

Evidence that viruses, particularly SIV, drove genetic adaptation in natural populations of eastern chimpanzees

Joshua M. Schmidt^{1,2}, Marc de Manuel³, Tomas Marques-Bonet^{3,4,5}, Sergi Castellano^{2,6,7} and Aida M. Andrés^{1,2}.

¹ Present address: UCL Genetics Institute, Department of Genetics, Evolution and Environment, University College London, London, UK.

² Max Planck Institute for Evolutionary Anthropology, Department of Evolutionary Genetics, Leipzig, Germany.

³ Institut de Biologia Evolutiva (Consejo Superior de Investigaciones Científicas–Universitat Pompeu Fabra), Barcelona, Spain.

⁴ National Centre for Genomic Analysis–Centre for Genomic Regulation, Barcelona Institute of Science and Technology, Barcelona, Spain.

⁵ Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain.

⁶ Present address: Genetics and Genomic Medicine Programme, Great Ormond Street Institute of Child Health, University College London (UCL), London, UK

⁷ Present address: UCL Genomics, London, UK

Corresponding author: Aida M. Andrés

1 **Abstract**

2

3 All four subspecies of chimpanzees are endangered. Differing in their demographic histories
4 and geographical ranges within sub-Saharan Africa, they have likely adapted to different
5 environmental factors. We show that highly differentiated SNPs in eastern chimpanzees are
6 uniquely enriched in genic sites in a way that is expected under recent adaptation. These sites
7 are enriched for genes that differentiate the immune response to infection by simian
8 immunodeficiency virus (SIV) in natural vs. non-natural host species. Conversely, central
9 chimpanzees exhibit selective sweeps at the cytokine receptors *CCR3*, *CCR9* and *CXCR6* –
10 paralogs of *CCR5* and *CXCR4*, the two major receptors utilized by HIV to enter human cells.
11 Thus, we infer that SIV may be eliciting distinctive adaptive responses in different chimpanzee
12 subspecies. Since central chimpanzee SIV is the source of the global HIV/AIDS pandemic,
13 understanding the mechanisms that limit pathogenicity of SIV in chimpanzees can broaden our
14 understanding of HIV infection in humans.

15 Chimpanzees (*Pan troglodytes*) are, alongside bonobos, human's closest living relatives – the
16 *Pan* and *Homo* lineages having diverged ~6Myr ago (Prado-Martinez et al. 2013). With a
17 genetic divergence of only ~1% (Consortium et al. 2005), *Pan* and *Homo* also share large
18 aspects of their physiology and behaviour, including susceptibility to some pathogens.
19 Studying chimpanzees can teach us about our species by putting recent human evolution in its
20 evolutionary context i.e. the mode and tempo of adaptation and the pressures driving it.

21
22 Selection imposed by pathogens has greatly shaped the long-term history of genetic adaptation
23 in the great apes, including chimpanzees and humans (Cagan et al. 2016, Enard et al. 2016).
24 The interest in recent human evolution (Sabeti et al. 2002, Voight et al. 2006, Sabeti et al. 2007,
25 Yi et al. 2010, Racimo 2016) means that we now also have good catalogues of the main targets
26 of local adaptation in many non-African human populations. Nevertheless, analyses of
27 genome-wide patterns of diversity suggest that adaptation via hard selective sweeps has had a
28 limited role in shaping human genomes. Complete selective sweeps involving non-
29 synonymous substitutions appear to have been rare (Hernandez et al. 2011) – but perhaps still
30 important (Enard, Messer, and Petrov 2014). Further, local adaptation has had little effect on
31 the patterns of population differentiation (Coop et al. 2009), unless inferences are boosted with
32 ancient DNA (Key et al. 2016). The focus on humans biases our view on the influence of
33 genetic adaptations in natural populations of primates, and we do not know whether positive
34 selection plays a similarly limited role in shaping other primate genomes. We aim to address
35 this limitation by exploring the recent adaptive history of chimpanzees.

36
37 There are four subspecies of chimpanzees, with common names reflecting their location in
38 western and central sub-Saharan Africa: eastern, central, Nigeria-Cameroon and western
39 (Figure 1). Each chimpanzee subspecies is currently endangered, with western chimpanzees

40 critically so (Humble et al. 2016). Subspecies are clearly differentiated, with divergence times
41 ranging from 450 kya to 100 kya, and estimated long-term N_e from 8,000 to 30,000 reflected
42 in varying levels of genetic diversity (Figure 1). There is a wide range of ecological variation
43 across the chimpanzee range, which spans over 5,000 km in sub-Saharan Africa and includes
44 deep forest and savanna-woodland mosaics. Pathogen incidence can also vary between these
45 groups, as seen recently with the lethal outbreaks of Anthrax (Leendertz et al. 2004) and Ebola
46 (Formenty et al. 1999), or the Simian immunodeficiency virus (SIV). SIV, the precursor of the
47 human immunodeficiency virus type 1 (HIV-1) virus that is responsible for the human AIDS
48 pandemic (Keele et al. 2006), is thought to be largely non-lethal to chimpanzees, although some
49 eastern chimpanzees can develop immunodeficiency, see (Rudicell et al. 2010, Keele et al.
50 2009)). Its prevalence is not uniform across the subspecies, and there is no evidence for
51 infections in western or Nigeria-Cameroon chimpanzees (Locatelli et al. 2016) but multiple
52 infections have been detected in communities of central and eastern chimpanzees (Locatelli et
53 al. 2016, Heuverswyn et al. 2007). Given the separate history and differential environment of
54 each subspecies, and the fact that each subspecies is an independent conservation unit, it is
55 crucial that we identify not only the genetic adaptations shared by all chimpanzees (Cagan et
56 al. 2016), but also the genetic differences conferring differential adaptation to each subspecies.
57
58 To do this, we investigated the signatures of recent genetic adaptation in the genomes of the
59 four subspecies. We show that only eastern chimpanzees have a clear genome-wide signal of
60 recent, local positive selection. This adaptation is potentially due to selection on immunity
61 related genes, with evidence consistent with selection imposed by viruses in general, and SIV
62 in particular. In contrast, putative adaptation to SIV in central chimpanzees seems mediated by
63 adaptation in a suite of cell-entry receptors, results which are suggestive of divergent paths of
64 adaptation to a common pathogen.

65

66 **Results**

67 **Genic enrichment in the distribution of derived allele frequency differences**

68 To investigate the influence of recent genetic adaptation in chimpanzee subspecies we
69 compared population differentiation at putatively functional sites (genic sites, defined as +-
70 2kb from protein-coding genes) to differentiation at non-functional sites (here non-genic).
71 Natural selection can only act on functional sites (or affect neutral sites tightly linked to
72 functional sites), so differences between functional and non-functional sites can be ascribed to
73 natural selection. After binning every SNP by its signed difference in derived allele frequency
74 between a pair of subspecies (δ), for each bin of δ we calculated the genic enrichment, defined
75 as the ratio of genic SNPs vs. all SNPs for each bin of δ , normalized by the global genic SNP
76 ratio (Coop et al. 2009, Hernandez et al. 2011, Key et al. 2016). This strategy has been deployed
77 in the study of human local adaptation (Coop et al. 2009, Hernandez et al. 2011, Key et al.
78 2016), and by not relying on the patterns of linked variation it is not strongly restricted to
79 particular modes of selection. The genic enrichment is greatest for SNPs with the largest δ ,
80 with the tail bins of δ exhibiting significantly greater genic enrichments than any other bin
81 (Figure 2). While not every genic SNP is in this bin due to positive selection, we expect these
82 SNPs, which show the largest frequency differences between subspecies in the genome, to be
83 strongly enriched in targets of positive selection that rose fast in frequency in one of the two
84 subspecies (Coop et al. 2009, Hernandez et al. 2011, Key et al. 2016).

85

86 The genic enrichment in the tails of δ is typically roughly symmetric (Figure 2a, symmetry
87 defined as overlapping δ tail bin genic enrichment 95% CIs), although the number of tail SNPs
88 and the magnitude of genic enrichment across subspecies pairs varies in accordance with their
89 N_e and divergence times (Figure 2 and Supplementary file 1). Calculated against western
90 chimpanzees, the subspecies with the lowest long-term N_e (de Manuel et al. 2016, Prado-

91 Martinez et al. 2013), the δ tail genic enrichment is the least, ranging from 1.05 to 1.09 (Figure
92 2a). A greater tail genic enrichment, 1.21 to 1.27, is seen for δ calculated using Nigeria-
93 Cameroon, the species with the second lowest long-term N_e (Figure 2a). This is comparable to
94 the magnitude of the genic enrichment in the tails of δ between human populations (Appendix
95 1, Appendix 1-figure 1; see (Coop et al. 2009, Hernandez et al. 2011, Key et al. 2016); the
96 genic enrichment across each bin of δ also resembles those observed in human populations
97 (Appendix 1, Appendix 1-figure 1; see (Coop et al. 2009, Hernandez et al. 2011, Key et al.
98 2016).

99
100 In marked contrast to these symmetric enrichments, we find a distinctive asymmetry between
101 the tail bin genic enrichments of central and eastern chimpanzees (Figure 2b). The central δ
102 tail exhibits a typical genic enrichment (1.25) but surprisingly, the eastern δ tail has a genic
103 enrichment (1.58) that is significantly greater than the central tail ($0.01 < P < 0.005$; weighted
104 200kb block jackknife, see Methods) and any other δ tail (all $P < 0.0001$; weighted 200kb block
105 jackknife).

106
107 The large confidence interval of the central chimpanzee δ tail genic enrichment is largely due
108 to the low number but high linkage of SNPs. For example, we found a highly unusual 200kb
109 genomic block on chromosome 3 that contains 70 highly differentiated alleles between central
110 and eastern chimpanzees, similarly distributed among the two tails (37 SNPs in the Central tail
111 and 33 SNPs in the Eastern tail). Concerned that this block could bias our results, we repeated
112 the enrichment analysis after excluding all SNPs contained within it. Removing this block
113 reduces the genic enrichment slightly in the Eastern tail (1.55) but substantially in the Central
114 tail (1.10) resulting in an even stronger asymmetry among the tails. Results are also robust to
115 the removal of the next largest block (with 27 SNPs in the two tails).

116

117

118 To directly quantify asymmetry of the eastern and central chimpanzee δ tail genic enrichments,
119 we tested if the \log_2 ratio of each pair of δ tail bin genic enrichments departs from zero.
120 Typically, the genic enrichment is greater for the subspecies with the higher long-term N_e .
121 However, the \log_2 ratios are similar and small (ranging from 0 and 0.055), except for eastern
122 vs. central, where it is 0.337 (95% CI, 0.153-0.521, 200kb weighted block jackknife). This is
123 six times larger than the highest ratio between other subspecies pairs (Figure 3; Supplementary
124 file 2, *p-value* δ western vs. central = 0.13 all other *p-values* ≤ 0.002 , *z-test*). The eastern vs.
125 central asymmetry in genic enrichment is thus a clear outlier (*p-value* $< 2.2e-16$, two-sided
126 Kolmogorov-Smirnov test).

127

128 **Background selection does not explain the δ tail asymmetry**

129

130 A certain level of δ tail bin genic enrichment (Figure 2) is, in principle, compatible with both
131 recent positive selection and background selection (BGS) (Coop et al. 2009), the latter because
132 linkage to sites under purifying selection reduces N_e locally in genic regions and increases the
133 effects of random genetic drift on neutral sites (Charlesworth, Morgan, and Charlesworth
134 1993). BGS can, for example, explain the δ tail bin genic enrichment in human populations,
135 suggesting that this pattern is not evidence for pervasive recent human adaptation (Coop et al.
136 2009, Hernandez et al. 2011, Key et al. 2016). To explore if BGS can explain our observations,
137 we used coalescent simulations (Methods) to estimate the expected reduction of neutral
138 diversity due to BGS, quantified as a genome-wide average B value (McVicker et al. 2009)
139 that best explains the genic enrichments across all bins of δ (the B value that minimizes the
140 summed square differences between observed and simulated enrichments across all pairwise δ
141 bins, Appendix 2). This B value is 0.888 – i.e. a reduction of diversity of ~ 11 per cent –

142 decreasing to 0.925 (weaker BGS) when excluding the δ tail bins, and increasing to 0.863
143 (stronger BGS) when fitting solely the twelve δ tail bins (Appendix 2, Supplementary file 4).
144 These values agree well with those inferred in humans using similar approaches (Hernandez et
145 al. 2011, Key et al. 2016). It is nonetheless clear that the B value of 0.863 that explains the δ
146 tail bin genic enrichments results in an extremely poor fit to the genic enrichments in all other
147 δ bins (Figure 2).

148

149 Previously, it was shown that BGS alone does not produce δ tail bin genic enrichment
150 asymmetries in comparisons of human populations (Key et al. 2016). We also find that BGS
151 does not result in significant eastern vs. central δ tail bin genic enrichment asymmetry.
152 Simulations show a slight asymmetry in the tail genic enrichment (Figs. 2B, 3) due to
153 differences in their demographic histories (Appendix 2, Supplementary file 5). Nevertheless,
154 no simulated value of B in the reasonable range identified above (0.925 – 0.863) results in a
155 tail genic enrichment \log_2 ratio that falls within the 95% CI of the observed ratio (Figure 3). In
156 contrast, the small (though statistically significant) asymmetries in other pairwise δ tail bin
157 genic enrichments are observed in simulations and thus fully explicable by demography and
158 BGS (Figure 3). Further, while we find one B value, $B = 0.850$, that results in a genic
159 enrichment that lies within the 95% CIs for both eastern and central chimpanzees, this B value
160 is a very poor fit to the genic enrichment in all other δ bins (Supplementary file 4) and cannot
161 explain the δ tail bin genic enrichment asymmetry (observed = 0.337, simulated B of 0.850 =
162 0.104).

163

164 Previous work has shown that background selection varies little among the great apes (Nam et
165 al. 2017). Theory suggests that the diversity reducing effect of BGS is independent of N_e , being
166 determined by the distribution of fitness effects (s), except for the narrow range of $N_e * s = 1$

167 (Nam et al., 2017), while previous work suggests that more than 80% of deleterious mutations
168 in chimpanzees have $N_e s \gg 1$ (Bataillon et al. 2015) Thus, the expectation is that the diversity
169 reducing effect of BGS should be the same across all four chimpanzee sub-species. Indeed, we
170 find comparable effects of background selection across subspecies: the relative reduction in
171 neutral variation linked to genes is comparable amongst chimpanzee subspecies (Appendix 3-
172 figure 1a), and neutral diversity has similar dependency on recombination rate and density of
173 functional features across subspecies (with the exception of western chimpanzees, Appendix
174 3-figure 1b). Further, using a population genetic statistical model (Corbett-Detig, Hartl, and
175 Sackton 2015) we estimate the same reduction in neutral diversity due to background selection
176 in each chimpanzee subspecies, at 11%, in the highest likelihood model (Appendix 3,
177 Supplementary file 6). Thus, despite their differing demographic histories (Figure 1), the
178 effects of BGS are very similar across each chimpanzee subspecies. This justifies using the
179 same average strength of BGS across subspecies above. Nevertheless, to explore if our
180 conclusions are robust to this assumption, we also modelled a greater strength of BGS in
181 eastern chimpanzees ($B = 0.825$, the value which best matches the eastern δ tail bin genic
182 enrichment) than in the other subspecies (B range 0.900-0.850). Stronger BGS in eastern
183 chimpanzees does not produce an eastern central δ tail bin asymmetry as large as that observed
184 in the genomes (log2 ratio range 0.120 – 0.146), further illustrating that BGS cannot explain
185 the greater tail genic enrichment in eastern chimpanzees (Figure 3-figure supplement 1).
186 Rather, this is most likely a signal of recent adaptation.

187

188 **Population-specific branch lengths with PBSnj**

189

190 Pairwise comparisons cannot determine which subspecies has changed. Direction, and
191 therefore biological meaning, to allele frequency difference can only be garnered by assuming
192 that derived alleles most often provide the basis for new adaptations. This approach is also

193 limited by the collapsing of the shared history of lineages. For example, in the Nigeria-
194 Cameroon vs. Eastern comparison, 22% of the SNPs in the Eastern δ tail are also in the Central
195 δ tail (for Nigeria-Cameroon vs. Central comparison), whereas only 3.5% (616 of 17,793) are
196 highly differentiated to both Nigeria and Central chimpanzees. Thus, δ summarises the allele
197 frequency change across several parts of the phylogeny, hampering the biological interpretation
198 of its tails.

199

200 To overcome this limitation, we developed a statistic that extends the widely used Population
201 Branch Statistic (PBS) (Yi et al. 2010). Briefly, large PBS values identify targets of positive
202 selection as SNPs with population-specific allele frequency differentiation, as these sites result
203 in unusually long branch lengths in pairwise F_{ST} -distance trees between three taxa. Small PBS
204 values are due to very short branches, for example due to purifying, shared balancing selection
205 or rare mutations. We extend this test to more than three taxa in the novel *PBSnj* statistic by
206 applying the Neighbor-Joining (NJ) algorithm on the matrix of the per-SNP pairwise F_{ST}
207 distances of the four subspecies (Methods, Appendix 4). This way, *PBSnj* allows us to jointly
208 compare the four subspecies and identify SNPs with very long branches (allele frequency
209 differentiation) in one subspecies only. Additional advantages of *PBSnj* are that it does not rely
210 on the specification of ancestral or derived states, and that the NJ algorithm does not require
211 specification of a phylogenetic tree describing the relationship amongst taxa (Appendix 4).

212

213 *PBSnj* allows us to determine within which lineage, eastern or central chimpanzees, allele
214 frequencies have changed to result in the asymmetric δ genic enrichment. Analogous to the δ
215 tail bins, we binned *PBSnj* scores and calculated the genic enrichment for each species *PBSnj*
216 tail (Figure 4A). The *PBSnj* eastern tail has significantly stronger genic enrichment than the
217 central tail (eastern: 1.36, central: 1.13, \log_2 ratio = 0.25, $p < 0.001$ estimated from weighted

218 200kb block jackknife, Figure 4B). This shows that the central vs. eastern asymmetry in the δ
219 tail bin genic enrichments (Figs. 2B, 3) is due to the drastic allele frequency rise of genic SNPs
220 in eastern chimpanzees since their divergence with central chimpanzees. Importantly, across
221 the range of B values (1.0 – 0.80), simulations show that eastern and central chimpanzee PBSnj
222 tail genic enrichments are expected to be equal (Figure 4B). In fact, BGS would need to be
223 much stronger in eastern chimpanzees than in central chimpanzees to produce the observed
224 levels of PBSnj tail genic enrichments. BGS with $B < 0.888$ would be required to produce the
225 genic enrichment exhibited in the eastern PBSnj tail, but $B = 0.888$ produces PBSnj tail genic
226 enrichments of equal or greater magnitude as those seen for central chimpanzees (and also
227 Nigeria-Cameroon and western, Appendix 5 and Appendix 6, and Supplementary file 9). Thus,
228 it is eastern chimpanzees that exhibit the greatest genic enrichment for highly differentiated
229 SNPs, an enrichment that (unlike in other subspecies) we cannot explain by background
230 selection alone. This suggests the greater enrichment in the PBSnj eastern tail is due to positive
231 selection, and by using the genomic blocks used to estimate the PBSnj tail Confidence Intervals
232 in Figure 4A, we estimate that an additional eight-19 population specific sweeps are sufficient
233 to explain this difference (Methods, Figure 4-figure supplement 1). Although this is a
234 conservative estimate, it shows that we do not require an unrealistically large number of
235 selective sweeps to explain the distinct pattern of eastern chimpanzees.

236

237 **Long-range LD and regulatory functions in the PBSnj eastern tail SNPs**

238

239 Further evidence that the PBSnj eastern tail genic enrichment is not due to background
240 selection would be provided by independent signatures such as the patterns of linkage
241 disequilibrium (LD) and the putative functional consequences of alleles. LD based tests of
242 positive selection are more robust to background selection than those based on population
243 differentiation (Enard, Messer, and Petrov 2014). We computed three haplotype-based

244 selection statistics that identify the signatures of positive selection within populations (*iHS*;
245 (Voight et al. 2006), *nSI* (Ferrer-Admetlla et al. 2014)) or between populations (*XP-EHH*
246 (Sabeti et al. 2007)). For each statistic, PBSnj eastern tail SNPs have a significantly higher
247 score than randomly sampled genic SNPs (mean *iHS* 0.69, mean *nSI* 0.94, and mean *XP-EHH*
248 mean 0.51, standardized for the genic background to have mean = 0 and sd = 1 for each statistic;
249 all $p < 0.0001$; re-sampling test; Supplementary file 10). Thus, SNPs specifically differentiated
250 in the eastern PBSnj tail have on average higher LD-based signatures of recent positive
251 selection than random genic SNPs, and also to a greater degree than all other subspecies' PBSnj
252 tails (Supplementary file 10) (two-sample *t*-tests, all *p-values* $\ll 0.0001$).

253

254 PBSnj tail genic SNPs are significantly enriched in exonic variants, but not in non-synonymous
255 (as compared with synonymous) ones, so less than 1% of PBSnj eastern tail SNPs result in
256 amino acid changes (observed = 0.84%; genic background = 0.18 %, $p < 0.001$, Supplementary
257 file 11 lists PBSnj eastern tail non-synonymous SNPs). Turning to regulatory changes, we used
258 regulomeDB (Boyle et al. 2012) to predict putative regulatory consequences of chimpanzee
259 SNPs from the sequence context and biochemical signatures of homologous human sites. The
260 PBSnj eastern tail genic SNPs are more likely to have strong evidence of regulatory function
261 (3.7% vs. 3.0%, permutation test $p = 0.012$) and less likely to have no ascribed regulatory
262 function (52.3% vs. 56.0%, permutation test $p = 0.0001$) than randomly sampled genic SNPs,
263 Supplementary file 12. In contrast, PBSnj central tail SNPs show no difference to the genic
264 background for either category (Supplementary file 12; Nigeria-Cameroon and western also
265 exhibit weaker but significant enrichments). Interestingly, PBSnj eastern tail SNPs do not differ
266 in functional constraint (as measured by *phastCons* scores (Siepel et al. 2005, see Methods)
267 from random genic SNPs (Supplementary file 13). This suggests that while likely enriched in
268 regulatory functions, these sites are not under particularly strong long-term constraint, perhaps

269 because they do not affect functions that have been tightly conserved over long evolutionary
270 times.

