1

# Environmentally-mediated selection parallels population divergence across a chimpanzee subspecies contact zone

4

5   Matthew W. Mitchell[1,2,*], Walker Alexander[3], Dana V. Mitchell[1], Adam H. Freedman[4],

6   Janina Dordel[1], Ryan J. Harrigan[5], Ahmet Sacan[3], Fabrice Kentatchime[1,6], Bryan S.

7   Featherstone[1], Ekwoge E. Abwe[1,7,8], Paul R. Sesink Clee[1], Abwe E. Abwe[8], Sabrina

8   Locatelli[9], Bethan J. Morgan[7,8,10], Bernard Fosso[11], Roger Fotso[11], Sarah A. Tishkoff[12],

9   Evan E. Eichler[13,14], Nicola M. Anthony[15], Thomas B. Smith[5,16], Mary Katherine

10  Gonder[1,6,*]

11  [1]Department of Biology, Drexel University, Philadelphia, PA, USA

12  [2]Coriell Institute for Medical Research, Camden, NJ, USA

13  [3]School of Biomedical Engineering, Science and Health Systems, Drexel University,
14  Philadelphia, PA, USA

15  [4]Faculty of Arts and Sciences Informatics Group, Harvard University, Cambridge, MA,
16  USA

17  [5]Center for Tropical Research, Institute of Environment and Sustainability, University of
18  California, Los Angeles, CA, USA

19  [6]Department of Ecology and Conservation Biology, Texas A&M University, College
20  Station, TX, USA

21  [7]San Diego Zoo Wildlife Alliance, Escondido, CA, USA

22  [8]Cameroon Biodiversity Association, Douala, Cameroon

23  [9]Institut de Recherche pour le Développement (IRD), Maladies Infectieuses et Vecteurs:
24  Ecologie, Génétique, Evolution et Contrôle (MIVEGEC) (IRD 224-CNRS 5290-
25  Université de Montpellier), Montpellier, France

26  [10]School of Natural Sciences, University of Stirling, Stirling, UK

27    [11]Wildlife Conservation Society Cameroon, Yaoundé, Cameroon

28    [12]Departments of Genetics and Biology, University of Pennsylvania, Philadelphia, PA,
29    USA

30    [13]Department of Genome Sciences, University of Washington School of Medicine,
31    Seattle, WA, USA

32    [14]Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA

33    [15]Department of Biological Sciences, University of New Orleans, New Orleans, LA, USA

34    [16]Department of Ecology and Evolutionary Biology, University of California, Los
35         Angeles, CA, USA

36
37    *Matthew W. Mitchell
38    Email: mmitchell@coriell.org

39
40    *Mary Katherine Gonder
41    katy.gonder@ag.tamu.edu

42

2

# Abstract

Species evolve from populations with ancestor-descendant relationships in a bifurcating

process shaped by geography, gene flow, genetic drift, and natural selection leading to

local adaptation to prevailing environmental and ecological conditions. Building on this

foundational understanding, we explored local adaptation in chimpanzees (*Pan

troglodytes*) at a key geographical intersection in Cameroon where the two main

chimpanzee phylogenetic lineages converge. The Nigeria-Cameroon chimpanzee (*P. t.

ellioti*) and central chimpanzee (*P. t. troglodytes*) last shared a common ancestor about

500 thousand years ago, with occasional gene flow between them. The evolutionary

processes driving their prolonged separation are not fully understood, but neutral

evolutionary mechanisms alone cannot account for the observed divergence pattern.

Cameroon is often referred to as 'Africa in miniature' because the Gulf of Guinea Forest,

Congo Basin Forest, and savanna converge there, forming an ecotone. Thus, this

contact zone between subspecies in Cameroon provides a unique natural laboratory

that enabled us to investigate how environmental variation and natural selection shape

divergence in chimpanzees. We developed a genome-wide panel of single-nucleotide

polymorphisms (SNPs) in 112 wild chimpanzees sampled in multiple habitats across

this contact zone. We augmented SNP discovery by sequencing eight new chimpanzee

genomes from Cameroon and analyzing them with previously published chimpanzee

genomes. We found that *P. t. ellioti* and *P. t. troglodytes* diverged from one another

around 478,000 years ago and occasionally exchange migrants. We identified 1,690

unique SNPs across 905 genes associated with 31 environmental variables that

describe the habitat. These genes are involved in essential biological processes,

66   including immune response, neurological development, behavior, and dietary

67   adaptations. This study highlights the importance of understanding the geographical

68   context of natural selection, paving the way for future studies to interpret evidence for

69   genetic variation with phenotypic traits and deepening our understanding of how

70   populations diverge in response to environmental pressures.

