

830 Supplements.

- 831 • Main Supplemental File: Figures S1 - S3, Tables S1-S4.
- 832 • File S1. Sample information. (File name: File_S1_sample_information.xlsx)
- 833 • File S2. Genome accessibility information. (File name:
834 File_S2_genome_accessibility.xlsx)
- 835 • File S3. Discoal parameter information. (File name:
836 File_S3_discoal_input_summary.xlsx)
- 837 • File S4. Classifier trial information. (File name:
838 File_S4_diploshic_classifier_summary.xlsx)
- 839 • File S5. Unmasked SNV count/fraction per window for VCF feature vectors. (File name:
840 File_S5_fvec_vcf_unmaskedsnpcount_unmaskedfrac_summary)
- 841 • File S6. Number of hard and soft sweep windows using higher probability thresholds.
842 (File name: File_S6_sweeptype_probability_cutoff_summary.xlsx)
- 843 • File S7. Genes included in sweep analysis (File name: File_S7_genes_to_include.xlsx)
- 844 • File S8. Sweep information. (File name: File_S8_selective_sweep_summary.xlsx)
- 845 • File S9. List of genes in hard and soft sweeps. (File name: File_S9_gene_lists.xlsx)
- 846 • File S10. Gene ontology analysis. (File name: File_S10_gene_ontology.xlsx)