271

272 **Biological functions of the PBSnj eastern tail SNPs**

273 To understand the biological mechanisms and putative selective factors driving the recent
274 adaptations in eastern chimpanzees, we investigated the genes containing the genic SNPs in
275 the PBSnj eastern tail (hereafter PBSnj eastern genes). Two Gene Ontology (GO) categories
276 (Ashburner et al. 2000, The Gene Ontology 2017) are significantly enriched ($p < 0.05$, False
277 Discovery Rate (FDR) < 0.05 ; GOWINDA; Supplementary files 15-18). The top category is
278 “cytoplasmic mRNA processing body assembly”, and three of the five PBSnj eastern genes in
279 this category (*DDX6* (Ayache et al. 2015), *ATXN2* (Nonhoff et al. 2007) and *DYNC1H1*
280 (Loschi et al. 2009)) are either key components of processing bodies (P-bodies) or regulate the
281 assembly or growth of P-bodies in response to stress. Selection on the immune system is
282 suggested by the second category, “antigen processing and presentation of peptide antigen via
283 MHC class I”. The signal in this category is due to six genes, of which only *HLA-A* is an MHC
284 gene, with the other genes being *B2M*, *ERAP1*, *PDI3*, *SEC13*, and *SEC24B*. With FDR < 0.1 ,
285 there are three more significant categories related with immunity: “T cell co-stimulation”,
286 “negative regulation of complement-dependent cytotoxicity”, and “type I interferon signalling
287 pathway”. There is thus a preponderance of immunity-related GO categories and genes
288 involved in anti-viral activity (see Discussion). Even the “cytoplasmic mRNA processing body
289 assembly” category is potentially linked to virus infection as P-bodies are cytoplasmic RNA
290 granules manipulated by viruses to promote viral survival and achieve infection (Tsai and
291 Lloyd 2014, Lloyd 2013). Interestingly, the PBSnj eastern genes are also enriched in three sets
292 of Viral Interacting Proteins (VIPs) (see Supplementary files 19-22), which are genes with no
293 annotated immune functions but that interact with viruses (Enard et al. 2016). As VIP sets do
294 not in the main contain classic or known immunity genes, this provides an independent signal

295 for the relevance of viruses in this gene set. Together, these results suggest that adaptation to
296 pathogens, and viruses in particular, may have had an important role in the recent adaptation
297 in eastern chimpanzees.

298

299 Amongst chimpanzee viruses, the simian immunodeficiency virus (SIV) is intensively studied
300 as it is the progenitor of the human immunodeficiency virus (HIV) that created the global
301 acquired immune deficiency syndrome (AIDS) pandemic. It is also of interest here because it
302 appears to only infect natural populations of eastern and central chimpanzees (Santiago et al.
303 2002, Santiago et al. 2003, Nerrienet et al. 2005, Boué et al. 2015), and because it has mediated
304 fast, recent adaptations in other natural hosts (Svardal et al. 2017). Svardal *et al.* (2017)
305 investigated a set of genes that change expression in response to SIV infection in SIV natural
306 hosts (vervet monkeys) but not in non-natural hosts that develop immunodeficiency
307 (macaques) (Jacquelin et al. 2014, Jacquelin et al. 2009), hereafter referred to as “natural host
308 SIV responsive genes”. Natural host SIV responsive genes are likely involved in the specific,
309 early immune response of natural hosts to SIV infection, which limits the effects of the virus
310 and prevents subsequent immunodeficiency. These genes show signatures of positive selection
311 in vervet monkeys, suggesting that ongoing adaptation to the virus in natural hosts can occur
312 (Svardal et al. 2017). Strikingly, the PBSnj eastern tail SNPs are significantly enriched in these
313 same natural host SIV responsive genes (Jacquelin et al. 2014, Jacquelin et al. 2009) (observed
314 118 genes, expected 100, p -value = 0.0195, GOWINDA, FDR < 0.1 see Methods, Table 1,
315 Supplementary file 23). In fact, the set of natural host SIV responsive genes can fully explain
316 the unique eastern signature: the asymmetry in the PBSnj tail is abolished when this set of
317 genes is removed from the analysis (genic enrichment in the eastern PBSnj tail decreases from
318 1.36 to 1.26, and the 95% confidence interval of this point estimate now overlaps those of
319 Nigeria-Cameroon and central chimpanzees (Methods). A reduction in the genic enrichment in

320 the PBSnj tail is expected, as it is enriched in natural host SIV responsive genes; but this
321 exercise allows us to show that in the absence of selection in natural host SIV responsive genes,
322 the signature of recent positive selection in eastern chimpanzees would not be exceptional. The
323 natural host response in vervet monkeys is associated with changes in the expression of these
324 natural host SIV responsive genes. In agreement with potential adaptations in gene expression,
325 the set of PBSnjE SNPs in the natural host SIV responsive genes is further enriched in sites
326 with a high likelihood having an inferred gene regulatory function ($p=0.0485$ when compared
327 with other PBSnj eastern tail genic SNPs, $p=0.0089$ with all genic SNPs) and strongly depleted
328 of sites with no predicted regulatory function ($p=0.0001$ when compared with other PBSnj
329 eastern tail genic SNPs, $p=0.0001$ with all genic SNPs, Supplementary file 24). While these
330 genes were not identified in chimpanzees, this suggests a similar mechanism of adaptation to
331 SIV (or to an unknown virus with a similar effect in gene expression) in vervet monkeys and
332 chimpanzees.

333

334 **Biological functions of the PBSnj central tail SNPs**

335 Despite having a larger long-term N_e than eastern chimpanzees, central chimpanzees do not
336 show a clear genomic signature of recent adaptation. Despite being naturally infected by SIV
337 and being the source of pandemic HIV, they show no clear indication of selection in SIV
338 responsive genes: the PBSnj central tail has a greater number of SNPs in SIV responsive genes
339 than expected (36 vs. 29), but the enrichment is non-significant ($p = 0.0756$; resampling test,
340 Table 1). Power to identify a significant enrichment may be hampered by the low number of
341 SNPs. However, highly differentiated SNPs in the PBSnj long branches of central chimpanzees
342 are significantly enriched in one GO category, “chemokine receptor activity”, due to SNPs in
343 *CCR3*, *CCR9* and *CXCR6* ($p = 0.00001$, FDR = 0.0197, GOWINDA). Each of these genes is
344 located within the large cluster of cytokine receptor genes on chromosome 3, but they appear
345 to be associated with different sweep events (Figure 4-supplement 3). These genes are of

346 interest because *CCR3* and *CXCR6* have paralogs (*CCR5* and *CXCR4*) that in humans are the
347 two most common co-receptors for HIV-1 cell entry (Berger 1997, Moore et al. 2004). Both
348 *CCR3* and *CXCR6* can be used to enter the cell by some SIV, HIV-1 and HIV-2 subtypes
349 (Nedellec et al. 2010, Gorry et al. 2007, Bron et al. 1997, Willey et al. 2003), and the SIV of
350 both Sooty mangabey (Elliott et al. 2015) and Vervet monkey (Wetzel et al. 2017) use *CXCR6*.
351 The breadth of co-receptors used by SIV in chimpanzees is unknown, but sequence changes in
352 the V3 section of the virus can modify the specificity of the co-receptors used by HIV (Gorry
353 et al. 2007). We note that one of the PBSnj tail SNPs in *CCR3* results in an amino acid
354 substitution (246 S/A) in transmembrane domain 6, and this region has been implicated in the
355 modulation of CCR5 activity (Steen et al. 2013). Thus, changes in these co-receptors may have
356 the potential to affect the entry of SIV in chimpanzee cells.

357

358 **Discussion**

359 Comparing whole genomes from the four subspecies of chimpanzees we find that the alleles
360 whose frequency rose quickly and substantially in particular chimpanzee subspecies, resulting
361 in strong genetic differentiation, are enriched in genic sites. Except for eastern chimpanzees,
362 these genic enrichments can be explained by the effects of BGS as is the case in human
363 populations (Coop et al. 2009, Hernandez et al. 2011, Key et al. 2016). Our PBSnj statistic
364 shows that this signature in eastern chimpanzees is due to SNPs whose frequency have changed
365 specifically in eastern chimpanzees since their divergence with central chimpanzees. These
366 sites tend to have high long-range LD, but few of them have significant LD signatures of
367 positive selection. Many of them are polymorphic in central chimpanzees, so it is likely that
368 many of these adaptations have occurred from standing genetic variation and consist of soft
369 sweeps (Hermisson and Pennings 2017). This would suggest that adaptation from standing
370 genetic variation is important throughout primate evolution, not just in recent human evolution
371 (Pritchard, Pickrell, and Coop 2010). Alternatively, some of these sites may be polymorphic in

372 central chimpanzees due to gene flow from eastern chimpanzees. The inferred chimpanzee
373 demography includes recurrent migration between eastern and central chimpanzees, in both
374 directions (de Manuel et al. 2016), requiring selection in eastern chimpanzees to be strong
375 enough to overcome the homogenising effect of gene flow.

376

377 These strongly differentiated alleles in eastern chimpanzees are enriched in sites with inferred
378 regulatory function, but not in sites that have been strongly constrained during mammalian
379 evolution. This agrees well with a role in adaptation to pathogens, which is often characterized
380 by fast arms-race evolution. The PBSnj eastern genes are enriched in several immune-related
381 categories, with many of them having known or potential virus-related functions. *OAS2* and
382 *RNASEL*, for example, are involved in foreign RNA degradation (Sadler and Williams 2008),
383 while *ERAPI* is a gene under long term balancing selection in humans (Andrés et al. 2010) that
384 is involved in MHC class I epitope presentation (Hearn, York, and Rock 2009). These are
385 plausible adaptations to viral infections in eastern chimpanzees. In fact, these PBSnj eastern
386 sites are located disproportionately in genes that differentiate the CD4 transcriptional response
387 to SIV in a natural host species that tolerates the virus from a non-natural host species that
388 develops immunodeficiency. Selection acting on this set of genes is sufficient to produce the
389 greater eastern signal. Two aspects of this enrichment are notable. First, these genes are
390 identified based on gene expression responses in vervet monkeys and macaques to SIV
391 infection (Jacquelin et al. 2014, Jacquelin et al. 2009), and are thus completely independent of
392 chimpanzee genetics. Second, the SIV responsive genes also show diversifying selection in
393 vervet monkeys (Svardal et al. 2017). Of note, these SNPs are strongly enriched in putative
394 regulatory functions, in agreement with putative adaptations through gene expression. This
395 suggests that SIV may continue to exert an important selective force in natural primate species,

396 which both vervet monkeys and eastern chimpanzees may respond to by shaping gene
397 expression.

398

399 How may this happen? The genes that are both SIV-responsive and contain PBSnj eastern tail
400 SNPs are significantly enriched in four GO categories (FDR < 0.1, GOWINDA,
401 Supplementary file 25). The top category is “type I interferon signalling pathway with four
402 genes (*IRF2*, *RNASEL*, *HLA-A* and *SP100*). This category is also significantly enriched in the
403 full set of PBSnj eastern tail SNPs. *OAS2* is also in this category but it is inducible in both
404 vervet and macaque shortly after SIV infection. *IRF2*, *RNASEL* and *SP100* are all upregulated
405 in the CD4 cells of vervet monkey but not of macaques one day post infection. This is relevant
406 as regulation of the interferon response is a key differentiator between natural and non-natural
407 SIV hosts (Harris et al. 2010) and the timing of interferon responses can be key in the
408 progression to AIDS in humans infected with HIV (Rotger et al. 2011, Utay and Douek 2016).
409 Another enriched category is “polycomb group (PcG) protein complex”. PcG complexes can
410 be involved in the epigenetic regulation of HIV-1 latency (Friedman et al. 2011, Khan et al.
411 2018), and three of the genes in this GO category, *PHC2*, *CBX7* and *KDM2B* encode
412 components of the same PcG complex, PCR1 (Khan et al. 2018). While in vervet monkeys
413 *CBX7* and *PHC2* are both downregulated in CD4 cells six days post infection, *KDM2B* is
414 upregulated 115 days post infection, which might hint at increased epigenetic control of SIV
415 in the chronic phase of infection.

416

417 Of course, it is also possible that other viruses elicited a selection response in eastern
418 chimpanzees, and in particular the SIV signature that we observe could be due to selection by
419 other ssRNA viruses. Possibilities include the viruses involved in the three significant sets of
420 VIPs, which are *Dengue virus* and the closely related *Bovine leukaemia virus* and *human*

421 *T-lymphotropic virus*. Nonetheless, SIV is a better candidate to explain our observations. The
422 set of genes that explains the PBSnj eastern signature, the natural host SIV responsive genes,
423 have a clear functional, direct involvement in response to SIV virus in African primates
424 (Svardal et al. 2017). These genes show signatures of recent positive selection also in vervet
425 monkeys (Svardal et al. 2017), suggesting that SIV is an important selective force even in
426 natural hosts. Chimpanzees have been classically described as natural SIV hosts, but some
427 reports suggest fitness consequences in populations of eastern chimpanzees infected with the
428 virus (Keele et al. 2009), with some infected individuals described as having an AIDS-like
429 pathology. It thus seems likely that the virus is a selective force in this subspecies. Thus, while
430 we cannot be completely certain that SIV is driving selection in eastern chimpanzees, this virus
431 is the best candidate considering all currently available evidence.

432
433 It is also probable that eastern chimpanzees have adapted to additional selective pressures
434 unrelated of viral pathogens or immunity. An obvious candidate would be life history traits.
435 For example, the gene *SKOR2*, which contains the fifth ranked eastern specific missense
436 polymorphism, has been associated with the timing of female puberty in human GWAS of age
437 of menarche (Pickrell et al. 2016). Unfortunately, the genetic basis of these traits is poorly
438 understood making it hard to contextualise this result.

439
440 Perhaps surprisingly, central chimpanzees have weaker signatures of natural selection despite
441 being the subspecies with the largest N_e (de Manuel et al. 2016, Prado-Martinez et al. 2013). A
442 few factors could blunt the evidence for positive selection in central chimpanzees, but none of
443 them are able to explain the observed difference in PBSnj tail genic enrichment between central
444 and eastern chimpanzees –including putative population substructure, gene flow from eastern
445 chimpanzees or introgression from bonobos (Appendix 7). Central chimpanzees also do not

446 have significant enrichment in SIV responsive genes despite, like eastern chimpanzees, being
447 naturally infected by SIV (Heuverswyn et al. 2007). However, central chimpanzees exhibit a
448 significant enrichment of highly differentiated SNPs in the chemokine genes *CCR3*, *CCR9* and
449 *CXCR6*, paralogs of which are involved with HIV cell entry in humans. *CCR3* and *CXCR6* are
450 used by SIV, HIV-1 and HIV-2 subtypes (Nedellec et al. 2010, Gorry et al. 2007, Bron et al.
451 1997, Willey et al. 2003, Wetzel et al. 2017, Elliott et al. 2015) and a SNP in *CCR3* may
452 modulate the activity of the channel. The signature of positive selection in *CXCR6* is interesting
453 because the SIV of natural hosts Sooty Mangabey (Elliott et al. 2015) and Vervet monkeys
454 (Wetzel et al. 2017) predominantly use *CXCR6* for host cell entry. This is in contrast with the
455 dominant *CCR5* usage in hosts such as humans and macaques that progress to AIDS. While it
456 is unclear which particular channels are used by SIV in each chimpanzee subspecies, the
457 evidence of selection in central chimpanzees in these receptors raises the intriguing possibility
458 that the two chimpanzee hosts have used strikingly distinct evolutionary responses to the virus:
459 limiting cell entry in central chimpanzees; modulation of gene expression response in eastern
460 chimpanzees. With the estimated time of infection being ~100,000 years ago, this could be due
461 to differential adaptation to a common selective pressure, or potential subspecies-specific
462 coadaptation between chimpanzee hosts and SIV. This last is an intriguing possibility that
463 warrants further investigation.

464

465 While our attention has focussed on eastern, and to a lesser extent, central chimpanzees, this is
466 not to say that positive selection has not acted on western and Nigeria-Cameroon chimpanzees.
467 Rather, what we conclude is that, as in the case of central chimpanzees, BGS under the inferred
468 demography of chimpanzees is adequate to explain the patterns of genetic differentiation in
469 these subspecies. We note however that the divergence time of these lineages makes tests of
470 allele frequency differentiation less well suited to identify adaptive loci than in eastern and

471 central chimpanzees. Alternative approaches, for example using intensive within subspecies
472 sampling, can help identify adaptive loci in these subspecies. Nevertheless, our results show
473 striking differences between the sister subspecies of eastern and central chimpanzees. Besides
474 helping us start to identify the genetic and phenotypic differences among subspecies, this
475 finding highlights the need for genetic studies and conservation efforts to account for functional
476 differentiation between subspecies and local populations across the entire chimpanzee range.

477

478 **Materials and Methods**

479 **Genotypes, haplotypes and genic regions**

480

481 We analysed the 58 chimpanzee genomes described in de Manuel *et. al.* (2016), with sample
482 sizes of: eastern 19, central 18, Nigeria-Cameroon 10, western 11 after excluding the hybrid
483 Donald. For most tests based on allele frequencies, we used the chimpanzee VCF file from de
484 Manuel *et al.*, (2016) after removing every SNP with at least one missing genotype across all
485 chimpanzees. For haplotype phasing, we also included the 10 bonobo genomes from (de
486 Manuel *et al.* 2016). To statistically phase haplotypes we used *Beagle (Browning and Browning*
487 *2007) v 4.1* (downloaded from https://faculty.washington.edu/browning/beagle/b4_1.html,
488 May 2016). We used default parameters without imputation, except that after the initial 10 burn
489 in iterations we performed 15 phasing iterations (default is five) using the following command
490 line: `java -Xmx12000m -jar beagle.03May16.862.jar gt=vcf out=vcf.phased impute=false`
491 `nthreads=1 niterations=15`

492 For the analysis of δ we chose to use the homologous human genome reference allele as the
493 ancestral state for chimpanzee SNPs. We used the human genome from the 1000 Genomes
494 project phase III human_g1k_v37.fasta, available from:
495 http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/human_g1k_v37.fasta.gz

496 We used the UCSC liftover utility to convert chimpanzee SNPs' coordinates from pantro 2.1.4
497 to human genome version 37 (hg19) coordinates, then used samtools faidx to retrieve the
498 human allele for that position.

499

500 We acknowledge that some AA inferences can be incorrect due to parallel mutations in the
501 human lineage or lineage sorting effects. To show that our result was robust to AA inference
502 method, we also used the homologous gorilla allele and parsimony of both the human and
503 gorilla allele extracted from 20 mammalian multiz alignment to the human genome hg38,
504 downloaded from UCSC
505 (<http://hgdownload.cse.ucsc.edu/goldenPath/hg38/multiz20way/maf/>), and the inferred
506 chimpanzee ancestral allele determined calculated using the EPO alignments and downloaded
507 from ENSEMBL, available at [ftp://ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-90/fasta/ancestral_alleles/pan_troglodytes_ancestor_CHIMP2.1.4_e86.tar.gz)
508 [90/fasta/ancestral_alleles/pan_troglodytes_ancestor_CHIMP2.1.4_e86.tar.gz](ftp://ftp.ensembl.org/pub/release-90/fasta/ancestral_alleles/pan_troglodytes_ancestor_CHIMP2.1.4_e86.tar.gz). Each of these
509 inference methods recovered the same signal of a significantly greater δ tail bin genic
510 enrichment for eastern vs. central chimpanzees, see Supplementary file 14. Again, we also note
511 that our new statistic PBSnj does not require inference of the ancestral allele.

512

513 We considered protein-coding genes on the autosomes (17,530 genes) and define 'genic sites'
514 by extending the transcription start and end coordinates from ENSEMBL biobank for
515 pantro2.1.4 by 2kb on each side.

516

517 **Genetic map**

518 For statistics that required a genetic map, we used the pan diversity genetic map (Auton et al.
519 2012) inferred from 10 western chimpanzees. We downloaded the
520 `chimp_Dec8_Haplotypes_Mar1_chr-cleaned.txt` files from

521 birch.well.ox.ac.uk/panMap/haplotypes/genetic_map. These files consist of SNPs and their
522 inferred local recombination rate. These map data were inferred from sequences aligned to the
523 pantro2.1.2 genome, so we used two successive liftover steps to convert the coordinates of sites
524 used to infer the genetic map to pantro2.1.4 coordinates: pantro2.1.2 to pantro2.1.3, then
525 pantro2.1.3 to pantro2.1.4. Two steps are required as there are no liftover chains relating
526 pantro2.1.2 to pantro2.1.4. Of the 5,323,278 autosomal markers, 33,263 were not lifted from
527 pantro2.1.2 to pantro2.1.3. The remaining 5,290,015 were also successfully converted to
528 pantro2.1.4 coordinates. After liftover we filtered sites that after the two steps were mapped to
529 unassigned scaffolds or the X chromosome, which left 5,289,844 SNPs. Next, we sorted loci
530 by position to correct cases where their relative order was scrambled. This left a final number
531 of 5,289,460 autosomal SNPs. Recombination rates were then recalculated by linear
532 interpolation between consecutive markers (marker x, marker y) using the average of their
533 estimated recombination rates (rate x, rate y). These recombination maps have been deposited
534 on Dryad (see Data availability).

535

536 **Signed difference in derived allele frequency (δ)**

537 Using the derived allele frequency of each SNP for each subspecies we calculated, for each
538 pair of chimpanzee subspecies, the signed difference in derived allele frequency (DAF)
539 between them: $\delta = \text{DAF}_{\text{pop1}} - \text{DAF}_{\text{pop2}}$; $\text{DAF}_{\text{pop1}} > \text{DAF}_{\text{pop2}} : \delta > 0$; $\text{DAF}_{\text{pop1}} < \text{DAF}_{\text{pop2}} : \delta < 0$;
540 $-1 \leq \delta \leq 1$. We bin δ into 10 bins of 0.2. The choice of subspecies assigned to pop1 or pop2
541 is arbitrary and has no bearing on the results. To ensure that both tail bins are identically wide,
542 we define them as Bin 1: $-1 \leq \delta \leq 0.8$ and Bin 10 as $0.79 < \delta \leq 1$. As a consequence, the
543 Bin 5 ($0.00 < \delta < 0.2$) is marginally narrower than the other bins (by 0.01), but it contains a
544 large number of sites and the slight size difference has negligible impact on the analyses. The
545 derived allele counts have been deposited on Dryad (see Data availability).

546 We estimate confidence intervals and infer *p-values* for δ genic enrichment using a weighted
547 block jackknife (Reich et al. 2009) utilising the method of Busing *et. al.* (Busing et al. 1999).
548 This has been used for analogous tests, as it accounts for linkage disequilibrium, which means
549 that SNPs in δ bins are not full independent of each other. We divide the genome into non-
550 overlapping 200kb windows to capture the blocking effect of LD. We then recalculate, for each
551 bin, the genic enrichment using a delete-1 window jackknife. We also weight the windows by
552 the total number of SNPs in them, to downweigh, within each bin, blocks with large numbers
553 of linked SNPs. We determine that two tails are differentially enriched if their 95% CIs of
554 enrichment do not overlap. For directly testing asymmetry (or in the case of PBSnj, equality)
555 using the \log_2 ratio, we use the same weighted block jackknife, and use the 95% CI as a two-
556 tailed test with $\alpha = 0.05$. Other enrichment and resampling tests are described in Methods
557 subsection “Statistics”.