71

# Author Summary

73   We investigated how local adaptation contributes to shaping the diversification of

74   chimpanzee subspecies at the geographical convergence point for the two major

75   branches of the chimpanzee phylogenetic tree. We analyzed genome-wide SNP

76   genotypes of 112 chimpanzees sampled from natural communities located in this

77   understudied area. We used tiered methods that identified 905 genes subject to

78   selection, each associated with one or more of 31 environmental predictors describing

79   the habitat. We found strong signals of selection in immune response genes that

80   separate *P. t. troglodytes* from *P. t. ellioti*, highlighting the important role of different

81   pathogen histories in their evolution. We also found evidence of selection in genes

82   associated with neurological development, behavior, and diet, that separate both the

83   subspecies and populations of *P. t. ellioti* that occupy different niches. These findings

84   suggest that ecological and cultural factors may also contribute to shaping the

85   diversification of chimpanzees across the contact zone.

86

# Introduction

Species consist of populations of reproductively compatible individuals with ancestor–descendant relationships that evolve through time [1]. Speciation may result from various factors. It may have a geographical dimension ranging from allopatry to sympatry with varying degrees of gene flow among populations, genetic drift, and natural selection [2]. Ecological factors often play a decisive role in this process through the local adaptation of populations to prevailing environmental conditions [3]. The fusion of genomics with ecological modeling has advanced the ability to identify loci under environmental selection. It contributes to understanding how species adapt to specific habitats and its impact on speciation [4, 5]. While this link has been studied in many taxa [6], it has been an especially strong focus in studies of human evolution. Human populations have adapted to a multitude of environments [5], disease landscapes [4, 7, 8], and diets including the ability to digest milk into adulthood [9], fatty acid digestion [10], foraging practices in tropical African rainforests [11], cereal-rich diets [12], and persistence in high-altitude environments [13-15].

By comparison, the factors that contribute to shaping the evolution of non-human great ape species are poorly understood. Genomic tools have contributed substantially to resolving the evolutionary relationships and histories of great ape species, subspecies, and some populations [16-20]. However, these studies generally assumed that neutral evolutionary processes (i.e., genetic drift) largely explain the partitioning of genetic variation in great apes. In particular, population genetic structure has been presented as evidence for allopatric speciation in 'Pleistocene Refugia,' among gorillas [21], isolation across conspicuous geographic boundaries like rivers [21-23], and

5

110    separation on different islands [20]. However, a growing body of evidence supports the

111    hypothesis that local adaptation due to natural selection occupies an essential role in

112    shaping the patterning of genetic variation and speciation in great apes [19, 24-26].

113         Among the great apes, chimpanzees (*Pan troglodytes*) have been particularly

114    well-studied, including analysis of genomes from a representative sample of captive

115    individuals [16, 18, 19] and population genetic studies of natural populations [22-24, 27].

116    The overall picture from these studies is that the species originated in western

117    equatorial Africa about 1mya. By 500kya in the Middle Pleistocene, two lineages began

118    to diverge from the ancestral *Pan* population: a western lineage composed of the

119    subspecies *P. t. verus* and *P. t. ellioti*, and a central/eastern lineage comprising *P. t.*

120    *troglodytes* and *P. t. schweinfurthii* (**Fig. 1a**). Major rivers, lakes, and the Dahomey Gap

121    are thought to have acted as dispersal barriers that separate the subspecies to different

122    degrees and timescales, potentially leading to allopatric speciation among chimpanzee

123    subspecies. Among these dispersal barriers, the Sanaga River in Cameroon stands out

124    (**Fig. 1b**). It separates the chimpanzee phylogenetic tree into its two main branches yet

125    remains permeable to occasional gene flow between *P. t. ellioti* and *P. t. troglodytes*

126    [18, 24, 27]. The Sanaga River has likely enabled some degree of allopatric divergence

127    due to genetic drift but the role that natural selection may have played in separating

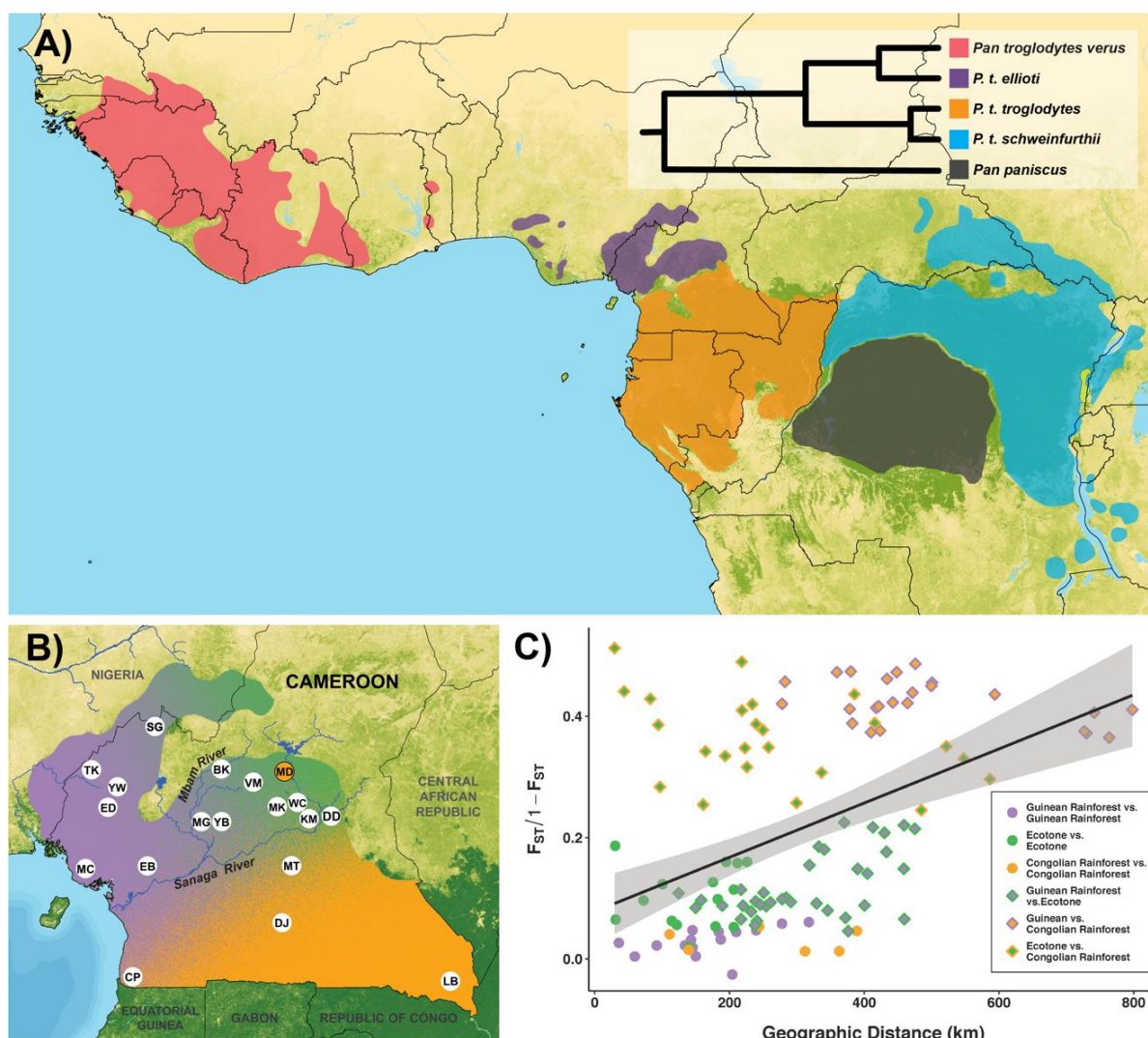128    these chimpanzee subspecies remains unknown.

**Fig 1. Chimpanzee evolutionary history across Africa and population structure in Cameroon.**
(A) Distribution and phylogeny of the genus *Pan*.
(B) Sampling locations of wild chimpanzee populations in Cameroon overlaid on spatial interpolation of population structure using SNPs from wild chimpanzees. The 'MD' sampling location is shaded orange to signify the presence of a *P. t. ellioti*/*P. t. troglodytes* F1 hybrid (CMMD06).
(C) Isolation-by-environment in wild chimpanzees in Cameroon. Correlation between 'linearized $F_{ST}$' and geographic distance (km) generated using SNPs from wild chimpanzees. Solid circles represent pairs of sampling locations from the same habitat. Dual-colored diamonds represent pairs of sampling locations from different habitats. Colors correspond to chimpanzee population origin: purple – *P. t. ellioti* (Rainforest), green – *P. t. ellioti* (Ecotone), and orange – *P. t. troglodytes*.