558

559 **Population Branch Statistic neighbour-joining.**

560 The Population Branch Statistic (PBS, Yi et al. 2010) is a test of population specific natural
561 selection. In the framework of a three-taxon distance tree, SNPs under selection specific to one
562 population are detected as those that result in longer than expected branch lengths (large allele
563 frequency differentiation). To generate the tree, for each site, the full distance matrix of
564 pairwise F_{ST} is computed. A three taxa tree is unrooted and has only one possible topology, so
565 simple algebra allows the calculation of each branch length in the tree. Extreme outliers in the
566 distribution of PBS are considered candidates of positive selection.

567 We introduce Population Branch Statistic neighbour-joining (PBSnj) as a simple method to
568 calculate population specific branch lengths when more than three taxa are being analysed. We
569 note that related Methods have recently appeared in the literature (Cheng, Xu, and DeGiorgio
570 2017, Racimo, Berg, and Pickrell 2018). Full details are in Appendix 4, but in brief, using the

571 full matrix of pairwise F_{ST} , F_{ST} values are transformed to units of drift time as $\ln(1-F_{ST})$ (Yi et
572 al. 2010). For fixed differences this transformation is mathematically undefined i.e. $\ln(0)$, and
573 $F_{ST}=1$ is replaced with the next largest observed F_{ST} value for a given population pair. Then
574 the Neighbor-Joining algorithm (Saitou and Nei 1987) is used to infer the tree topology and
575 calculate branch lengths. This overcomes errors in the inferred length of external branches due
576 to misspecification of a fixed tree topology. To enable a binning scheme of PBSnj values that
577 is comparable between subspecies, these scores are further normalised to be on the 0-1 scale.
578 These data have been deposited on Dryad (see Data availability). F_{ST} for PBSnj was calculated
579 using the estimator described in (Bhatia et al. 2013) because there are unequal sample sizes for
580 the subspecies, and the classical Weir and Cockerham estimator can be biased with unequal
581 sample sizes (Bhatia et al. 2013). To calculate genic enrichments along the PBSnj distribution
582 we bin SNPs in PBSnj bins 0.2 units wide. As for δ analyses, we use the 200 kb weighted block
583 jack-knife to estimate confidence and significance levels. We provide a source code file,
584 written in R, to calculate PBSnj (“PBSnj_method.R”).

585

586 **Model of Chimpanzee demographic history**

587 The most detailed exploration of chimpanzee demography comes from the work of de Manuel
588 *et. al.* (2016). This paper describes the 58 chimpanzee full genome sequences we use here, and
589 estimation of their inferred demographic model. As this paper took a primary interest in
590 investigating chimpanzee-bonobo post speciation gene flow, and to reduce the number of
591 parameters to be estimated, models were inferred using either Nigeria-Cameroon or western
592 chimpanzees, but not both. Thus, de Manuel *et. al.* (2016) provides “bonobo, eastern, central,
593 Nigeria-Cameroon” and “bonobo, eastern, central, western” models. These are referred to,
594 respectively, as ‘becn’ and ‘becw’ models below.

595

596 For this investigation we use a merged demographic history. To begin the construction of this
597 model, we recognised that there is little gene flow involving western chimpanzees in the ‘becw’
598 model, but that gene flow events are a key determinant of patterns of chimpanzee genetic
599 diversity and differentiation in the ‘becn’ model. We therefore used the ‘becn’ model as a
600 scaffold to which parameters relating to western chimpanzees (bottlenecks, expansions and N_e
601 estimates) from the ‘becw’ model are “grafted” in, to create a merged ‘becnw’ model. To make
602 sure that the N_e of western chimpanzees was appropriately scaled, all N_e s 1000 ya pastwards
603 for western chimpanzees specified in the ‘becw’ model were normalized by multiplying by the
604 ratio of the inferred N_e s of central chimpanzees specified from 1000 ya pastwards in the ‘becn’
605 and ‘becw’ models: scaled western $N_e = \text{western } N_e * 3.66914400056 / 4.3158739382$. Present
606 western N_e was normalised by the ratio of the present central N_e : scaled western $N_e = \text{western}$
607 $N_e * 0.3092 / 0.30865$.

608

609 Initially, we used the split time of the western and Nigeria-Cameroon lineages of 250ky
610 reported by de Manuel *et. al.* which was estimated from sequence divergence data, but this
611 gave a bad fit to F_{ST} values, being substantially lower than observed (Supplementary file 3).
612 We addressed this by increasing the western/Nigeria-Cameroon divergence time in proportion
613 to the ratio of model:observed western/Nigeria-Cameroon F_{ST} . i.e. $F_{ST}^{\text{Observed}} / F_{ST}^{\text{Model}} =$
614 $\text{timeX} / 250\text{kya} \Rightarrow \text{timeX} = F_{ST}^{\text{Observed}} / F_{ST}^{\text{Model}} \times 250\text{kya}$. We adjust the observed F_{ST} by
615 -0.008 – to capture the average difference between model versus observed F_{ST} values for
616 central/eastern/Nigeria-Cameroon chimpanzees. This simple calculation results in an adjusted
617 time of 267kya for the western/Nigeria-Cameroon split. F_{ST} values for this new model show a
618 much better fit to observed values (Supplementary file 3), and it is this model that we use for
619 all subsequent modelling of genic enrichments and the effects of background selection.

620

621 To determine model fit above, we calculated all pairwise average F_{ST} values for the simulated
622 data and compared them to the empirical F_{ST} estimates. For each scenario, we simulated
623 1,000,000 2kb fragments (2 Gb of sequence).

624

625 All simulations of neutral diversity and background selection were performed with *msms*
626 (Ewing and Hermisson 2010), and following de Manuel *et. al.* assuming a mutation rate of
627 $1.2e^{-8}$ and recombination rate $0.96e^{-8}$, with the following command line:

628

```
629 msms 116 1 -t 0.96048 -r 0.768384 2001 -I 5 0 38 36 20 22 0 -n 1 0.0742 -n 2 0.3181 -n 3  
630 0.3092 -n 4 0.0386 -n 5 0.08114434 -m 1 2 0 -m 1 3 0 -m 1 4 0 -m 2 1 0 -m 2 3  
631 1.8181960943074 -m 2 4 0 -m 3 1 0 -m 3 2 2.02290154800773 -m 3 4 0 -m 4 1 0 -m 4 2 0 -m  
632 4 3 0 -m 5 1 0 -m 1 5 0 -m 5 2 0 -m 2 5 0 -m 5 3 0 -m 3 5 0 -m 4 5 0 -m 5 4 0 -en 0.001 1  
633 1.83290809268 -en 0.001 2 1.161030985567 -en 0.001 3 3.66914400056 -en 0.001 4  
634 1.23640124358 -en 0.001 5 0.9132505 -em 0.020875 1 2 0 -em 0.020875 1 3 0 -em 0.020875  
635 1 4 0 -em 0.020875 2 1 0 -em 0.020875 2 3 1.8181960943074 -em 0.020875 2 4  
636 1.12888460726286 -em 0.020875 3 1 0 -em 0.020875 3 2 2.02290154800773 -em 0.020875 3  
637 4 0.514005225416364 -em 0.020875 4 1 0 -em 0.020875 4 2 0.61034918826118 -em 0.020875  
638 4 3 2.77081002950074 -em 0.042025 1 2 0 -em 0.042025 1 3 0.0447270935214584 -em  
639 0.042025 1 4 0.00204350937063846 -em 0.042025 2 1 0 -em 0.042025 2 3 1.8181960943074  
640 -em 0.042025 2 4 1.12888460726286 -em 0.042025 3 1 0.0340892941439601 -em 0.042025  
641 3 2 2.02290154800773 -em 0.042025 3 4 0.514005225416364 -em 0.042025 4 1  
642 0.00878072013784504 -em 0.042025 4 2 0.61034918826118 -em 0.042025 4 3  
643 2.77081002950074 -en 0.104325 2 0.0402577179646081 -en 0.104325 3 0.192594746352967  
644 -en 0.106325 3 8.73162876459514 -ej 0.106325 2 3 -em 0.106325 1 2 0 -em 0.106325 1 3  
645 0.0177338314347154 -em 0.106325 1 4 0.00204350937063846 -em 0.106325 2 1 0 -em
```

646 0.106325 2 3 0 -em 0.106325 2 4 0 -em 0.106325 3 1 0.00723425109237692 -em 0.106325 3
647 2 0 -em 0.106325 3 4 0.193855714034029 -em 0.106325 4 1 0.00878072013784504 -em
648 0.106325 4 2 0 -em 0.106325 4 3 0.00771007640703268 -en 0.21195 5 0.1223036 -en
649 0.214175 5 0.194964 -en 0.267475 4 1.23640124358 -en 0.267475 5 0.194964 -ej 0.2675 5 4
650 -en 0.41955 1 0.158405393915496 -en 0.42155 1 0.299481445247702 -en 0.473075 4
651 0.0306317427630759 -en 0.475075 4 2.79429564470655 -en 0.480625 4
652 0.0872103733618782 -em 0.480625 1 2 0 -em 0.480625 1 3 0.0177338314347154 -em
653 0.480625 1 4 0.00204350937063846 -em 0.480625 2 1 0 -em 0.480625 2 3 0 -em 0.480625 2
654 4 0 -em 0.480625 3 1 0.00723425109237692 -em 0.480625 3 2 0 -em 0.480625 3 4
655 0.193855714034029 -em 0.480625 4 1 0.00878072013784504 -em 0.480625 4 2 0 -em
656 0.480625 4 3 0.00771007640703268 -en 0.482625 3 1.66920782430592 -ej 0.482625 4 3 -em
657 0.482625 1 2 0 -em 0.482625 1 3 0.241282075772286 -em 0.482625 1 4 0 -em 0.482625 2 1
658 0 -em 0.482625 2 3 0 -em 0.482625 2 4 0 -em 0.482625 3 1 0.0101771164248256 -em
659 0.482625 3 2 0 -em 0.482625 3 4 0 -em 0.482625 4 1 0 -em 0.482625 4 2 0 -em 0.482625 4 3
660 0 -en 1.5988 3 0.00336130452736601 -en 1.6008 3 1.47105091660349 -ej 1.6008 1 3 -em
661 1.6008 1 2 0 -em 1.6008 1 3 0 -em 1.6008 1 4 0 -em 1.6008 2 1 0 -em 1.6008 2 3 0 -em 1.6008
662 2 4 0 -em 1.6008 3 1 0 -em 1.6008 3 2 0 -em 1.6008 3 4 0 -em 1.6008 4 1 0 -em 1.6008 4 2 0
663 -em 1.6008 4 3 0

664

665 As a further assessment of the fit of the model, we plotted the observed and simulated site
666 frequency spectrum (SFS), Figure 2-figure supplement 1. In general, the model fit is good,
667 being poorest for singletons (too high) and high frequency derived sites (too low). This is likely
668 due to effects of selection on the genome, which is not incorporated into the neutral
669 demographic model. We note too, that this model was computed using only the allele counts

670 from regions of the genome under weak/no selection as inferred from GERP scores, further
671 explaining the reduced fit at these two site classes.

672

673 **Simulations of chimpanzee genetic data under neutrality and background selection.**

674

675 We used *msms* to perform coalescent simulations of chimpanzee demography. To simulate the
676 effects of background selection (BGS) we modified the estimates of effective population size
677 (N_e) from the demographic model by multiplying them by a scaling factor, which represents
678 the B score or effective reduction in N_e due to BGS. 0.8, for example, reduces the N_e and hence
679 expected neutral diversity to 80% the level seen for neutral sites unlinked to regions under
680 purifying selection (Key et al. 2016). We simulated non-genic regions with $B=1$, and genic
681 regions with various strengths of BGS. We used B in the range 1-0.8, incremented by 0.025,
682 with additional 0.0125 increments between 0.9 – 0.85. For neutral regions and for each B we
683 simulated 25 million 2.0 kb loci. After processing and calculating allele frequencies, we
684 performed δ and PBSnj genic enrichments as described previously. To estimate a BGS strength
685 that best matched the observed δ genic enrichments, we performed a simple sum of squared
686 differences, summed for each δ genic enrichment bin for each pairwise comparison.

687

688 **Estimating the number of extra eastern chimpanzee adaptive events.**

689

690 We use the structure of the block jack-knife to estimate the number of adaptive events that are
691 needed to result in the PBSnj eastern tail genic enrichment being greater than that of central
692 chimpanzees or generated by BGS. Recall that to estimate the error variance on the genic
693 enrichment in each bin of PBSnj, we divided the genome into non-overlapping 200 kb blocks.
694 For each block we have the number of genic and non-genic SNPs per bin of the PBSnj
695 distribution.

696

697 For eastern chimps, there are 3475 genic SNPs contained within 842 blocks (i.e. 168 MB) in
698 the PBSnj rhs tail i.e. with a PBSnj scaled length ≥ 0.8 .

699 Of these, there are 468 blocks containing only 1 SNP i.e. 56% of blocks, 82 blocks with 10 or
700 more outlier genic SNPs. i.e. 10% of blocks, with a maximum block count of 111 genic SNPs
701 (Figure 4-figure supplement 1a).

702

703 We rank blocks by the number of genic SNPs that are outliers. Iterating over this sorted list we
704 remove blocks and recalculate the enrichment for genic SNPs. We define matching as the
705 number of iterations required to reduce the tail bin genic enrichment to below a target value.
706 We chose to order by the number of eastern tail genic SNPs as this results in a monotonically
707 decreasing genic enrichment with each block being removed.

708

709 **Haplotype/LD based tests of selection**

710

711 *iHS* is the ratio of the extended haplotype homozygosities (EHH) of derived versus ancestral
712 alleles at polymorphic loci. As EHH is measuring linkage disequilibrium (LD), a larger value
713 indicates greater LD for the derived allele. Under neutrality, derived allele frequency is
714 correlated with allele age, so *iHS* scores are standardised in bins of derived frequency.
715 Standardised scores have a mean of 0 and standard deviation of 1. Outliers are typically defined
716 as standardised *iHS* > 2 . nSL is a related statistic, but it calculates haplotype homozygosity as
717 the number of matching SNPs rather than genetic distance. This approach is less biased towards
718 regions of low recombination and is reportedly more sensitive to the detection of soft sweeps
719 (Ferrer-Admetlla et al. 2014). XP-EHH compares the homozygosity of focal haplotypes
720 between populations and we only performed XP-EHH calculations for the sister taxa:
721 central/eastern and Nigeria-Cameroon/western.

722

723 **Measures of conservation and effects on gene regulation**

724

725 We used *phastCons* (Siepel et al. 2005) to infer highly conserved sites. We used the 20
726 mammalian multiz alignment to the human genome hg38, downloaded from UCSC
727 (<http://hgdownload.cse.ucsc.edu/goldenPath/hg38/multiz20way/maf/>). To reduce the chance
728 that polymorphism in chimpanzees affects inference of conservation, we removed both the
729 chimp and bonobo reference genomes from these alignments. We estimated the phylogenetic
730 models from fourfold degenerate (nonconserved model) and codon first position sites
731 (conserved model). We then predicted base conservation scores and conserved fragments using
732 the following options: --target-coverage 0.25 --expected-length 30. Resultant conserved
733 elements covered 69.24% of the human exome, or an enrichment of 17.27. We note that
734 although we attempted to remove the *Pan* branch from our estimates, it is impossible to
735 completely avoid the use of these genomes, for example, when converting predicted conserved
736 elements from hg38 to pantro2.1.4. These results have been deposited on Dryad (see Data
737 availability).

738

739 We used regulomeDB (Boyle et al. 2012) to identify putatively regulatory role of genomic
740 sites. Due to the close phylogenetical relationship between chimpanzees and humans, we argue
741 that in lieu of any functional data for chimpanzees, inferred function from homologous
742 positions in the human genome is a useful proxy for function in the chimpanzee genome. To
743 obtain regulomeDB information for variable chimpanzee positions we used liftover to map
744 SNP coordinates from pantro2.1.4 to hg19, keeping positions that reciprocally mapped to
745 homologous chromosomes. Alan Boyle then kindly provided regulomeDB annotations for
746 these positions. In regulomeDB, lower scores reflect higher confidence in regulatory function.
747 We modified scores on the basis that scores 1a-f are given for positions that are human eQTLs,
748 which we do not use as they refer to the specific allele change in humans rather than to the

749 function of the site. Without eQTLs, scores 1a-c and 2a-c reflect the same biochemical
750 signatures and location within transcription factor motifs. Thus, we combine these scores in to
751 a new “high confident” regulatory function category. Our “non-regulatory” category includes
752 positions with regulomeDB scores of 6 or 7, which have no evidence of being regulatory. We
753 did not use sites with intermediate scores.

754

755 **Gene set enrichment analyses**

756

757 We used GOWINDA (Kofler and Schlötterer 2012) to test for enrichments in Gene Ontology
758 (GO) categories, which corrects for clustering and gene length biases. We used either GO
759 categories or custom gene lists as candidate gene sets. GO categories for humans were obtained
760 from the GO consortium (The Gene Ontology 2017, Ashburner et al. 2000) , while gene sets
761 were manually created from published sets of Viral Interaction Proteins (Enard et al. 2016) and
762 a set of genes that are differentially expressed in CD4 cells after SIV infection in the natural
763 SIV host vervet monkey but not in that non-natural host macaque (Svardal et al. 2017, Jacquelin
764 et al. 2014, Jacquelin et al. 2009).

765

766 GOWINDA has an input file format which enables flexible usage of nonstandard gene sets.
767 Genes are defined in a gtf file. We created a gtf from the ENSEMBLE gene definitions, but
768 restricted these to genes with clear 1-1 orthologs with humans. Our gtf file contained 16,198
769 of 17,530 protein coding genes. This gene set has been deposited on Dryad (see Data
770 availability). Additional inputs are the PBSnj tail SNP set, and the background SNP set (of
771 which the candidates are a subset). For all gene set enrichments, the background SNPs set was
772 the full genome-wide set of genic variants for which PBSnj could be calculated.

773

774 GOWINDA was designed to reduce false positives that result from gene length bias (the
775 probability of randomly containing an outlier SNP increases with gene length) and clustering
776 of genes (such as paralogs) that share function. It achieves this by using resampling of
777 background SNPs, which is the genome wide set of SNPs considered in a test. We use the --
778 mode *gene* switch. In this case, background SNPs are randomly sampled until the number of
779 overlapping genes matches the total number of genes overlapping the PBSnj tail SNP set.
780 Empirical *p-values* are estimated for each GO category, as the proportion of resamples which
781 contain the same or greater number of genes than the PBSnj tail SNP set, per GO category (for
782 each random background sample a pseudo *p-value* per GO category is also likewise calculated).
783 FDR at each *p-value*, *p*, is then estimated as the number of observed *p-values* less than or equal
784 to *p*, R_{obs} , divided by the total number of resamples with a *p-value* less than *p* R_{exp} i.e. $FDR =$
785 R_{obs} / R_{exp} .

786

787 It is important to note that only genic background SNPs that are within the candidate set of
788 genes (e.g. genes with GO definitions) are used in the random sampling. For the GO
789 enrichment, after filtering for gene sets with at least 3 genes, the GO definition file contains
790 definitions for 15649 genes, and 95% of genic background SNPs are used for resampling. This
791 is important, as therefore GOWINDA cannot be used to directly test for enrichment in a single
792 or small set of candidate gene sets. Providing one category, for example, would reduce the
793 background SNP set to only those background SNPs in the genes in that category. Resampling
794 can only ever return the same number of genes in this case. Thus, for VIPs and for the SIV
795 gene set, we included an additional category, which is the full set of genes in the gtf file (“all
796 gene set”). This has no effect on empirical *p-value* estimation. Its effect on FDR correction is
797 limited as R_{obs} is unchanged. For a candidate *p-value*, the all gene set will not be lower or equal

798 to it unless the candidate *p-value* is itself 1. Thus R_{obs} is unchanged. The effect on R_{exp} is hard
799 to determine, but for small empirical *p-values* should be proportionately small.

800

801 There are 98 VIP gene sets in (Enard et al. 2016), reduced to 53 when filtered for those
802 containing at least 3 genes. For these and for the GO categories we used an $FDR < 0.1$ as a cut-
803 off when discussing significant categories. There is only one SIV response genes set, so we
804 only report the empirical *p-value* and treat $p-value < 0.05$ as significant. Note that this
805 procedure does not allow the calculation of an FDR for the SIV set, nor over the family of tests
806 (SIV gene set enrichment in all four subspecies) but we tested a strong *a priori* expectation that
807 given the eastern PBSnj tail genes are enriched for viral immunity genes, this would be due to
808 ververt SIV response genes. However, to estimate such an FDR, we used a resampling scheme:
809 For each gene in the genome, we assign a weight, which is the proportion of SNPs in that gene
810 compared to the genome as a whole. This is to correct for gene length bias. We make the
811 intersect of all the SIV genes in each PBSnj tail. We then do weighted resampling from all
812 genes in the genome to create sets of genes as large as the intersect set, and calculate an
813 empirical *p-value* for each subspecies, as defined above. These empirical *p-values* are highly
814 similar to those provided by GOWINDA, suggesting that our weighting scheme effectively
815 controls for gene length bias. We then calculate the FDR for each empirical *p-value*, with R_{exp}
816 summed over all four subspecies.

817 **Natural Host SIV responsive genes underpin the eastern PBSnj tail genic enrichment.**

818 We wanted to test if selection on natural host SIV responsive genes could be the reason that
819 Eastern chimpanzees exhibit the strongest signal of genetic adaptation. Our simple test is to
820 hypothetically propose that if selection had not acted on the natural host SIV responsive genes
821 then those genes would not contribute a SNP to the PBSnj eastern genic tail. Thus, we removed
822 the genic tail SNPs from the 118 genes that are natural host SIV responsive and have SNPs in

823 the outlier bin of the eastern PBS scores. However, we don't remove the genic SNPs within
824 these genes that are in any of the other subspecies. This means we will affect the eastern genic
825 enrichment, but not the enrichment of other subspecies. We argue that this answers the question
826 "what would the eastern genic enrichment be if selection had not acted on these genes in eastern
827 chimpanzees?"

828 **Statistics**

829 To test enrichment in LD statistics, *phastCons* scores and regulomeDB scores we use random
830 resampling tests. For a candidate set of SNPs sized n , we randomly draw the same number of
831 genic SNPs. For LD statistics we calculated the mean score. For *phastCons* and regulomeDB
832 we calculate the proportion of SNPs in a category. For the LD tests, all SNP scores are
833 normalised so that genic SNPs have mean = 0 and sd = 1, within each bin of derived allele
834 frequency. Thus, tail SNPs with a high score have a higher score compared to other genic SNPs
835 with the same derived allele frequency.

836

837 For all resampling tests, *p-values* are estimated as $1 + n \text{ resamples} \geq \text{observed}$ (or $\leq \text{observed}$
838 as appropriate) / $1 + n \text{ resamples}$. Adding 1 to both the numerator and denominator ensures
839 that resampling p-values do not equal 0, which is a downward biased estimate given finite
840 resampling.

841 **Data availability**

842 Data generated in the course of this investigation and relevant for the interpretation of the
843 results presented here have been deposited with dryad:
844 <https://emea01.safelinks.protection.outlook.com/?url=https%3A%2F%2Fdatadryad.org%2Freview%3Fdoi%3Ddoi%3A10.5061%2Fdryad.2b3p518&data=02%7C01%7Cj.schmidt%40ucl.ac.uk%7C5b349af8089f4134a17b08d6a65d8526%7C1faf88fea9984c5b93c9210a11>

847 [d9a5c2%7C0%7C0%7C636879317941617285&sdata=90ktL2z4I6XxuAX7F5qxMbWx](https://doi.org/10.1101/582411)
848 [dhoWEFpP8F3KoR8vJg0%3D&reserved=0](#)

849 **Appendix 1**

850 **Signed differences in derived allele frequency (δ) amongst human populations.**

851 We were interested in comparing the recent adaptive history of chimpanzees and humans.
852 Previously (Coop et al. 2009) found that that those SNPs with the greatest allele frequency
853 difference between populations of modern humans were enriched for genic variants.
854 Subsequent work presented by (Hernandez et al. 2011) and (Key et al. 2016) replicated these
855 findings. To present consistent analyses and make more specific comparisons with
856 chimpanzees, we replicated the analyses of signed differences in derived allele frequency (δ)
857 in three human populations: Yoruba in Ibadan, Nigeria (YRI); Japanese in Tokyo, Japan (JPT);
858 British in England and Scotland (GBR). We choose JPT and GBR because their pairwise F_{ST}
859 is 0.10 i.e. approximately the same as for eastern and central chimpanzees. The YRI- vs. non-
860 African pairwise comparisons have amongst the largest F_{ST} values of all comparisons among
861 1000 Genomes populations, ~ 0.15 (data not shown). By down sampling each population to n
862 $= 10$ or $n = 20$ individuals, we can also assess the impact of sample size, considering that the
863 range of chimpanzee samples is 10 -19. We used genotype data from the 1000 genomes phase
864 III (Genomes Project et al. 2015), without any filtering of genotypes. We used the annotated
865 ancestral allele in this same dataset (which are derived themselves from EPO alignments) to
866 polarise derived allele frequencies.

867

868 We find that sample size has moderate effects on the determination of δ tail bin genic
869 enrichment, except for GBR vs. JPT, which have few SNPs with a large frequency difference,

870 and for which the high resolution afforded with sample $n=91$ is required to ascertain a
871 significant tail bin enrichment (Figure S1a).

872

873 When comparing the 95% confidence intervals of tail bin genic enrichments we find, consistent
874 with previous results, that δ tail bin genic enrichments are symmetrical for human populations
875 when ancient DNA information is not incorporated (Key et al. 2016). This finding is consistent
876 across all sample sizes (Figure S1a), suggesting that sample size is not a contributor to the
877 stronger genic enrichment in central vs. eastern chimpanzees.

878

879 To further explore asymmetry among the genic enrichment in the two tails of δ , we repeat the
880 calculation of tail bin genic enrichment \log_2 ratios as in Main Text Figure 2. In nearly all cases
881 the enrichment is symmetric (Figure S1b). The only significant asymmetry is for an increased
882 genic enrichment in Yoruba vs. Japanese when the sample size is 91. We note that despite
883 being significant, this asymmetry is only half that observed in the comparison between eastern
884 chimpanzees and central chimpanzees. No human comparison thus shows signatures that
885 compare to those between eastern and central chimpanzees.

886 **Appendix 2**

887 **Estimating the strength of background selection required to explain δ bin genic** 888 **enrichments.**

889 Previously it has been shown that background selection (BGS) can result in genic enrichment
890 in sites with large frequency differences between populations (Coop et al. 2009, Hernandez et
891 al. 2011, Key et al. 2016). To find the strength of BGS (measured as a B score, a fraction of
892 the expected neutral diversity) that could explain genic enrichments observed in chimpanzees,
893 we simulated 25 million 2kb loci for non-genic ($B = 1$) and genic regions. For genic regions

894 we used a range of B ($1 - 0.8$) in 0.025 steps, except between $0.9 - 0.85$ for which we used a
895 step size of 0.0125. While the strength of purifying and background selection varies among
896 genes and genomic regions, this global (average) inference allows us to make comparisons at
897 the genome scale. We used the sum of squared differences between simulated and observed
898 genic enrichments for δ bins to ascertain which B provides the best fit to observed genic
899 enrichments.

900

901 We find that when fitting all δ bins for all pairwise δ , the best fit is provided by $B = 0.888$
902 (Supplementary file 4), which indicates a reduction in neutral diversity levels of 11% in genes
903 when compared with non-genes.

904

905 We repeated this exercise using data on the 12 δ tail bins alone. Doing so allows us to infer the
906 strength of BGS required to fully explain the genic enrichment in the most highly differentiated
907 SNPs, which likely harbour targets of positive selection. While assuming no influence of
908 positive selection is unrealistic, this allows us to explore whether background selection alone
909 could, in theory, explain our observations. The best fitting B in this case is 0.863, or a 14%
910 reduction in genic diversity levels due to BGS (Supplementary file 4).

911

912 In contrast, by excluding the 12 δ tail bins, the fit of observed to simulated genic enrichments
913 is less likely to be reduced due to the influence of targets of positive selection. The best fitting
914 B in this case is 0.925 (Supplementary file 4).

915

916 Comparing the relative order of magnitudes of the Sum of Squares shows that the worst fit of
917 simulated and observed genic enrichment are seen when attempting to fit all δ bins. This is an
918 indication that BGS is not the only force affecting drift and diversity levels in genes, and

919 combined with the observation that the best fit is $B=0.925$ when δ tail bins are excluded
920 suggests positive selection is contributing to the genic enrichments in the δ tails.
921 We also checked if a greater genic enrichment in eastern vs. central chimpanzees is expected
922 given the demographic history of chimpanzees and/or the effects of BGS. For each value of B
923 we modelled above, we also calculated the \log_2 ratio of the eastern and central δ tail bin.
924 No value of B in the range $1 - 0.8$ results in an asymmetry in genic enrichment between eastern
925 and central chimpanzees as great as that observed in the genomic data (max $B = 0.850, 0.103$;
926 observed 0.34). No large asymmetry is generated under the demographic model without BGS
927 ($B = 1$). Both results suggest that no combination of BGS strength can produce the difference
928 in eastern and central δ tail bin genic enrichment observed.

929 **Appendix 3**

930 **Evidence for, and explanatory power of, differing strengths of BGS amongst** 931 **chimpanzees.**

932 We note first the evidence suggesting that background selection varies little among the great
933 apes despite their large differences in N_e (Nam et al. 2017) and despite the stronger purifying
934 selection in larger N_e subspecies (Bataillon et al. 2015). Background selection is expected to
935 reduce diversity in genic regions more than in non-genic ones by removing variants linked to
936 deleterious alleles, but the action of this type of selection appears independent of N_e (Nam et
937 al. 2017). Background selection is instead determined by the distribution of fitness effects for
938 deleterious alleles, which is likely similar among the great apes owing to their generally
939 conserved gene location and function (Nam et al. 2017). Further, simulations show that the rate
940 of selective sweeps explains the larger reduction of diversity around genes in species with
941 larger N_e (Nam et al. 2017). Thus, the diversity reducing effect of background selection should
942 be the same across all four chimpanzee sub-species. We tested this comparing the levels of

943 scaled neutral diversity (π / divergence to macaque) between chimpanzee sub-species as a
944 function of the distance to the nearest gene (normalized to the lowest diversity seen for each of
945 the sub-species, Appendix 3-figure 1a). We confirmed that the relative reduction in neutral
946 variation linked to genes is the same across sub-species (both Appendix 3-figure 1a), and that
947 the nucleotide distance from genes at which neutral diversity reaches equilibrium is also
948 similar. In addition, we find that the average genomic diversity in central, eastern and Nigeria-
949 Cameroon chimpanzees has similar dependency on recombination rate and density of
950 functional features (gene coding and gene untranslated sequences and non-coding conserved
951 elements (Appendix 3-figure 1b) suggesting yet again that background selection is comparable
952 among them. Note however that functional categories appear to be worse predictors of diversity
953 levels in western chimpanzees than the other subspecies (95% CI of the bootstrap distributions
954 of *rho*, the partial *spearman's* correlation controlling from recombination rate, do not overlap).
955 We have not investigated this further, but it is possibly due to the fact that the genetic map is
956 based on a sample of western chimpanzees, is therefore most accurate for this subspecies with
957 the effect of smaller residuals in the regression of diversity on recombination rate.

958

959 Lastly, we turn to a population genetic statistical model able to estimate the reduction in neutral
960 diversity due to background selection (Corbett-Detig, Hartl, and Sackton 2015). Full details for
961 this model are given in Corbett-Detig *et. al.* (2015), but we briefly recapitulate the main points.
962 The effect of BGS is estimated as the population scaled mutation rate ($4N_e\mu$, θ) scaled by a
963 parameter G , that models BGS as a local reduction in N_e with G allowed to vary in windows
964 along the genome in proportion to the per window fraction of functional sites. The effect of
965 selective sweeps or Hitch Hiking (HH) is also estimated as θ divided by the population scaled
966 rate of sweeps, $2Nv$. Following the implementation of this model by Corbett-Detig *et. al.* (2015)
967 we calculated average neutral diversity in 500 kb windows and used the number of bp in exons

968 for functional density. We ran the `compute_gk` package on this data to estimate the effects of
969 linked selection, as described by Corbett-Detig *et. al.* (2015). Using this model, we estimated
970 the same reduction in neutral diversity in each chimpanzee subspecies (11% reduction in the
971 highest likelihood model; Supplementary file 6), indicating equivalent levels of background
972 selection among sub-species.

973

974 Despite there being no evidence to suggest that there are differences in the strength of BGS
975 between eastern and central chimpanzees, it is useful to determine if such a putative asymmetry
976 in BGS strength could lead to eastern chimpanzees having a greater tail bin genic enrichment
977 than central chimpanzees. To investigate this possibility, we performed simulations of BGS,
978 where the strength of BGS was stronger in eastern chimpanzees ($B = 0.825$) than in all other
979 chimpanzees (B range: $0.900 - 0.850$). We find that while a greater eastern B marginally
980 increases the relative magnitude of eastern chimpanzee δ tail bin genic enrichment, none of the
981 simulated ratios are within the 95% of the eastern vs. central δ tail bin \log_2 ratio. This further
982 reinforces stronger BGS in eastern than in other chimpanzees would result in differences in δ
983 tail bin genic enrichment.

984 **Appendix 4**

985

986 **Population branch statistics**

987

988 In the two-population case (Pop A, Pop B), a “scan” for targets of population selection can be
989 performed by identifying outliers – e.g. the top 5% of sites – in the genome wide distribution
990 of per site pairwise F_{ST} values, if one assumes that these outliers are likely enriched for true
991 targets of positive selection (the empirical distribution could also be compared to simulated
992 values). As pairwise F_{ST} is a summary of the joint site frequency spectrum (SFS) of two

993 populations, these outliers are the sites with the greatest site (allele) frequency difference. If
994 one considers this as an unrooted two population tree (i.e. a straight line), outliers are simply
995 those sites with the longest branch lengths. The problem of course is that there is no
996 directionality to F_{ST} , but one assumes that the population with the highest derived allele
997 frequency is the one in which selection has acted.

998

999 The Population Branch Statistic (PBS), introduced by (Yi et al. 2010) in their study that
1000 identified *EPAS1* as under selection in Tibetans, extends the pairwise F_{ST} case by the addition
1001 of a third population, Pop C. PBS is a function of the three possible pairwise F_{ST} values amongst
1002 three populations (AB, AC, and BC). As in the two-population case, there is only one unrooted
1003 tree relating three populations, with each population connected to the central node. Therefore,
1004 each population can be assigned a unique branch length or PBS value (Appendix 4-figure 1a).
1005 The branch length is indicative of the population specific change in allele frequency, and targets
1006 of positive selection can be identified as outliers. Thus, PBS overcomes the issue of assigning
1007 directionality to allele frequency differences between populations, although with the
1008 assumption that selection occurs in one branch only.

1009

1010 We wanted to analyse the joint frequency spectrum of the four chimpanzee subspecies, and
1011 used PBS as an inspiration to develop a new statistic, PBS_{nj}. We analyse a simple four
1012 population model, with two groups of sister taxa A,B and C,D sharing a common ancestor
1013 AB,CD. Split times for AB,CD, A,B and C,D are 0.2, 0.1, 0.1 scaled time units respectively,
1014 and population size is $10e^3$ throughout. We performed 2 million simulations of a 2kb locus,
1015 with mutation rate = $1.2e^{-8}$ and recombination rate = $0.96e^{-8}$, and sampling 50 chromosomes
1016 per population. The msms (Ewing and Hermisson 2010) command line used is:

1017

1018 msms 200 1 -t 0.96048 -r 0.768384 -I 4 50 50 50 50 -n 1 1 -n 2 1 -n 3 1 -n 4 1 -en 0.1 2 1 -en
1019 0.1 4 1 -ej 0.1 1 2 -ej 0.1 3 4 -en 0.2 4 1 -ej 0.2 2 4.

1020

1021 We take A as the focal population. There are three possible combinations of F_{ST} values to
1022 calculate the branch length leading to population A: ABC, ABD, and ACD, denoted PBS_{ABC}
1023 *etc.* We note that in the Tibetan PBS example, populations were chosen so that one was clearly
1024 ancestral: Danish is the outgroup to Tibetan and Han. This highlights that while the underlying
1025 tree is unrooted and the Tibetan branch represents allele frequency change since their split with
1026 Han Chinese, in reality the Danish branch is a compound branch length combining branches
1027 leading from the basal Eurasian common ancestor to the Danish and the basal Eurasian
1028 common ancestor to the common ancestor of Tibetans and Han Chinese. In this sense, the
1029 Danish PBS branch would not represent population specific selection events *per se*, and its
1030 length is not an indication of selection events in the Danish. This indicates that the ability of
1031 PBS to truly distinguish population specific allele frequency changes is dependent on the
1032 configuration of populations included in its computation. To show this is true, we plot the rank
1033 correlations of the three different PBS statistics possible for PopA. PBS_{ABC} and PBS_{ABD} are
1034 highly correlated (*spearman's rho* = 0.82, Appendix 4-figure 1b) but both are poorly correlated
1035 with PBS_{ACD} , which is a compound branch length in our model (*spearman's rho* = 0.46 and
1036 0.47; PBS_{ACD} vs. PBS_{ABC} and PBS_{ABD} , Appendix 4-figure 1b).

1037

1038 That PBS_{ABC} and PBS_{ABD} are not perfectly correlated indicates that each contains independent
1039 information in delimiting the branch length of PopA, and illustrates the motivation in producing
1040 a statistic that draws upon the full four population F_{ST} matrix.

1041

1042 In deriving this statistic, we note that PBS is just a simple algebraic function of the matrix of
1043 pairwise F_{ST} values. To find PBS_{ABC} , for example: $PBS_{ABC} = (\text{distance}_{AB} + \text{distance}_{AC} -$
1044 $\text{distance}_{BC})/2$. An alternative method for finding distances in a phylogeny is the Neighbor-
1045 Joining algorithm (NJ) (Saitou and Nei 1987). Without giving the full details, NJ proceeds by
1046 calculating a Q matrix from the input distance matrix, creating a node by grouping the two taxa
1047 with the smallest Q , and re-calculating distances with respect to the new node. In this sense,
1048 branch lengths are a by-product of the NJ procedure, but nonetheless, by recording these branch
1049 lengths for each SNP NJ tree across the genome, we can generate a distribution of branch
1050 lengths analogous to PBS. For this reason, we name this proposed statistic PBSnj. While the
1051 details and actual distances calculated differ, PBS and PBSnj both define a distance for each
1052 branch in a tree, and the correlation between three-population PBS and PBSnj branch lengths
1053 suggests that these two methods are near identical in their results (spearman's $\rho = 0.995$).

1054 Extending PBS to more than three populations require fixing a topology. In the four population
1055 case, branch length A could be calculated as: $PBS_{ABCD} = (PBS_{ABC} + PBS_{ABD}) / 2$, but this
1056 assumes that the tree at each site follows the species tree ((A,B), (C,D)). It also “hides” the
1057 presence of an internal branch implicit in a bifurcating four taxa tree. While more complicated
1058 sets of algebraic functions could be combined to solve this or other conundrums, it is enough
1059 to point out that nj does not assume a topology (it is after all a topology finder) and that its
1060 algebraic rules are consistent no matter the number of taxa, the only change being the number
1061 of repetitions of the algorithm. Thus, we conclude that PBSnj is the more natural method to
1062 use. Lastly, while we have not considered it here, in theory PBSnj is extendable to any number
1063 of taxa.

1064

1065 We also do not consider the internal branch as in this investigation we are only interested in
1066 the selection pressures that differentiate extant populations of chimpanzees. Furthermore,

1067 interpretation of the direction of the internal branch in the four taxa case relies again on
1068 assuming that the derived allele is the target of selection.

1069 The schematic for calculating PBSnj is as follows:

- 1070 1. for each site, calculate the full F_{ST} matrix.
- 1071 2. apply the Neighbor-Joining algorithm on the F_{ST} matrix, i.e. generate a nj-tree.
- 1072 3. for each site, record the branch length for each taxa in the nj-tree.

1073

1074 Following the original description of PBS (Yi et al. 2010), we transform the F_{ST} values into
1075 units of drift time: $-\ln(1-F_{ST})$. As this is undefined for $F_{ST} == 1$, we substitute $F_{ST} == 1$ for the
1076 next lowest possible pairwise F_{ST} value. So that branch lengths exhibit the same range,
1077 following (Malaspinas et al. 2016) we standardised branch lengths by the total length of the
1078 tree, e.g. $PBS_{njA_scaled} = PBS_{njA} / (1 + PBS_{njA} + PBS_{njB} + PBS_{njC} + PBS_{njD} + PBS_{njINTERNAL})$.

1079 Lastly to perform genic enrichment tests analogous to derived allele frequency difference we
1080 re-scale so that values are within the range 0-1. This implies values of $PBS_{nj} \geq 0.8$ (which
1081 we use as our cut-off or PBS_{nj} genic tail bin) are those equal to 80% or more of the max
1082 possible values of PBS_{nj} , are not a quantile cut-off and can therefore contain a differing number
1083 of sites per taxa.

1084

1085 As a simple illustration of the effectiveness of PBS_{nj} to identify population specific changes
1086 in allele frequency, we asked how well the statistics identify Pop A specific allele frequency
1087 change. We plot the derived allele frequency in each of Pops A-D, for those sites for which
1088 $PBS_{njA_scaled} \geq 0.8$ (Appendix 4-figure 1c). As a comparison, we do the same for PBS_{ABC} ,
1089 PBS_{ABD} and PBS_{ACD} . PBS_{njA} clearly delineates those sites specifically differentiated in Pop A.
1090 PBS_{ACD} is the worst test statistic, as Pop B allele frequencies are nearly uniformly distributed
1091 in the range of 0-1 despite these sites being identified as pop A outliers. PBS_{ABC} and PBS_{ABD}

1092 offer a substantial improvement, but note there is a tendency for a more uniform distribution
1093 of allele frequencies in the population not included in the calculation of PBS_{ABC} and PBS_{ABD}
1094 (D and C respectively). Note too, the point masses near 0 for Pop A, and near 1 for Pops B-D
1095 in PBS_{njA_scaled} which represent those sites where PopA has a very low derived allele frequency
1096 i.e. are ancestral allele outliers.

1097 **Appendix 5**

1098 **The relationship between divergence times and N_e and the effects of BGS.**

1099 Demography varies greatly amongst the chimpanzee sub-species, with a wide range of pairwise
1100 divergence times and effective population sizes N_e (see Main Figure 1). This means that the
1101 total genetic drift between, for example, Western and Central chimpanzees is much greater
1102 than that between Central and Eastern chimpanzees. It is unknown how these differences in
1103 drift times effects genic enrichments in bins of either the signed differences in derived allele
1104 frequency (δ) or PBS_{nj} . To explore this, we again used a simple four population model
1105 (described Appendix 4). To model the effect of background selection (BGS) we can scale N_e
1106 by a value B, such that $B= 0.9$, for example, represents BGS that reduces linked neutral
1107 diversity by 10%. Genic regions were simulated using $B = 0.9$. We also allowed either the basal
1108 split time, or the split time of pop1 and pop2 to increase, therefore widening the range of
1109 divergence times.

1110

1111 For each scenario, we simulated 50 chromosomes per deme for a 2kb locus, for 2 million
1112 replicates, using a mutation rate of 1.2×10^{-8} and recombination rate of 0.96×10^{-8} .

1113

1114 **The effect of divergence time and N_e on δ tail bin SNP n.**

1115

1116 Increasing the divergence time increases the number of SNPs in both genic and non-genic tail
1117 bins, as is expected due to a greater variance in allele frequency due to genetic drift. While
1118 intuitive, it is important to demonstrate as it shows that the number of SNPs in tail bin is not
1119 itself an indication of the statistical support for selection. Increasing divergence time also
1120 reduces the genic:non-genic SNP n ratio: from ~1.5 at time = 0.1 down to ~1.08 at time = 0.4.
1121

1122 Changes in N_e also greatly affect the number of δ tail bin SNPs. We varied the simulated
1123 PopB:PopA ratio to be either 0.9, 0.5 or 0.1. On this time scale, an N_e ratio of 0.9 has a modest
1124 impact on the number of Pop2 tail bin SNPs. However, ratios of 0.5 and 0.1 result in a dramatic
1125 increase in both Pop2 genic and non-genic tail bin SNPs. Of course, this mirrors the result of
1126 increasing divergence time – for the same evolutionary time, lower N_e results in greater drift.
1127 What is also apparent is that that the lowered Pop2 N_e also results in an increase in the Pop1 δ
1128 tail bin counts. When Pop2 $N_e = 0.1$, the ratio of genic to non-genic δ tail SNPs is ~ 1 for both
1129 populations. We posit that these two factors – increased divergence and lower effective
1130 population size – explain the lower genic enrichments seen for δ calculated with western
1131 chimpanzees. A secondary point is the implication that the genic enrichment produced by a
1132 given strength of BGS decreases with drift time.