129

130          Natural selection has numerous opportunities to contribute to genetic divergence

131     that may vary between subspecies or populations in different habitats. Life history traits

132   and pathogen defense stand out as likely candidates for establishing among-population

133   divergence due to local adaptation Among these, the role of pathogens is best

134   understood. Differences in pathogen presence and prevalence have long been

135   associated with genotypic differences among great apes, especially chimpanzees. For

136   instance, wild chimpanzee populations are infected to different degrees with several

137   disease-causing pathogens, including malaria [28], Ebola [29], and viruses like simian

138   immunodeficiency virus (SIV) [30]. In the case of SIV and similar viruses, it is relatively

139   well established that these pathogens have exerted selective pressure on chimpanzees,

140   particularly the central and eastern subspecies [31-33]. Interestingly, Cameroon is a

141   unique disease landscape for chimpanzees, especially concerning the puzzling

142   distribution of SIVcpz. Unlike *P. t. troglodytes* and *P. t. schweinfurthii*, SIVcpz has not

143   been found in *P. t. ellioti* or *P. t. verus*, despite extensive sampling [34-36] (**Fig. 1a**).

144       Secondly, each chimpanzee subspecies occupies a distinct set of environmental

145   niches [37], creating opportunities for adaptation to local environmental conditions.

146   Although little is known about the links among genotypes, phenotypes, and

147   environmental conditions, chimpanzees in arid environments are more efficient in salt

148   removal than their counterparts in more humid forested environments [38]. However,

149   the role of local adaptation to specific environments remains largely unexplored yet is

150   perhaps the most intriguing avenue of investigation in their evolution. Chimpanzees, like

151   humans, have complex social systems and behaviors and maintain diverse cultural

152   traditions [39]. Similarly, cultural variation among chimpanzee communities may lead to

153   localized gene-culture co-evolution, potentially facilitating adaptation to diverse habitats

154   [37] that are vulnerable to human encroachment [40]. Habitat variation and resource

155    availability, specifically food types, are also known to affect chimpanzee socioecological

156    patterns directly [41], yet whether this variation translates into heritable genetic

157    differences remains speculative.

158        We investigated how local adaptation has influenced the evolution of

159    chimpanzees in Cameroon, a key region where the western and central/eastern

160    lineages of chimpanzees converge. Despite the wealth of research on the contributions

161    of neutral evolutionary processes to the genetic variation found in wild chimpanzees, the

162    contribution of natural selection remains a significant knowledge gap that our study

163    aimed to fill. We employed a two-tier approach to identify genic regions under selection

164    from a comprehensive analysis of natural chimpanzee communities sampled intensively

165    across Cameroon. First, we used whole-genome sequencing (WGS) data from 24

166    previously published chimpanzee genomes [16], along with eight newly sequenced

167    genomes of individuals from Cameroon to create and annotate a map of genomic

168    regions under natural selection from this expanded sample of complete genomes of

169    chimpanzees originating from Cameroon. Second, we used the analysis of this

170    expanded sample of genomes to create a genome-wide panel of ancestry-informative

171    putatively neutral SNPs, as well as SNPs that fell within signals of positive selection

172    (inferred with the WGS data) and, thus, were good candidates for performing tests to

173    assess local adaptation. We genotyped these SNPs in 112 wild chimpanzees sampled

174    across multiple habitats in Cameroon, encircling the contact zone between *P. t. ellioti*

175    and *P. t. troglodytes*, and that represent the diversity of habitats occupied by

176    chimpanzees across the contact zone [42], including the northern extent of the Congo

177    Basin Forest, the lowland and montane Gulf of Guinea Forest, and the forest/savanna

9

178   ecotone that bridges these two forest ecosystems. Finally, we used these SNP panels

179   to investigate the relationship between individual SNPs and a suite of SNPs

180   representing the genome to understand the relationship between allele frequencies and

181   environmental variability. Our objective was to assess whether environmental pressures

182   from differing ecologies have influenced allele frequency variation across these wild

183   populations.

184

185   # Results

186   We used WGS data from 24 previously published chimpanzee genomes [16],

187   along with eight newly sequenced genomes of individuals from Cameroon from the

188   Limbe Wildlife Center referred to hereafter as 'captive chimpanzees' (**Fig. S1** and **Table**

189   **S1**). We used the captive chimpanzee dataset and previously published data [16] to

190   create an annotated map of genomic regions under natural selection. Second, we used

191   a genome-wide panel of SNPs in 112 wild chimpanzees sampled across multiple

192   habitats across Cameroon (**Fig. 1b**) to develop a high-resolution, spatially explicit map

193   of allele frequencies to understand the link between habitat variation and loci under

194   selection.