1133 1134 **The effect of divergence time and N_e on PBSnj tail bin genic enrichment**

1135
1136 Divergence time and N_e impact δ tail genic enrichment of both populations. This is because it
1137 conflates allele frequency change occurring in two populations. In contrast, PBSnj is able to
1138 determine the allele frequency change that occurs specifically in one branch of a phylogeny.
1139 To show the effect of N_e on the PBSnj tail genic enrichment, we plot the genic enrichment
1140 assuming BGS of B = 0.9 but varying the N_e of Pop2 (Appendix 5-figure 1). Given a relative
1141 $N_e = 1$, the genic enrichment in the PBSnj tail bin = 1.20, and there is no effect in reducing the

1142 N_e to 0.9 (Appendix 5-figure 1). Below $N_e = 0.9$, the genic enrichment drops precipitously to
1143 1.16 for $N_e = 0.5$ and 1.06 for $N_e = 0.1$. We suggest that this shows that BGS has a greater
1144 impact when divergence times are shorter and N_e relatively large, that is when most of the
1145 variation between lineages is still segregating. Longer divergence times and lower N_e results
1146 in a greater number of fixed differences between lineages, and BGS does not impact the
1147 divergence in genic regions to the extent that it reduces diversity and distorts the SFS of
1148 segregating variation.

1149
1150 This result is the motivation for comparing only central and eastern chimpanzee PBSnj tail
1151 genic enrichments. Not only is their pairwise divergence time the lowest amongst the
1152 chimpanzees, but given their relative N_e , we would not expect N_e to be the reason that eastern
1153 chimpanzees exhibit a greater PBSnj tail genic enrichment. Indeed, simulations recapitulating
1154 the demographic history of chimpanzees suggest that BGS produces equal genic enrichments
1155 for eastern and central chimpanzees. As well as expecting a similar level of drift in each of
1156 their branches, given a constant rate of adaptive evolution, we would also expect a similar
1157 number of adaptive events to contribute to the genic enrichment.

1158 **Appendix 6**

1159 **Estimating the strength of background selection required to explain PBSnj tail genic** 1160 **enrichments in chimpanzees.**

1161 We also determined how background selection can affect the PBSnj statistic amongst
1162 chimpanzees and found that they each have a unique value of B which best explains their PBSnj
1163 tail bin genic enrichment. We explain this by positive selection differentially influencing the
1164 tail of each species, as we do not expect (Nam et al. 2017) or observe differences in the effects

1165 of background selection across species. However, only eastern chimpanzees require a B
1166 stronger than 0.888 to achieve the observed PBSnj tail genic enrichment.

1167 Of critical importance for interpreting the greater PBSnj tail genic enrichment for eastern
1168 compared to that for central chimpanzees is the observation that across all values of B tested,
1169 simulated genic enrichments are approximately identical for these two subspecies (Figure 4-
1170 supplement 2b). Thus, demography and BGS should not produce the observed pattern. Results
1171 from our generalised four population model also indicate that the relatively small difference
1172 between eastern and central N_e are also not a likely explanation. In fact, such differences should
1173 result in a higher enrichment for central chimpanzees who have the larger N_e . We again posit
1174 that this is evidence for a greater rate of adaptive events along the eastern branch than that for
1175 the central branch.

1176

1177 In contrast, for any given strength of BGS, the simulated eastern and central genic enrichments
1178 are always greater than those of western and Nigeria-Cameroon chimpanzees. Our explanation
1179 for this is as follows: most tail SNPs for Nigeria-Cameroon, and especially for western
1180 chimpanzees, are actually fixed differences to all other chimpanzees. On the other hand, most
1181 eastern and central PBSnj tail SNPs are polymorphisms shared between these two sub-species.
1182 In addition, results from our general four population model, indicate that by increasing the
1183 lineage specific drift by increasing divergence time and/or decreasing N_e , the genic enrichment
1184 caused by BGS decreases. Again, we suggest that this indicates that BGS is more important for
1185 polymorphism than divergence. Finally, we conclude that only the eastern vs. central PBSnj
1186 tail bin comparison is informative in judging the significance of the eastern PBSnj tail genic
1187 enrichment or the likelihood that this can be explained by BGS.

1188 **Appendix 7**

1189

1190 **Demography and the evidence of positive selection in central chimpanzees**

1191

1192 We use this section to also discuss the factors that could blunt the evidence for positive
1193 selection in central chimpanzees. Population structure within the sampled central chimpanzees
1194 could reduce the apparent number of highly differentiated alleles. Fixed beneficial alleles in
1195 two divergent central chimpanzee populations would, for example, look to be segregating at
1196 intermediate frequencies if these were both sampled equally. Population structure within
1197 chimpanzee subspecies has been extensively analysed by both (Prado-Martinez et al. 2013)
1198 and (de Manuel et al. 2016). de Manuel et. al. (2016) present the results of numerous analyses
1199 within their supplementary material, including results from sNMF (Frichot et al. 2014),
1200 fineSTRUCTURE (Lawson et al. 2012), and ADMIXTURE (Alexander, Novembre, and
1201 Lange 2009) (respectively figures S11, S13 and S14 in de Manuel et. al. (2016)). For each of
1202 these analyses, there is less structuring of the sampled central chimpanzees compared to the
1203 sampled eastern chimpanzees, which in contrast often appear as a cline of variation from
1204 Tanzania and the south of the Democratic Republic of the Congo (DRC) through to Uganda
1205 and northern DRC. This suggests that unaccounted-for population structure is not a reason for
1206 weaker genic enrichment of differentiated alleles in central chimpanzees.

1207

1208 Another possible blunting mechanism is gene flow. When simulating neutral evolution and
1209 BGS, we used coalescent simulations using demographic parameters previously described in
1210 de Manuel et. al. (2016). This model includes inferred gene flow amongst Pan lineages,
1211 including that of the bonobo – central chimpanzee introgression. However, what these
1212 simulations do not address is the possibility that alleles selected in central chimpanzees were
1213 constantly stopped from reaching fixation due to the introduction of bonobo alleles until
1214 cessation of this gene flow ~40 kya. We note first, that gene flow is a general barrier to local

1215 adaptation, and that the rate of migration and strength of selection are the two key parameters
1216 determining the likelihood of reaching fixation. If gene flow into central chimpanzees was too
1217 great or selection too weak, then this could reduce the genic enrichment in population specific,
1218 highly differentiated alleles – but this would reflect the biological reality of reduced local
1219 adaptation in this subspecies. Secondly, we highlight that while bonobo introgression into
1220 central chimpanzees did occur, the scale of this gene flow is dwarfed by the ongoing, and near
1221 symmetrical, gene flow between central and eastern chimpanzees (migration into central
1222 chimpanzees from bonobo was ~ 1.6% of the ongoing rate from eastern chimpanzees, and ~
1223 1.9% of the ongoing rate of migration from central into eastern chimpanzees). These rates of
1224 migration would pose a greater barrier to adaptive population differentiation, and the signal in
1225 eastern chimpanzees is identified despite this.

1226 **Acknowledgments.**

1227 We thank: Fabrizio Mafessoni, Linda Vigilant, Mimi Arandjelovic, Paolo Gratton, Hjalmar
1228 Kühl and Lauren White of the Max Planck Institute for Evolutionary Anthropology for helpful
1229 discussions and/or comments on the manuscript; Alan Boyle for providing regulomeDB scores;
1230 Hannes Svoldal for extensive discussions and comments on the manuscript. This work is
1231 (partly) funded by the NIHR GOSH BRC. The views expressed are those of the author(s) and
1232 not necessarily those of the NHS, the NIHR or the Department of Health.

1233 **References**

1234 Alexander, D. H., Novembre, J., & Lange, K. (2009). Fast model-based estimation of ancestry
1235 in unrelated individuals. *Genome Res*, 19(9), 1655-1664. doi:10.1101/gr.094052.109
1236 Andrés, Aida M., Dennis, Megan Y., Kretzschmar, Warren W., Cannons, Jennifer L., Lee-Lin,
1237 Shih Queen, Hurle, Belen, . . . Green, Eric D. (2010). Balancing selection maintains a
1238 form of ERAP2 that undergoes nonsense-mediated decay and affects antigen
1239 presentation. *PLoS Genetics*, 6(10), 1-13. doi:10.1371/journal.pgen.1001157

- 1240 Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., . . . Sherlock,
1241 G. (2000). Gene ontology: tool for the unification of biology. The Gene Ontology
1242 Consortium. *Nat Genet*, 25(1), 25-29. doi:10.1038/75556
- 1243 Auton, Adam, Fledel-Alon, Adi, Pfeifer, Susanne, Venn, Oliver, Ségurel, Laure, Street, Teresa,
1244 . . . McVean, Gil. (2012). A fine-scale chimpanzee genetic map from population
1245 sequencing. *Science*, 336(6078), 193-198. doi:10.1126/science.1216872
- 1246 Ayache, J., Benard, M., Ernoult-Lange, M., Minshall, N., Standart, N., Kress, M., & Weil, D.
1247 (2015). P-body assembly requires DDX6 repression complexes rather than decay or
1248 Ataxin2/2L complexes. *Molecular Biology of the Cell*, 26(14), 2579-2595.
1249 doi:10.1091/mbc.E15-03-0136
- 1250 Bataillon, Thomas, Duan, Jinjie, Hvilsom, Christina, Jin, Xin, Li, Yingrui, Skov, Laurits, . . .
1251 Schierup, Mikkel H. (2015). Inference of purifying and positive selection in three
1252 subspecies of chimpanzees (*Pan troglodytes*) from exome sequencing. *Genome Biology
1253 and Evolution*, 7(4), 1122-1132. doi:10.1093/gbe/evv058
- 1254 Berger, E. A. (1997). HIV entry and tropism: the chemokine receptor connection. *AIDS*, 11
1255 *Suppl A*, S3-16.
- 1256 Bhatia, G., Patterson, N., Sankararaman, S., & Price, A. L. (2013). Estimating and interpreting
1257 FST: the impact of rare variants. *Genome Res*, 23(9), 1514-1521.
1258 doi:10.1101/gr.154831.113
- 1259 Boué, Vanina, Locatelli, Sabrina, Boucher, Floriane, Ayouba, Ahidjo, Butel, Christelle,
1260 Esteban, Amandine, . . . Liégeois, Florian. (2015). High rate of simian
1261 immunodeficiency virus (SIV) infections in wild chimpanzees in northeastern Gabon.
1262 *Viruses*, 7(9), 4997-5015. doi:10.3390/v7092855
- 1263 Boyle, Alan P., Hong, Eurie L., Hariharan, Manoj, Cheng, Yong, Schaub, Marc A., Kasowski,
1264 Maya, . . . Snyder, Michael. (2012). Annotation of functional variation in personal
1265 genomes using RegulomeDB. *Genome Research*, 22(9), 1790-1797.
1266 doi:10.1101/gr.137323.112
- 1267 Bron, R., Klasse, P. J., Wilkinson, D., Clapham, P. R., Pelchen-Matthews, A., Power, C., . . .
1268 Marsh, M. (1997). Promiscuous use of CC and CXC chemokine receptors in cell-to-
1269 cell fusion mediated by a human immunodeficiency virus type 2 envelope protein. *J
1270 Virol*, 71(11), 8405-8415.
- 1271 Browning, S. R., & Browning, B. L. (2007). Rapid and accurate haplotype phasing and
1272 missing-data inference for whole-genome association studies by use of localized
1273 haplotype clustering. *Am J Hum Genet*, 81(5), 1084-1097. doi:10.1086/521987
- 1274 Busing, Frank M. T. A., Meijer, Erik, Leeden, Rien Van Der %J Statistics, & Computing.
1275 (1999). Delete-m Jackknife for Unequal m. 9(1), 3-8. doi:10.1023/a:1008800423698
- 1276 Cagan, Alexander, Theunert, Christoph, Laayouni, Hafid, Santpere, Gabriel, Pybus, Marc,
1277 Casals, Ferran, . . . Andrés, Aida M. (2016). Natural selection in the great apes.
1278 *Molecular Biology and Evolution*, 33(12), 3268-3283. doi:10.1093/molbev/msw215
- 1279 Charlesworth, B., Morgan, M. T., & Charlesworth, D. (1993). The Effect of Deleterious
1280 Mutations on Neutral Molecular Variation. *Genetics*, 134(4), 1289-1303.
- 1281 Cheng, X., Xu, C., & DeGiorgio, M. (2017). Fast and robust detection of ancestral selective
1282 sweeps. *Mol Ecol*, 26(24), 6871-6891. doi:10.1111/mec.14416
- 1283 Consortium, The Chimpanzee Sequencing and Analysis, Waterson, Robert H., Lander, Eric S.,
1284 & Wilson, Richard K. (2005). Initial sequence of the chimpanzee genome and
1285 comparison with the human genome. *Nature*, 437(7055), 69.
1286 doi:doi:10.1038/nature04072
- 1287 Coop, Graham, Pickrell, Joseph K., Novembre, John, Kudaravalli, Sridhar, Li, Jun, Absher,
1288 Devin, . . . Pritchard, Jonathan K. (2009). The role of geography in human adaptation.
1289 *PLoS Genetics*, 5(6), e1000500-e1000500. doi:10.1371/journal.pgen.1000500

- 1290 Corbett-Detig, Russell B., Hartl, Daniel L., & Sackton, Timothy B. (2015). Natural Selection
1291 Constrains Neutral Diversity across A Wide Range of Species. *PLoS Biology*, 13(4),
1292 e1002112-e1002112. doi:10.1371/journal.pbio.1002112
- 1293 de Manuel, Marc, Kuhlwilim, Martin, Frandsen, Peter, Sousa, Vitor C., Desai, Tariq, Prado-
1294 Martinez, Javier, . . . Marques-Bonet, Tomas. (2016). Chimpanzee genomic diversity
1295 reveals ancient admixture with bonobos. *Science*, 354(6311), 477-481.
1296 doi:10.1126/science.aag2602
- 1297 Elliott, Sarah T. C., Wetzell, Katherine S., Francella, Nicholas, Bryan, Steven, Romero, Dino
1298 C., Riddick, Nadeene E., . . . Kirchhoff, F. (2015). Dualtropic CXCR6/CCR5 Simian
1299 Immunodeficiency Virus (SIV) Infection of Sooty Mangabey Primary Lymphocytes:
1300 Distinct Coreceptor Use in Natural versus Pathogenic Hosts of SIV.
1301 doi:10.1128/JVI.01236-15
- 1302 Enard, D., Cai, L., Gwennap, C., & Petrov, D. A. (2016). Viruses are a dominant driver of
1303 protein adaptation in mammals. *eLife*, 5. doi:10.7554/eLife.12469
- 1304 Enard, David, Messer, Philipp W., & Petrov, Dmitri A. (2014). Genome-wide signals of
1305 positive selection in human evolution. *Genome Research*, 24(6), 885-895.
1306 doi:10.1101/gr.164822.113
- 1307 Ewing, G., & Hermisson, J. (2010). MSMS: a coalescent simulation program including
1308 recombination, demographic structure and selection at a single locus. *Bioinformatics*,
1309 26(16), 2064-2065. doi:10.1093/bioinformatics/btq322
- 1310 Ferrer-Admetlla, Anna, Liang, Mason, Korneliussen, Thorfinn, & Nielsen, Rasmus. (2014).
1311 On detecting incomplete soft or hard selective sweeps using haplotype structure.
1312 *Molecular Biology and Evolution*, 31(5), 1275-1291. doi:10.1093/molbev/msu077
- 1313 Formenty, Pierre, Boesch, Christophe, Wyers, Monique, Steiner, Claudia, Donati, Franca,
1314 Dind, Frédéric, . . . Le Guenno, Bernard. (1999). Ebola Virus Outbreak among Wild
1315 Chimpanzees Living in a Rain Forest of Cote d'Ivoire. *The Journal of Infectious*
1316 *Diseases*, 179(s1), S120-S126. doi:10.1086/514296
- 1317 Frichot, E., Mathieu, F., Trouillon, T., Bouchard, G., & Francois, O. (2014). Fast and efficient
1318 estimation of individual ancestry coefficients. *Genetics*, 196(4), 973-983.
1319 doi:10.1534/genetics.113.160572
- 1320 Friedman, J., Cho, W. K., Chu, C. K., Keedy, K. S., Archin, N. M., Margolis, D. M., & Karn,
1321 J. (2011). Epigenetic Silencing of HIV-1 by the Histone H3 Lysine 27
1322 Methyltransferase Enhancer of Zeste 2 ν . In *J Virol* (Vol. 85, pp. 9078-9089).
- 1323 Genomes Project, Consortium, Auton, Adam, Brooks, Lisa D., Durbin, Richard M., Garrison,
1324 Erik P., Kang, Hyun Min, . . . Abecasis, Gonçalo R. (2015). A global reference for
1325 human genetic variation. *Nature*, 526(7571), 68-74. doi:10.1038/nature15393
- 1326 Gorry, P. R., Dunfee, R. L., Mefford, M. E., Kunstman, K., Morgan, T., Moore, J. P., . . .
1327 Gabuzda, D. (2007). Changes in the V3 region of gp120 contribute to unusually broad
1328 coreceptor usage of an HIV-1 isolate from a CCR5 Delta32 heterozygote. *Virology*,
1329 362(1), 163-178. doi:10.1016/j.virol.2006.11.025
- 1330 Harris, L. D., Tabb, B., Sodora, D. L., Paiardini, M., Klatt, N. R., Douek, D. C., . . . Estes, J.
1331 D. (2010). Downregulation of robust acute type I interferon responses distinguishes
1332 nonpathogenic simian immunodeficiency virus (SIV) infection of natural hosts from
1333 pathogenic SIV infection of rhesus macaques. *J Virol*, 84(15), 7886-7891.
1334 doi:10.1128/JVI.02612-09
- 1335 Hearn, A., York, I. A., & Rock, K. L. (2009). The Specificity of Trimming of MHC Class I-
1336 Presented Peptides in the Endoplasmic Reticulum1. *J Immunol*, 183(9), 5526-5536.
1337 doi:10.4049/jimmunol.0803663

- 1338 Hermisson, J., & Pennings, P. S. (2017). Soft sweeps and beyond: understanding the patterns
1339 and probabilities of selection footprints under rapid adaptation. *Methods in Ecology*
1340 *and Evolution*, 8(6), 700-716. doi:10.1111/2041-210x.12808
- 1341 Hernandez, Ryan D., Kelley, Joanna L., Elyashiv, Eyal, Melton, S. Cord, Auton, Adam,
1342 McVean, Gilean, . . . Przeworski, Molly. (2011). Classic selective sweeps were rare in
1343 recent human evolution. *Science*, 331(6019), 920-924. doi:10.1126/science.1198878
- 1344 Heuverswyn, Fran Van, Li, Yingying, Bailes, Elizabeth, Neel, Cecile, Lafay, Benedicte, Keele,
1345 Brandon F., . . . Peeters, Martine. (2007). Genetic diversity and phylogeographic
1346 clustering of SIVcpzPtt in wild chimpanzees in Cameroon. *Virology*, 368(1), 155-171.
1347 doi:10.1016/j.virol.2007.06.018
- 1348 Humle, T., Maisels, F., Oates, J. F., Plumptre, A., & Williamson, E. A. (2016). *Pan troglodytes*
1349 (*errata version published in 2016*). Retrieved from
1350 <http://dx.doi.org/10.2305/IUCN.UK.2016-2.RLTS.T15933A17964454.en>
- 1351 Jacquelin, B., Mayau, V., Targat, B., Liovat, A. S., Kunkel, D., Petitjean, G., . . . Muller-
1352 Trutwin, M. C. (2009). Nonpathogenic SIV infection of African green monkeys induces
1353 a strong but rapidly controlled type I IFN response. *J Clin Invest*, 119(12), 3544-3555.
1354 doi:10.1172/JCI40093
- 1355 Jacquelin, B., Petitjean, G., Kunkel, D., Liovat, A. S., Jochems, S. P., Rogers, K. A., . . . Muller-
1356 Trutwin, M. (2014). Innate immune responses and rapid control of inflammation in
1357 African green monkeys treated or not with interferon-alpha during primary SIVagm
1358 infection. *PLoS Pathog*, 10(7), e1004241. doi:10.1371/journal.ppat.1004241
- 1359 Keele, Brandon F., Jones, James Holland, Terio, Karen A., Estes, Jacob D., Rudicell, Rebecca
1360 S., Wilson, Michael L., . . . Hahn, Beatrice H. (2009). Increased mortality and AIDS-
1361 like immunopathology in wild chimpanzees infected with SIVcpz. *Nature*, 460(7254),
1362 515-519. doi:10.1038/nature08200
- 1363 Keele, Brandon F., Van Heuverswyn, Fran, Li, Yingying, Bailes, Elizabeth, Takehisa, Jun,
1364 Santiago, Mario L., . . . Hahn, Beatrice H. (2006). Chimpanzee reservoirs of pandemic
1365 and nonpandemic HIV-1. *Science*, 313(5786), 523-526. doi:10.1126/science.1126531
- 1366 Key, Felix M., Fu, Qiaomei, Romagne, Frederic, Lachmann, Michael, & Andres, Aida M.
1367 (2016). Human adaptation and population differentiation in the light of ancient
1368 genomes. *Nature Communications*, 7, 10775-10775. doi:10.1038/ncomms10775
- 1369 Khan, Sheraz, Iqbal, Mazhar, Tariq, Muhammad, Baig, Shahid M., & Abbas, Wasim. (2018).
1370 Epigenetic regulation of HIV-1 latency: focus on polycomb group (PcG) proteins.
1371 *Clinical Epigenetics*, 10(1), 14-14. doi:10.1186/s13148-018-0441-z
- 1372 Kofler, Robert, & Schlötterer, Christian. (2012). Gowinda: Unbiased analysis of gene set
1373 enrichment for genome-wide association studies. *Bioinformatics*, 28(15), 2084-2085.
1374 doi:10.1093/bioinformatics/bts315
- 1375 Lawson, D. J., Hellenthal, G., Myers, S., & Falush, D. (2012). Inference of population structure
1376 using dense haplotype data. *PLoS Genet*, 8(1), e1002453.
1377 doi:10.1371/journal.pgen.1002453
- 1378 Leendertz, Fabian H., Ellerbrok, Heinz, Boesch, Christophe, Couacy-Hymann, Emmanuel,
1379 Mätz-Rensing, Kerstin, Hakenbeck, Regine, . . . Pauli, Georg. (2004). Anthrax kills
1380 wild chimpanzees in a tropical rainforest. *Nature*, 430(6998), 451-452.
1381 doi:10.1038/nature02722
- 1382 Lloyd, R. E. (2013). Regulation of Stress Granules and P-Bodies During RNA Virus Infection.
1383 *Wiley Interdiscip Rev RNA*, 4(3), 317-331. doi:10.1002/wrna.1162
- 1384 Locatelli, Sabrina, Harrigan, Ryan J., Sesink Clee, Paul R., Mitchell, Matthew W., McKean,
1385 Kurt A., Smith, Thomas B., & Gonder, Mary Katherine. (2016). Why are Nigeria-
1386 Cameroon chimpanzees (*Pan troglodytes ellioti*) free of SIVcpz infection? *PLoS ONE*,
1387 11(8), e0160788-e0160788. doi:10.1371/journal.pone.0160788

- 1388 Loschi, M., Leishman, C. C., Berardone, N., & Boccaccio, G. L. (2009). Dynein and kinesin
1389 regulate stress-granule and P-body dynamics. In *J Cell Sci* (Vol. 122, pp. 3973-3982).
- 1390 Malaspinas, A. S., Westaway, M. C., Muller, C., Sousa, V. C., Lao, O., Alves, I., . . . Willerslev,
1391 E. (2016). A genomic history of Aboriginal Australia. *Nature*, *538*(7624), 207-214.
1392 doi:10.1038/nature18299
- 1393 McVicker, G., Gordon, D., Davis, C., & Green, P. (2009). Widespread Genomic Signatures of
1394 Natural Selection in Hominid Evolution. *PLoS Genetics*, *5*(5). doi:ARTN e1000471
1395 10.1371/journal.pgen.1000471
- 1396 Moore, J. P., Kitchen, S. G., Pugach, P., & Zack, J. A. (2004). The CCR5 and CXCR4
1397 coreceptors--central to understanding the transmission and pathogenesis of human
1398 immunodeficiency virus type 1 infection. *AIDS Res Hum Retroviruses*, *20*(1), 111-126.
1399 doi:10.1089/088922204322749567
- 1400 Nam, Kiwoong, Munch, Kasper, Mailund, Thomas, Nater, Alexander, Greminger, Maja
1401 Patricia, Krützen, Michael, . . . Schierup, Mikkel Heide. (2017). Evidence that the rate
1402 of strong selective sweeps increases with population size in the great apes. *Proceedings*
1403 *of the National Academy of Sciences*, *114*(7), 1613-1618.
1404 doi:10.1073/pnas.1605660114
- 1405 Nedellec, R., Coetzer, M., Shimizu, N., Hoshino, H., Polonis, V. R., Morris, L., . . . Mosier, D.
1406 E. (2010). Virus entry via the alternative coreceptors CCR3 and FPRL1 differs by
1407 human immunodeficiency virus type 1 subtype. *Journal of Viral Entry*, *4*(1), 33-33.
1408 doi:10.1128/JVI.00780-09
- 1409 Nerrienet, E., Santiago, M. L., Foupouapouognigni, Y., Bailes, E., Mundy, N. I., Njinku, B., .
1410 . . . Ayouba, A. (2005). Simian Immunodeficiency Virus Infection in Wild-Caught
1411 Chimpanzees from Cameroon. *Journal of Virology*, *79*(2), 1312-1319.
1412 doi:10.1128/JVI.79.2.1312-1319.2005
- 1413 Nonhoff, U., Ralser, M., Welzel, F., Piccini, I., Balzereit, D., Yaspo, M. L., . . . Krobitsch, S.
1414 (2007). Ataxin-2 Interacts with the DEAD/H-Box RNA Helicase DDX6 and Interferes
1415 with P-Bodies and Stress Granules. In *Mol Biol Cell* (Vol. 18, pp. 1385-1396).
- 1416 Pickrell, Joseph K., Berisa, Tomaz, Liu, Jimmy Z., Ségurel, Laure, Tung, Joyce Y., & Hinds,
1417 David A. (2016). Detection and interpretation of shared genetic influences on 42 human
1418 traits. *Nature Genetics*, *48*(7), 709-717. doi:10.1038/ng.3570
- 1419 Prado-Martinez, Javier, Sudmant, Peter H., Kidd, Jeffrey M., Li, Heng, Kelley, Joanna L.,
1420 Lorente-Galdos, Belen, . . . Marques-Bonet, Tomas. (2013). Great ape genetic diversity
1421 and population history. *Nature*, *499*(7459), 471-475. doi:10.1038/nature12228
- 1422 Pritchard, Jonathan K., Pickrell, Joseph K., & Coop, Graham. (2010). The Genetics of Human
1423 Adaptation: Hard Sweeps, Soft Sweeps, and Polygenic Adaptation. *Current Biology*,
1424 *20*(4), R208-R215. doi:10.1016/j.cub.2009.11.055
- 1425 Racimo, Fernando. (2016). Testing for ancient selection using cross-population allele
1426 frequency differentiation. *Genetics*, *202*(2), 733-750. doi:10.1534/genetics.115.178095
- 1427 Racimo, Fernando, Berg, Jeremy J., & Pickrell, Joseph K. (2018). Detecting polygenic
1428 adaptation in admixture graphs. *Genetics*, *208*(4), 1565-1584.
1429 doi:10.1534/genetics.117.300489
- 1430 Reich, D., Thangaraj, K., Patterson, N., Price, A. L., & Singh, L. (2009). Reconstructing Indian
1431 population history. *Nature*, *461*(7263), 489-494. doi:10.1038/nature08365
- 1432 Rotger, Margalida, Dalmau, Judith, Rauch, Andri, McLaren, Paul, Bosinger, Steven E.,
1433 Martinez, Raquel, . . . Telenti, Amalio. (2011). Comparative transcriptomics of extreme
1434 phenotypes of human HIV-1 infection and SIV infection in sooty mangabey and rhesus
1435 macaque. *Journal of Clinical Investigation*, *121*(6), 2391-2400. doi:10.1172/JCI45235
- 1436 Rudicell, Rebecca S., Jones, James Holland, Wroblewski, Emily E., Learn, Gerald H., Li,
1437 Yingying, Robertson, Joel D., . . . Wilson, Michael L. (2010). Impact of simian

- 1438 immunodeficiency virus infection on chimpanzee population dynamics. *PLoS*
1439 *Pathogens*, 6(9). doi:10.1371/journal.ppat.1001116
- 1440 Sabeti, Pardis C., Reich, David E., Higgins, John M., Levine, Haninah Z. P., Richter, Daniel
1441 J., Schaffner, Stephen F., . . . Lander, Eric S. (2002). Detecting recent positive selection
1442 in the human genome from haplotype structure. *Nature*, 419(6909), 832-837.
1443 doi:10.1038/nature01140
- 1444 Sabeti, Pardis C., Varilly, Patrick, Fry, Ben, Lohmueller, Jason, Hostetter, Elizabeth, Cotsapas,
1445 Chris, . . . Stewart, John. (2007). Genome-wide detection and characterization of
1446 positive selection in human populations. *Nature*, 449(7164), 913-918.
1447 doi:10.1038/nature06250
- 1448 Sadler, A. J., & Williams, B. R. G. (2008). Interferon-inducible antiviral effectors. *Nat Rev*
1449 *Immunol*, 8(7), 559-568. doi:10.1038/nri2314
- 1450 Saitou, N., & Nei, M. (1987). The neighbor-joining method: a new method for reconstructing
1451 phylogenetic trees. *Mol Biol Evol*, 4(4), 406-425.
1452 doi:10.1093/oxfordjournals.molbev.a040454
- 1453 Santiago, Mario L., Lukasik, Magdalena, Kamenya, Shadrack, Li, Yingying, Bibollet-Ruche,
1454 Frederic, Bailes, Elizabeth, . . . Hahn, Beatrice H. (2003). Foci of endemic simian
1455 immunodeficiency virus infection in wild-living eastern chimpanzees (*Pan troglodytes*
1456 *schweinfurthii*). *Journal of Virology*, 77(13), 7545-7562. doi:10.1128/JVI.77.13.7545-
1457 7562.2003
- 1458 Santiago, Mario L., Rodenburg, Cynthia M., Kamenya, Shadrack, Bibollet-Ruche, Frederic,
1459 Gao, Feng, Bailes, Elizabeth, . . . Hahn, Beatrice H. (2002). SIVcpz in wild
1460 chimpanzees. *Science*, 295(5554), 465-465. doi:10.1126/science.295.5554.465
- 1461 Siepel, Adam, Bejerano, Gill, Pedersen, Jakob S., Hinrichs, Angie S., Hou, Minmei,
1462 Rosenbloom, Kate, . . . Haussler, David. (2005). Evolutionarily conserved elements in
1463 vertebrate, insect, worm, and yeast genomes. *Genome Research*, 15(8), 1034-1050.
1464 doi:10.1101/gr.3715005
- 1465 Steen, A., Thiele, S., Guo, D., Hansen, L. S., Frimurer, T. M., & Rosenkilde, M. M. (2013).
1466 Biased and Constitutive Signaling in the CC-chemokine Receptor CCR5 by
1467 Manipulating the Interface between Transmembrane Helices 6 and 7*. In *J Biol Chem*
1468 (Vol. 288, pp. 12511-12521).
- 1469 Svardal, Hannes, Jasinska, Anna J., Apetrei, Cristian, Coppola, Giovanni, Huang, Yu, Schmitt,
1470 Christopher A., . . . Nordborg, Magnus. (2017). Ancient hybridization and strong
1471 adaptation to viruses across African vervet monkey populations. *Nature Genetics*,
1472 49(12), 1705-1713. doi:10.1038/ng.3980
- 1473 The Gene Ontology, Consortium. (2017). Expansion of the Gene Ontology knowledgebase and
1474 resources. *Nucleic Acids Res*, 45(D1), D331-D338. doi:10.1093/nar/gkw1108
- 1475 Tsai, W. C., & Lloyd, R. E. (2014). Cytoplasmic RNA Granules and Viral Infection. *Annu Rev*
1476 *Virol*, 1(1), 147-170. doi:10.1146/annurev-virology-031413-085505
- 1477 Utay, N. S., & Douek, D. C. (2016). Interferons and HIV Infection: The Good, the Bad, and
1478 the Ugly. *Pathog Immun*, 1(1), 107-116. doi:10.20411/pai.v1i1.125
- 1479 Voight, Benjamin F., Kudaravalli, Sridhar, Wen, Xiaoquan, & Pritchard, Jonathan K. (2006).
1480 A map of recent positive selection in the human genome. *PLoS Biology*, 4(3), 0446-
1481 0458. doi:10.1371/journal.pbio.0040072
- 1482 Wetzel, K. S., Yi, Y., Elliott, S. T. C., Romero, D., Jacquelin, B., Hahn, B. H., . . . Collman, R.
1483 G. (2017). CXCR6-Mediated Simian Immunodeficiency Virus SIVagmSab Entry into
1484 *Sabaeus* African Green Monkey Lymphocytes Implicates Widespread Use of Non-
1485 CCR5 Pathways in Natural Host Infections. In *J Virol* (Vol. 91).
- 1486 Willey, S. J., Reeves, J. D., Hudson, R., Miyake, K., Dejucq, N., Schols, D., . . . Clapham, P.
1487 R. (2003). Identification of a subset of human immunodeficiency virus type 1 (HIV-1),

1488 HIV-2, and simian immunodeficiency virus strains able to exploit an alternative
1489 coreceptor on untransformed human brain and lymphoid cells. *J Virol*, 77(11), 6138-
1490 6152.
1491 Yi, Xin, Liang, Yu, Huerta-Sanchez, Emilia, Jin, Xin, Cuo, Zha Xi Ping, Pool, John E., . . .
1492 Wang, Jian. (2010). Sequencing of 50 human exomes reveals adaptation to high
1493 altitude. *Science*, 329(5987), 75-78. doi:10.1126/science.1190371
1494

1495 **Tables**

1496 **Table 1:** VIP gene enrichment in the PBSnj eastern tail.

VIRUS	P-VALUE	FDR P-VALUE
BLV	0.0015	0.0239
DENV	0.0025	0.0239
HTLV	0.0145	0.0780

1497

1498

1499 **Table 2:** SIV responsive gene enrichment in subspecies PBSnj tails.

Subspecies	Observed	Expected	P-VALUE
Eastern	118	99	0.0198
Central	36	29	0.0739

1500

1501 **Figure legends**

1502 **Figure 1. *The geographic distribution and population history of chimpanzees.***

1503 **A**, The ranges of each chimpanzee subspecies within western and central sub-Saharan Africa.
1504 Range data from (Humble et al. 2016). **B**, Phylogenetic relationships amongst chimpanzees and
1505 the timing of their population divergence, modified from (De Manuel et al. 2016). **C**,
1506 Heterozygosity, reflective of relative differences in effective population sizes. Box plots show
1507 median central interquartile range, whiskers the upper and lower interquartile range. Points
1508 show individual heterozygosity. For all panels, colour designates subspecies: Blue = western,
1509 red = Nigeria-Cameroon, green = central, orange = eastern. The heterozygosity counts, and
1510 code for plotting panel C are contained in Figure 1–Source Data 1.

1511

1512 **Figure 2: *Genic enrichment in bins of signed difference in derived allele frequency (δ).***

1513 **A**, X-axis: δ is computed as the difference in derived allele frequency, for each pair of
1514 chimpanzee subspecies. Tail bins (the last bin in either end of δ) contain those SNPs with the
1515 largest allele frequency differences. Numbers are of the genic SNPs in each tail bin. Y-axis:
1516 genic enrichment in each δ bin (Methods). **B**, Genic enrichment eastern and central chimpanzee
1517 δ , plotted separately due to a different Y-axis limit. NC = Nigeria-Cameroon. The asterisk
1518 shows significance of the asymmetry in the genic enrichment (* = 0.01). Shading represents
1519 the 95% CI (i.e. alpha = 0.05 for a two-tailed test) estimated by 200kb weighted block
1520 jackknife. Grey dashed lines represent simulations under increasing levels of background
1521 selection that best match different aspects of the data: lightest to darkest shades: B= 0.925
1522 (excluding δ tail bins), 0.888 (all δ bins), and 0.863 (only δ tail bins). The observed and BGS
1523 simulated genic enrichments, and code for plotting are contained in Figure 2–Source Data 1.

1524

1525 **Figure 2- figure supplement 1: *Observed and Simulated Site Frequency Spectra.***

1526 We plot the Site Frequency Spectrum (SFS) for each chimpanzee subspecies. X axes: derived
1527 allele count. Y axes: proportion. Black: observed. Green: simulated.
1528 Simulated counts come from 25 million 2kb loci simulated with *msms*, using the chimpanzee
1529 demography specified in Methods.

1530

1531 **Figure 3: Direct quantification of δ tail bin genic enrichment asymmetry.**

1532 The asymmetry of the genic enrichments in the δ tails is measured by taking their \log_2 ratio,
1533 thus 0 indicates a symmetric enrichment (equal enrichment in both δ tails). NC = Nigeria-
1534 Cameroon. Dot = observed asymmetry. Horizontal lines represent confidence intervals
1535 estimated by 200kb weighted block jackknife (light = 95%, black = 99%, i.e. alpha = 0.05 or
1536 0.01 for a two-tailed test). Grey vertical marks represent the δ tail asymmetry in simulations,
1537 under increasing levels of background selection that best match different aspects of the data:
1538 lightest to darkest shades: B= 0.925 (excluding δ tail bins), 0.888 (all δ bins), and 0.863 (only
1539 δ tail bins). The observed and BGS simulated genic enrichment tail bin \log_2 ratios, and code
1540 for plotting are contained in Figure 3–Source Data 1.

1541

1542 **Figure 3-figure supplement 1: Stronger eastern BGS does not result in observed levels of δ**
1543 ***tail bin genic enrichment asymmetry.***

1544 The asymmetry of the genic enrichments in the δ tails is measured by taking their \log_2 ratio,
1545 thus 0 indicates a symmetric enrichment (equal enrichment in both δ tails). We created
1546 coalescent simulations in which the strength of BGS was greater in eastern chimpanzees than
1547 other subspecies. For eastern chimpanzees we chose a fixed B = 0.825, as this B provided the
1548 best fit the eastern δ tail genic enrichment. All other subspecies had the same B, in the range
1549 of 0.900 – 0.850. A larger difference in B between subspecies results in a slight increase in
1550 asymmetry, but none of the simulated differences in BGS result in the observed asymmetry.

1551 Point = observed asymmetry. Horizontal lines represent confidence intervals estimated by
1552 200kb weighted block jackknife (light = 95%, black = 99%, i.e. alpha = 0.05 or 0.01 for a two-
1553 tailed test). Grey vertical marks represent the δ tail asymmetry in simulations, under increasing
1554 levels of difference in background selection between eastern and other chimpanzees: lightest
1555 to darkest shades: All $B_{\text{eastern}} = 0.825$; $B_{\text{others}} = 0.850, 0.863, 0.850, 0.888, 0.900$.

1556

1557 **Figure 4: Genic enrichment in bins of PBSnj in eastern and central chimpanzees.** **A** X-axes:
1558 PBS scaled to take values in the range 0 -1. Y-axes: Genic enrichment computed as described
1559 in Figure 2. Shading represents the 95% CI (i.e. alpha = 0.05 for a two-tailed test) estimated
1560 by 200kb weighted block jackknife. **B**: \log_2 ratio of the eastern and central PBSnj tail (PBS \geq
1561 0.8) genic enrichment. **A,B** Grey dashed (**A**) or vertical (**B**) lines represent the PBSnj genic
1562 enrichment in simulations, under increasing levels of background selection that best match
1563 different aspects of δ , as described in Figs. 2 and 3: lightest to darkest shades: $B = 0.925$
1564 (excluding δ tail bins), 0.888 (all δ bins), and 0.863 (only δ tail bins). The observed and BGS
1565 simulated PBSnj genic enrichments, tail bin \log_2 ratios, and code for plotting are contained in
1566 Figure 4–Source Data 1.

1567

1568 **Figure 4-figure supplement 1: Number of adaptive events in eastern chimpanzees.**

1569 **a**, Most 200kb blocks contain few PBSnj eastern outlier SNPs, but there is an extended right
1570 hand tail. **b**, we ranked blocks by the number of PBSnj eastern tail SNPs, then iteratively
1571 removed outlier genic SNPs. This results in a monotonically decreasing genic enrichment, and
1572 the removal of eight blocks is required to reduce the genic enrichment of the PBSnj eastern tail
1573 to overlap the 95% CI of the PBSnj central tail, and 19 blocks to reduce it below the level of
1574 the point estimate of the central PBSnj tail. We could alternatively order the windows by the
1575 total number of outlier SNPs i.e. without regard to genic vs. non-genic. Doing so increases our

1576 estimated range of sweeps to 15-26. But we note that the genic enrichment does monotonically
1577 decrease with block removal (c). This is partly due to the arbitrary nature of the definition of
1578 genic, as it implies that there are some 200 kb blocks that have more non-genic than genic
1579 outlier SNPs contained within them, and this may very well change if the definition of genic
1580 was changed from transcription start and end sites \pm 2kb. (d) Lastly, we randomly shuffled
1581 the removal order of the 200-kb blocks. We did so for 1000 random shuffles of the block order
1582 (A single random shuffle is shown). We find that the median number of blocks (i.e. sweeps)
1583 across random shuffles is 165 to match the upper 95% CI of the central chimpanzee estimate
1584 (middle 90% quantile range 114-221; min = 78, max = 278) increased to 273 to match the
1585 central chimpanzee point estimate of genic enrichment (middle 90% quantile range 214-329;
1586 min = 162, max = 381). Such a procedure is likely an overestimate, as most of the removal
1587 steps are those removing 1 to 9 genic outlier SNPs (panel a), resulting in minimal reduction of
1588 the genic enrichment.

1589

1590 **Figure 4-figure supplement 2: Scaled PBSnj bin genic enrichment for all chimpanzee**

1591 *subspecies*.

1592 **a**, X-axes: PBS scaled to take values in the range 0 -1, per subspecies. Y-axes: Genic
1593 enrichment computed as described in Fig. 2. Shading represents the 95% CI (i.e. $\alpha = 0.05$
1594 for a two-tailed test) estimated by 200kb weighted block jackknife. Grey dashed lines represent
1595 the PBSnj genic enrichment in simulations, under increasing levels of background selection
1596 that best match different aspects of δ , as described in Figs. 2 and 3: lightest to darkest shades:
1597 $B = 0.925$ (excluding δ tail bins), 0.888 (all δ bins), and 0.863 (only δ tail bins). **b**, BGS does
1598 not result in eastern and central chimpanzees differing in the PBSnj tail bin genic enrichment.
1599 X axis is the \log_2 ratio of the PBSnj tail genic enrichment eastern / pop2. Grey shaded ticks
1600 represent the PBSnj genic enrichment in simulations, under increasing levels of background

1601 selection that best match different aspects of δ , as described in Figs. 2 and 3: lightest to darkest
1602 shades: B= 0.925 (excluding δ tail bins), 0.888 (all δ bins), and 0.863 (only δ tail bins).

1603

1604 **Figure 4-figure supplement 3: Number of sweeps in the chromosome 3 chemokine receptor**
1605 **cluster of central chimpanzees.**

1606 X axis: position along chromosome 3 (Mb). Plotted in the upper panel are the PBSnj central
1607 scores in the region encompassing *CCR3*, *CCR9*, and *CXCR6*. An independent cluster of high
1608 PBSnj scores is associated with each candidate gene. Each point represents one PBSnj score,
1609 colour has an alpha = 30% to reduce over plotting.

1610 Haplotypes are plotted in the central panel. Yellow ticks are derived alleles, blue are ancestral,
1611 while white is space so that each tick aligns with PBSnj scores. Inspection indicates that there
1612 is a degree of haplotype scrambling between each of the candidate genes. Lastly, we depict the
1613 genes in this region in the lower panel.

1614

1615 **Appendix 1 Figure 1: Genic enrichment in bins of signed difference in derived allele**
1616 **frequency (δ), for human populations from the 1000 Genomes Phase III.**

1617 **a**, X-axis: δ is computed as the difference in derived allele frequency, for each pair of
1618 populations. Tail bins (the last bin in either end of δ) contain those SNPs with the largest allele
1619 frequency differences. Numbers are of the genic SNPs in each tail bin. Y-axis: genic
1620 enrichment in each δ bin, computed as described in Methods. Shading represents the 95% CI
1621 (i.e. alpha = 0.05 for a two-tailed test) estimated by 200kb weighted block jackknife, **b**, The
1622 asymmetry of the genic enrichments in the δ tails is measured by taking their \log_2 ratio, thus 0
1623 indicates a symmetric enrichment (equal enrichment in both δ tails). Dot = observed
1624 asymmetry, with size indicating the relative sample size (10, 20, 91 individuals). Horizontal

1625 lines represent confidence intervals estimated by 200kb weighted block jackknife (light = 95%,
1626 black = 99%, i.e. alpha = 0.05 or 0.01 for a two-tailed test).

1627

1628 **Appendix 3 Figure 1: *The effect of background selection on patterns of neutral diversity in***
1629 ***chimpanzees.***

1630 **a**, Diversity levels at neutral sites as a function of the distance to the nearest gene. We
1631 calculated scale diversity (π / divergence to macaque) in bins of distance to genic regions. We
1632 then rescaled scaled diversity for each subspecies so that the diversity was in the range 0-1. **b**,
1633 To further explore the effects of BGS on chimpanzee genomes we checked the correlation of
1634 density of functional sites with neutral diversity (π). We used windows 500kb spaced at least
1635 1MB apart in the genome. Here, ρ is the spearman rank partial correlation between windowed
1636 diversity and density of functional sites per window, controlling for recombination rate (the
1637 average rate per window). Each dot represents a bootstrap replicate (random sample of 500 kb
1638 windows). We calculated the partial ρ for each bootstrap. Box plots show the median and
1639 interquartile ranges of the bootstrap replicates.

1640

1641 **Appendix 4 Figure 1: *Deriving the PBS_{nj} statistic.***

1642 **a**, PBS is just a simple arithmetic function of pairwise F_{ST} values for a group of three taxa or
1643 populations. **b**, The configuration or choice of populations determines the information content
1644 of PBS. In each panel are the spearman's ρ correlations between different PBS
1645 configurations, and between PBS and our new statistic PBS_{nj} for a simple four population
1646 model (described in Appendix 4). In each case Pop A is the focal population. PBS_{ABC} and
1647 PBS_{ABD} are highly correlated but not identical indicating that incorporating both Pops C and D
1648 would refine the identification of Pop A specific differentiated variants. PBS_{njA}, which utilises
1649 information from all four populations is more highly correlated with both PBS_{ABC} and PBS_{ABD}

1650 than they are with each other. Alpha = 10% for plotted points to reduce over saturation. **c**, For
1651 each statistic, we plot the site frequency spectrum (SFS) for each of the four populations for
1652 sites identified as outliers in Pop A. PBSnj clearly finds those sites differentiated in Pop A, and
1653 better than either PBS_{ABC} and PBS_{ABD}. In the standard PBS, the SFS in the species not included
1654 in the PBS configuration has a more uniform distribution, indicating that some sites identified
1655 as PBS frequency outliers in Pop A are not true population specific outliers.

1656

1657 **Appendix 5 Figure 1: *Effect of reduced N_e on PBSnj genic enrichments.***

1658

1659 In a simple four population model, we modelled genic regions as those with a $B = 0.9$. In
1660 population 2, we simulated four effective population size ratios (1, 0.9, 0.5, 0.1). N_e ratios of
1661 0.5 and 0.1 result in a reduced genic enrichment given the same strength of background
1662 selection. X-axes: PBS scaled to take values in the range 0 -1, per subspecies. Y-axes: Genic
1663 enrichment.

1664 **Source data files**

1665

1666 Figure 1–Source Data 1: Genome wide average heterozygosity counts for chimpanzees.

1667 Figure 2–Source Data 1: Observed and BGS simulated Genic enrichment in bins of signed
1668 difference in derived allele frequency.

1669 Figure 3–Source Data 1: Observed and BGS simulated Genic log 2 ratios of tail bin enrichment,
1670 signed difference in derived allele frequency.

1671 Figure 3-figure supplement 1–Source Data 1: Observed and BGS simulated, with greater
1672 strength of BGS in eastern chimpanzees, Genic log 2 ratios of tail bin enrichment, signed
1673 difference in derived allele frequency.

1674 Figure 4–Source Data 1: Genic enrichment in bins of PBSnj in eastern and central chimpanzees

1675 Figure 4-figure supplement 2-Source data 1 Genic enrichment in bins of PBSnj in all four
1676 chimpanzee sub-species.
1677

1678 **Supplementary Files**

1679

1680 Supplementary File 1: Signed difference in derived allele frequency, genic and non-genic tail
1681 counts.

1682 Supplementary File 2: Observed and simulated δ bin genic enrichments

1683 Supplementary File 3: Observed and model chimpanzee subspecies F_{ST}

1684 Supplementary File 4: Fit of simulated to observed genic enrichments across δ bins.

1685 Supplementary File 5: \log_2 ratio of eastern and central chimpanzee δ tail bin genic enrichments
1686 with different strengths of background selection.

1687 Supplementary File 6: Model-based reduction of neutral diversity in chimpanzee sub-species.

1688 Models are tested for their ability to explain diversity as a function of distance to functional
1689 sites.

1690 Supplementary File 7: Effect of divergence on δ tail SNP number.

1691 Supplementary File 8: Effect of N_e on δ tail SNP number.

1692 Supplementary File 9: Fitting BGS to match observed PBSnj tail genic enrichments.

1693 Supplementary File 10: PBSnj tail SNP haplotype statistic scores.

1694 Supplementary File 11: Non-synonymous PBSnj eastern tail SNPs.

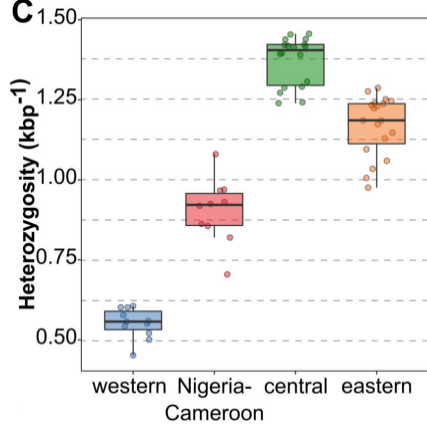
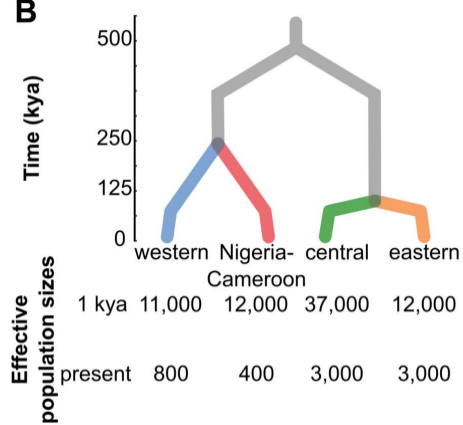
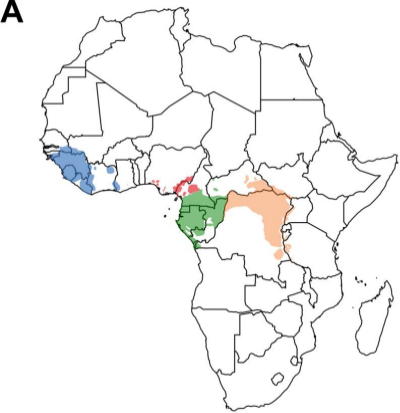
1695 Supplementary File 12: PBSnj tail SNP regulomeDB enrichments.

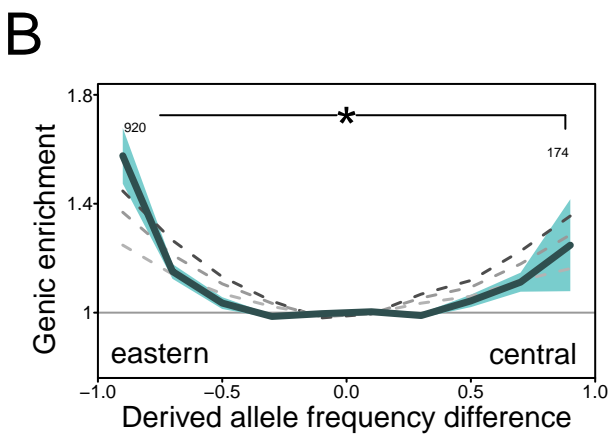
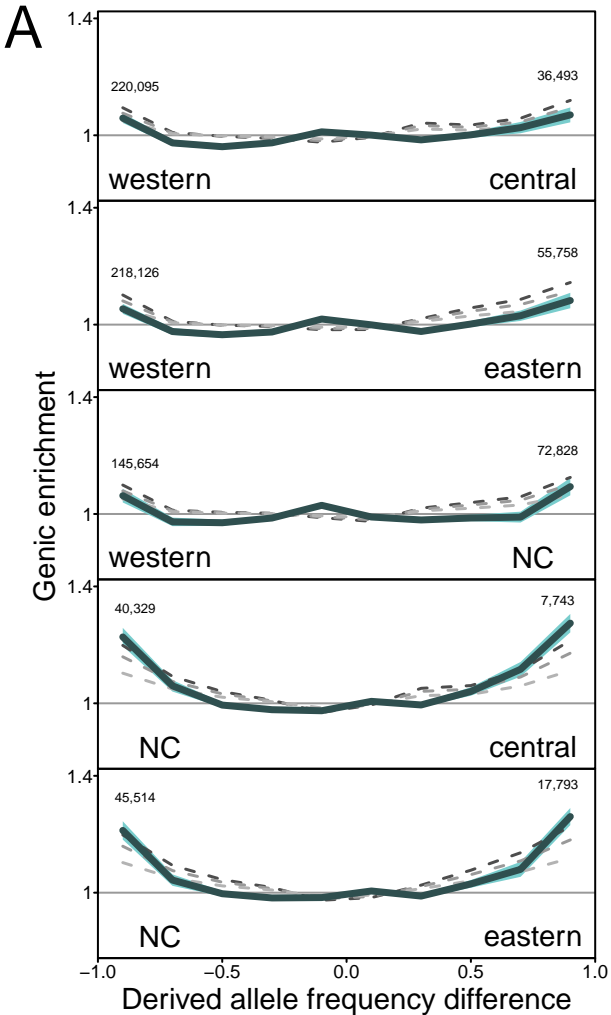
1696 Supplementary File 13: PBSnj tail SNP conservation/phastCons score enrichments.

1697 Supplementary File 14: Effect of Ancestral Allele estimation on eastern vs. central chimpanzee
1698 δ bin genic enrichments.

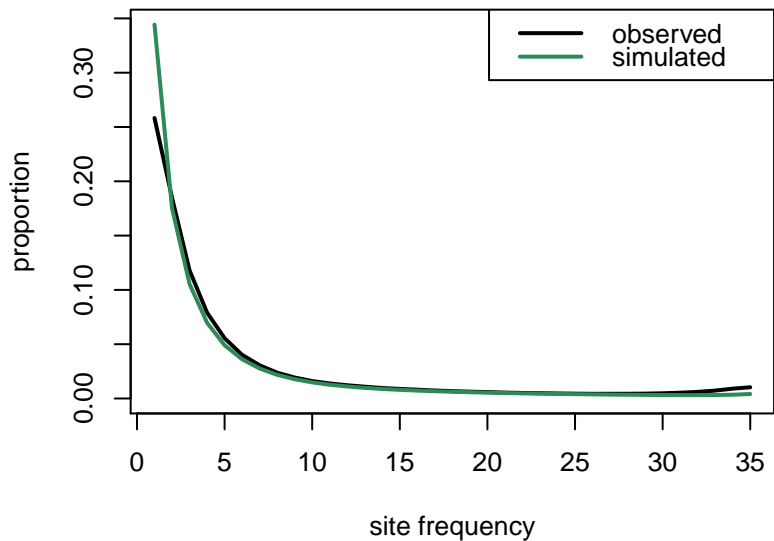
1699 Supplementary File 15: PBSnj Eastern GO enrichment.

- 1700 Supplementary File 16: PBSnj Central GO enrichment.
- 1701 Supplementary File 17: PBSnj Nigeria-Cameroon GO enrichment.
- 1702 Supplementary File 18: PBSnj Western GO enrichment.
- 1703 Supplementary file 19: PBSnj Eastern VIP enrichment.
- 1704 Supplementary file 20: PBSnj Central VIP enrichment.
- 1705 Supplementary file 21: PBSnj Nigeria-Cameroon VIP enrichment.
- 1706 Supplementary file 22: PBSnj Western VIP enrichment.
- 1707 Supplementary file 23: SIV responsive gene enrichment tests.
- 1708 Supplementary File 24: PBSnj eastern tail and SIV responsive genes SNP regulomeDB
1709 enrichment
- 1710 Supplementary file 25: PBSnj Eastern and SIV gene GO enrichment
- 1711 **Source code.**
- 1712 *PBSnj_method.R*
1713 Contains function for calculating PBSnj from the full pairwise F_{ST} matrix of four populations.
1714

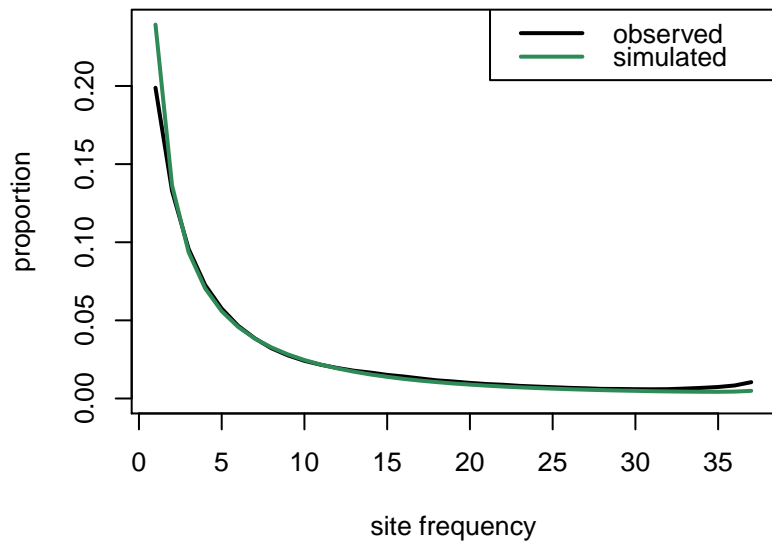




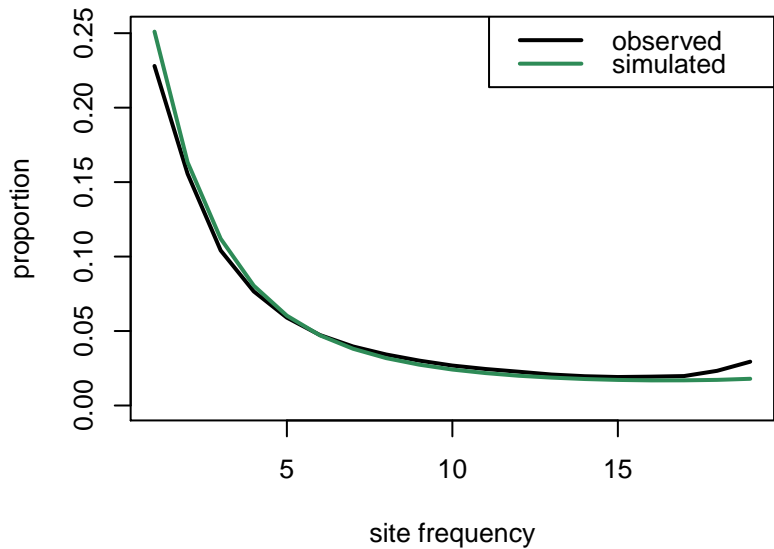
central simulated vs. observed SFS



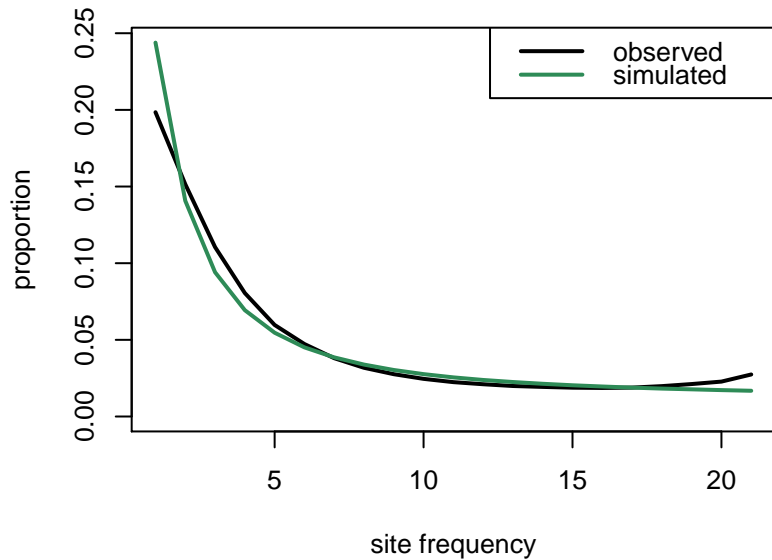
eastern simulated vs. observed SFS

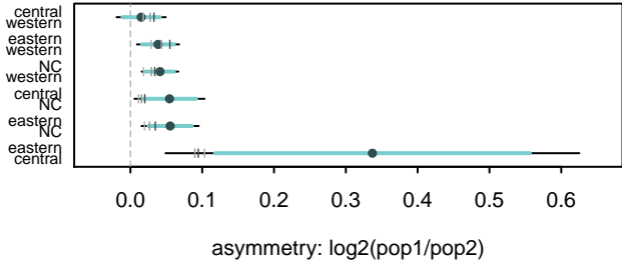


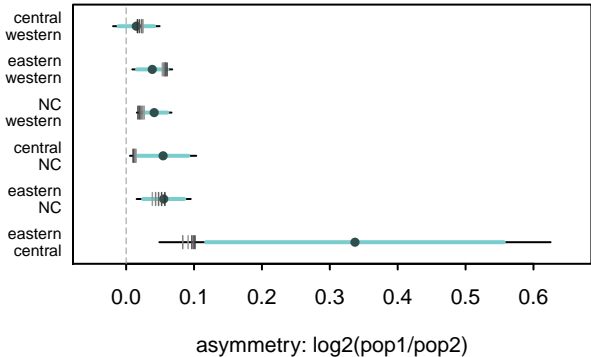
nigeria simulated vs. observed SFS

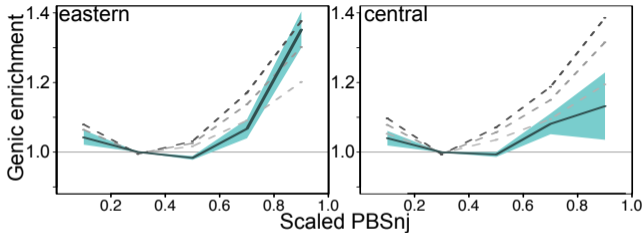
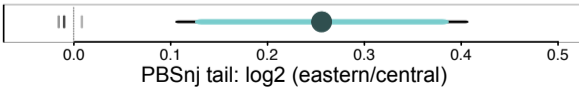


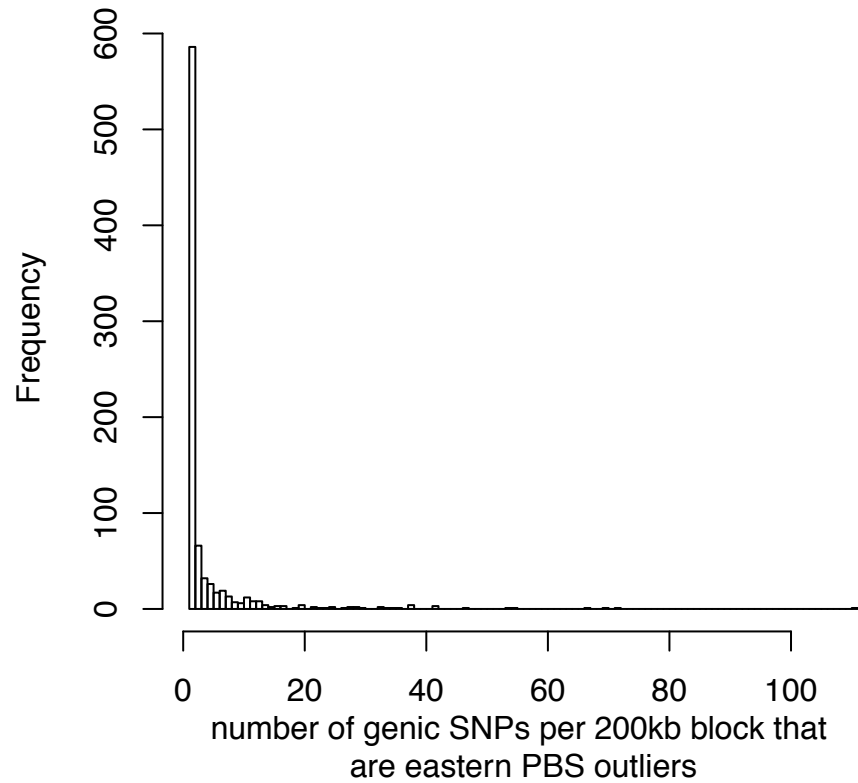
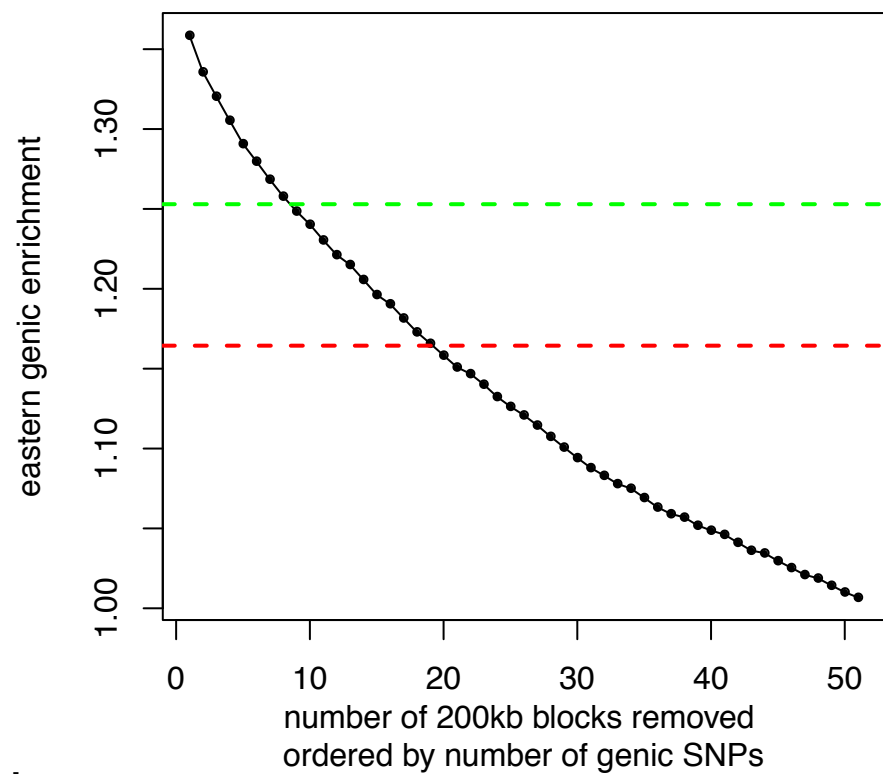
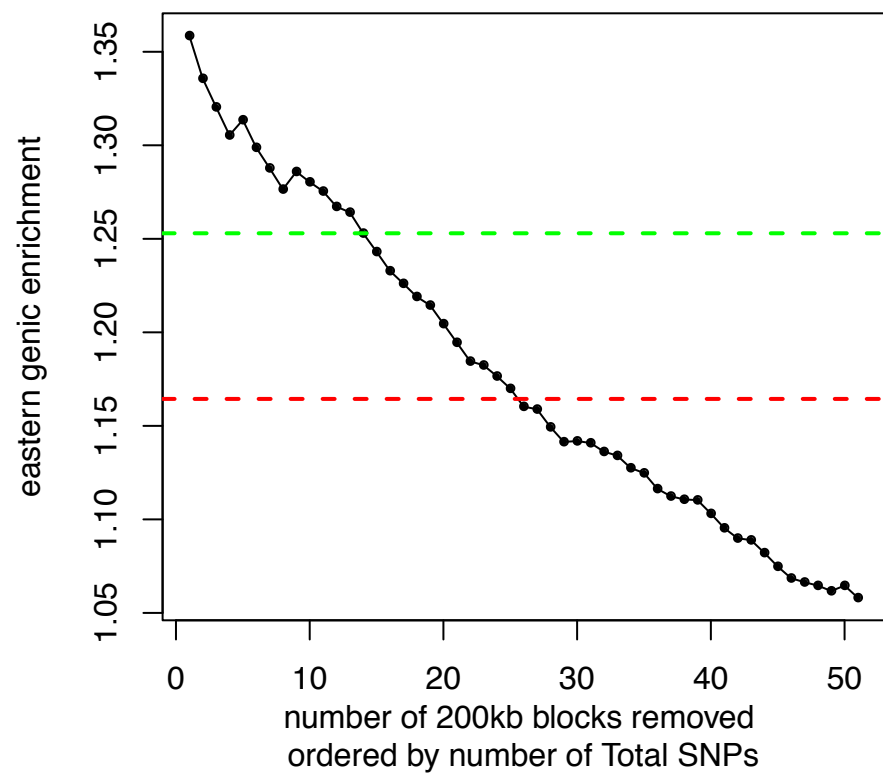
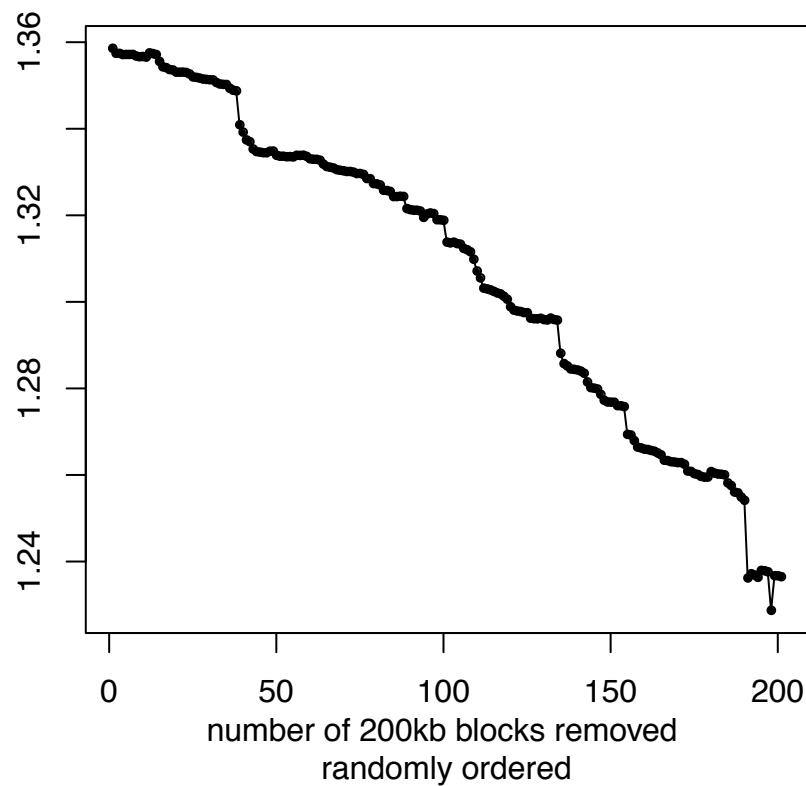
western simulated vs. observed SFS

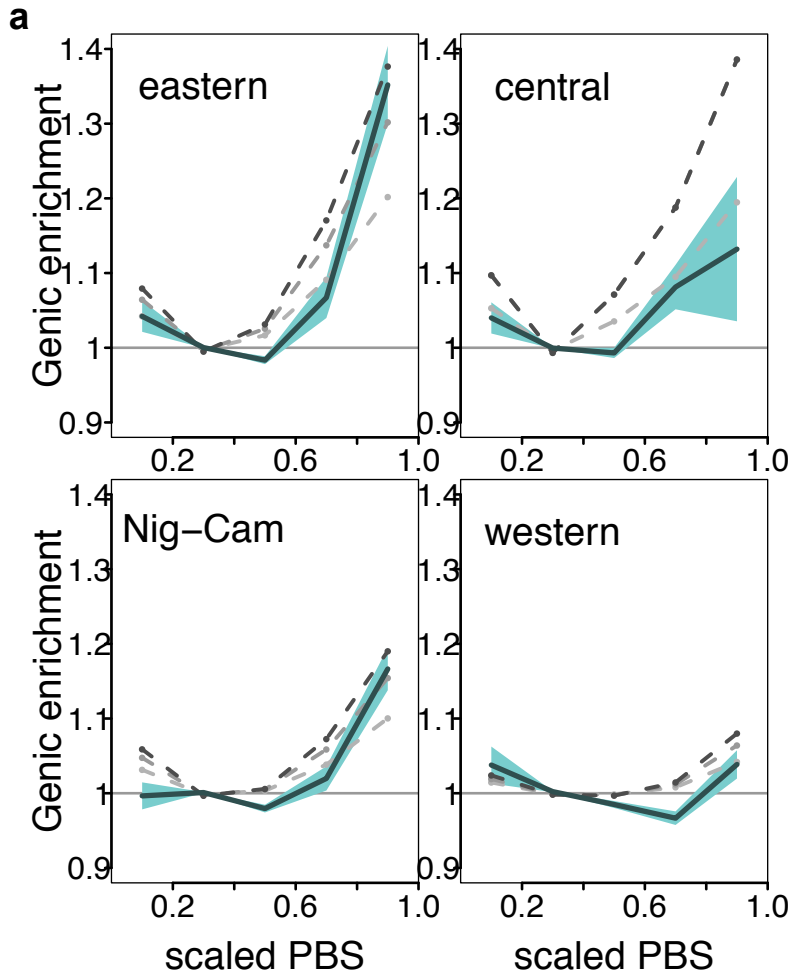




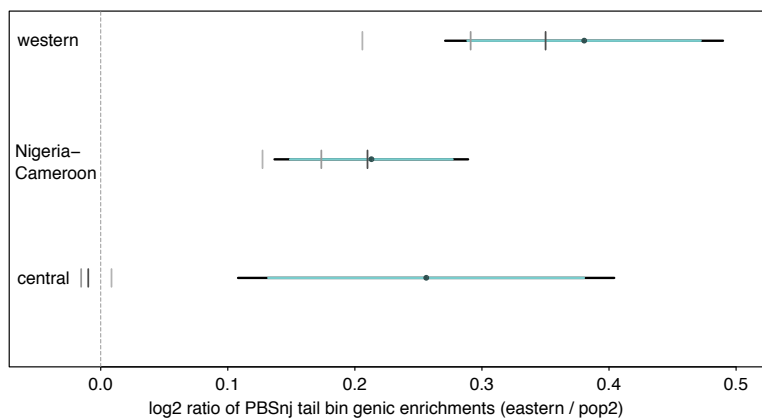


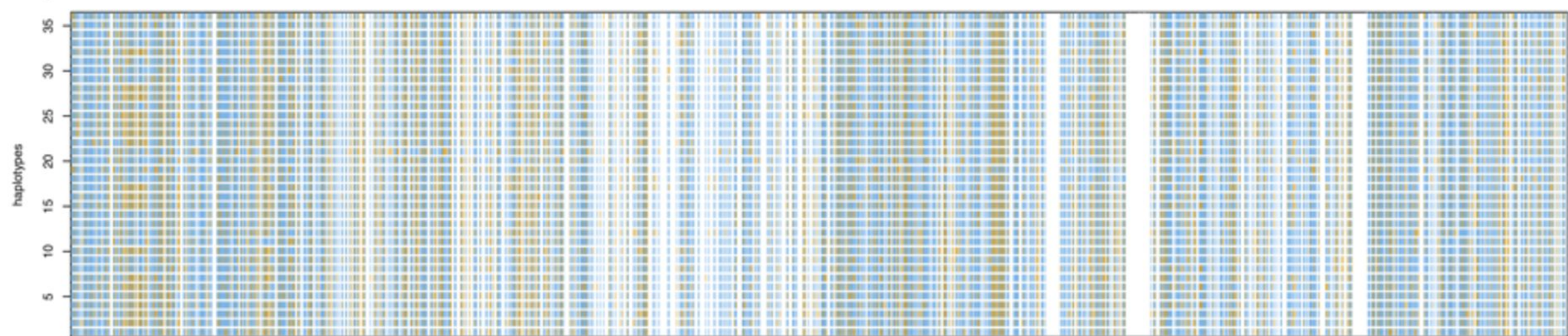
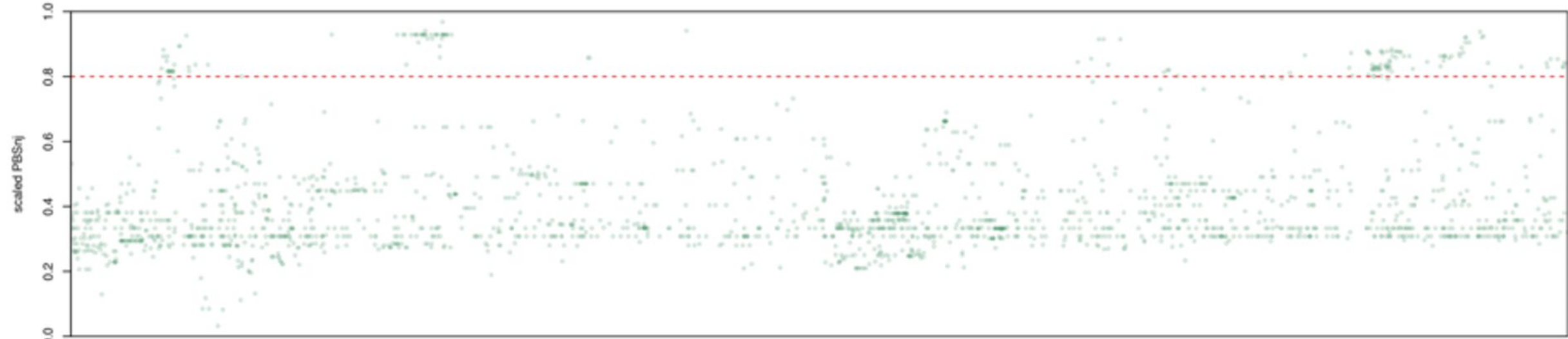
A**B**

a**b****c****d**

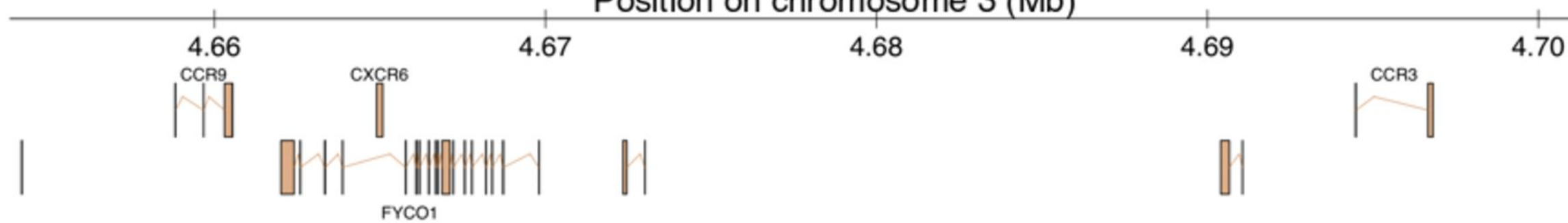


b

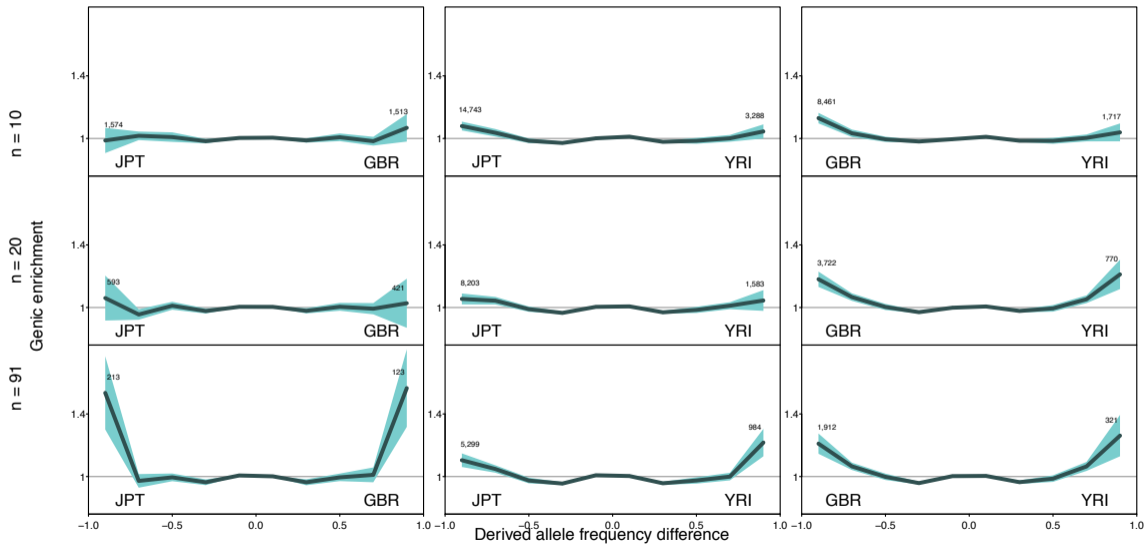




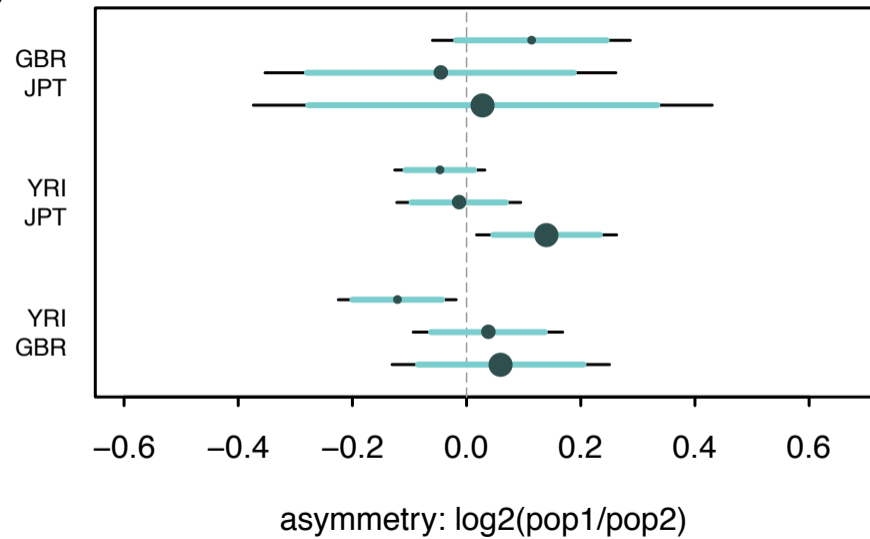
Position on chromosome 3 (Mb)

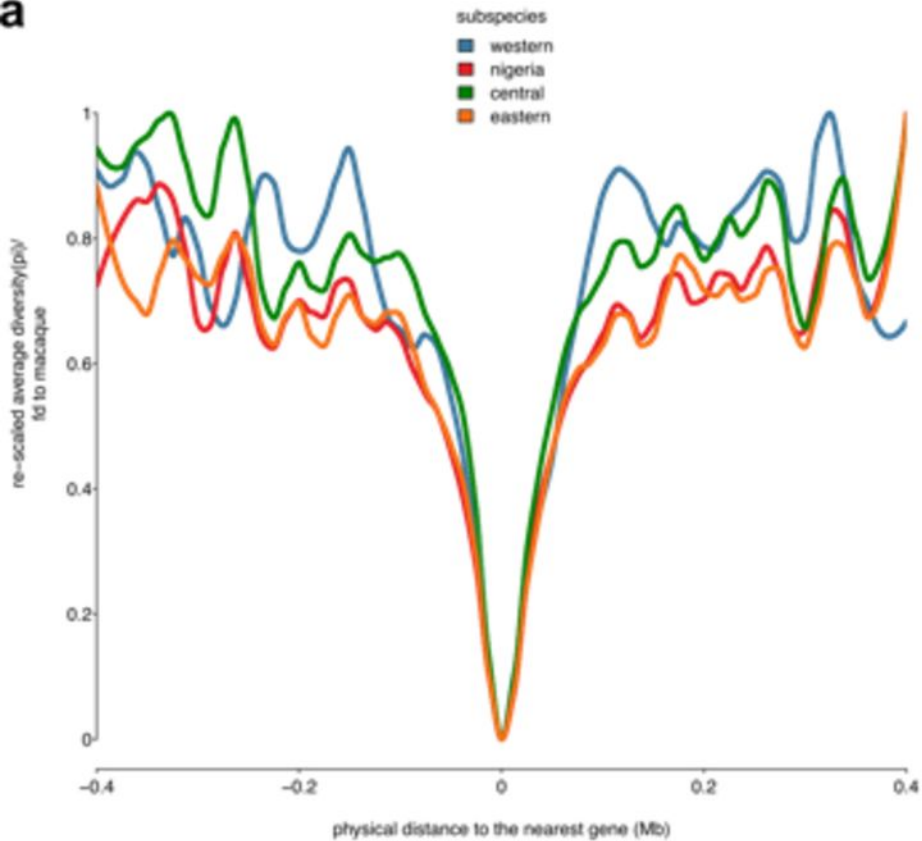
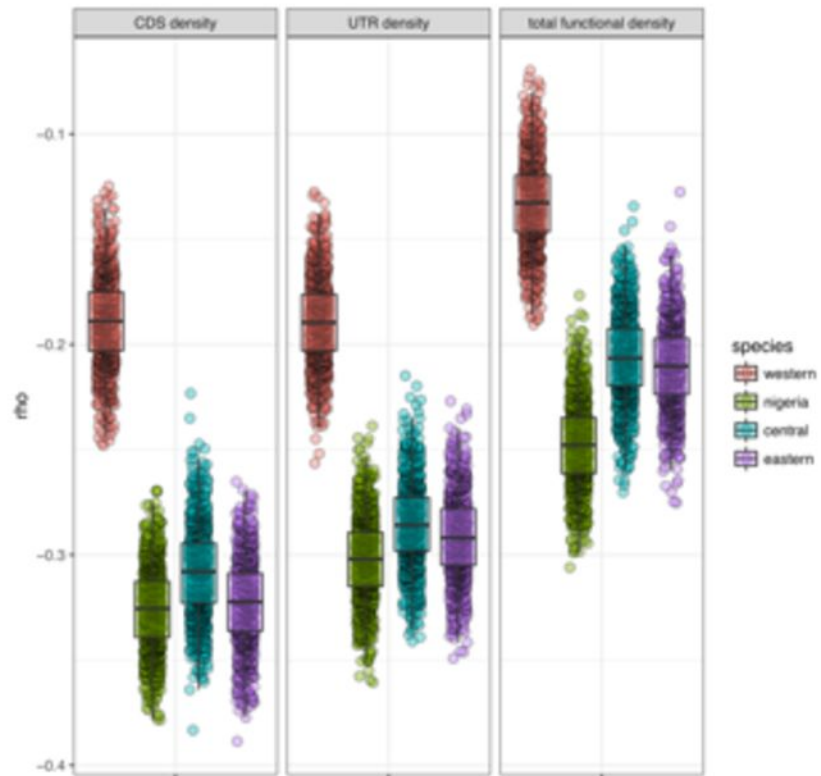


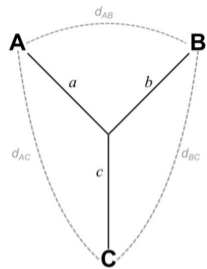
a



b



a**b**

a

$$a = [d_{AB} + d_{AC} - d_{BC}]/2$$

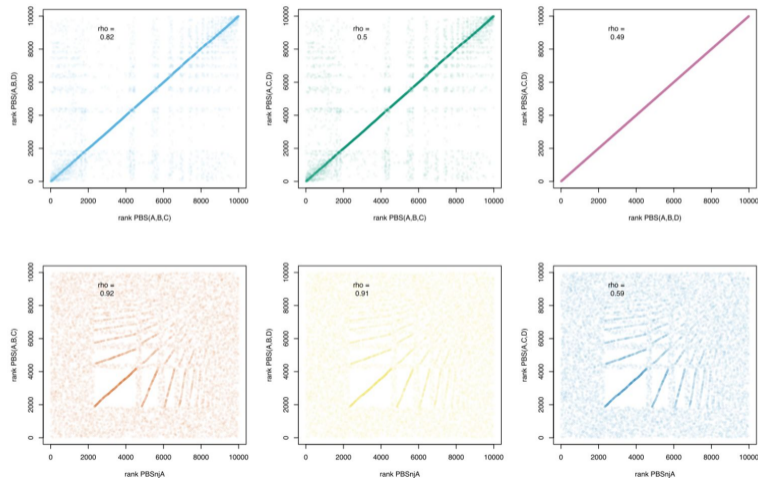
$$b = [d_{AB} + d_{BC} - d_{AC}]/2$$

$$c = [d_{AC} + d_{BC} - d_{AB}]/2$$

A, B... Taxa/populations

$d_{AB} d_{AC} \dots$ Pairwise distances ($-\ln(1-F_{ST})$)

a, b, \dots Branch lengths

b**c**