195   ## Captive chimpanzee genome analysis and SNP discovery

196   ### Developing SNP datasets

197   We identified SNPs from 32 chimpanzee genomes across all four subspecies,

198   which included eight newly sequenced genomes from the contact zone between the

199   western and central/eastern chimpanzee lineages. After quality filtering, we retained

200   12,754,225 high-quality SNPs. Based on this initial whole-genome SNP set, two

201   datasets were created. The first dataset was thinned for linkage disequilibrium (LD),

202   retaining only SNPs with $r^2 \leq 0.1$, which resulted in 1,113,142 SNPs retained. The

203   second dataset was thinned to include only SNPs that followed our neutrality criteria

204   (**Methods**), resulting in 147,700 SNPs. **S1 Text** provides additional details on

205   heterozygosity (**Fig. S2**) and population cluster analyses (**Figs. S3 and S4**).

206   **Genome scans for signals of selection and defining genomic 'outlier'**

207   **regions**

208            To calculate a test statistic for cross-population extended haplotype

209   homozygosity (XP-EHH) and integrated haplotype score (iHS), SNP-based results were

210   summarized into windows following Pickrell et al. [43], but chromosomes were split into

211   100kb windows and SNPs were binned in 100 SNP increments. We merged windows

212   indicating positive selection for each method and population. The analysis identified

213   regions specific to the two lineages, and those shared among the Western and

214   Central/Eastern lineages were analyzed separately. **Table 1** summarizes outlier regions

215   for the XP-EHH and iHS and combined outlier tests. In the Western lineage, we found

216   335 outlier windows stretching 83.5 Mb with 695 candidate genes. The Central/Eastern

217   lineage had 318 windows stretching 81 Mb with 682 genes. We found 25 windows over

218   13.6 Mb with 80 candidate genes shared between lineages. We plotted the distribution

219   of the outlier regions on individual chromosomes (**Fig. 2a**).

220

**Table 1. Summary of captive chimpanzee whole genome "outlier" regions.**

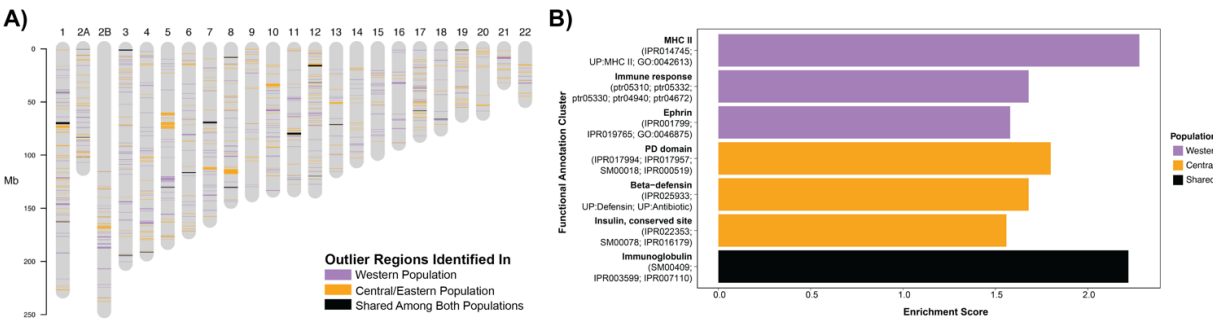| | XP-EHH | | iHS | | Combined | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Western lineage | Central/Eastern lineage | Western lineage | Central/Eastern lineage | Western lineage | Central/Eastern lineage | Shared |
| Windows found | 257 | 270 | 118 | 95 | 335 | 318 | 25 |
| Base pairs in windows | 53,753,675 | 55,149,110 | 34,000,000 | 32,200,000 | 83,453,675 | 81,049,110 | 13,600,085 |
| Protein coding genes used in enrichment analysis | 563 | 563 | 152 | 144 | 695/610 | 682/593 | 80/70 |



**Fig 2. Natural selection in chimpanzees.**
(A) Regions under selection found using captive chimpanzee genomes plotted on individual chromosomes.
(B) Functional enrichment clustering of genes and pathways under selection in chimpanzees. Enriched functional annotation clusters (based on genes in outlier regions) including their respective enrichment score. The name of one functional annotation of each cluster was taken to represent the complete cluster.

While all chromosomes are affected by selective sweeps some chromosomes show more regions under selection in one lineage or the other. The most extreme example in chimpanzees is chromosome 20, showing 6 times as many genetic regions under selection in the Central/Eastern lineage than in the Western lineage. Less extreme examples are found on chromosomes 8, 9, 13, and 19 with 2-fold more genome space showing evidence of selective sweeps in the Central/Eastern than the

233    Western lineage. In the Western lineage chromosomes 16, 15, 18, 3, and 11 show 5-,

234    4-, 3-, 2-, and 2-fold more genome space to be under selection than the Central/Eastern

235    lineage, respectively.

236        While the number of regions under selection in the Western as well as the

237    Central/Eastern lineage was equal on chromosomes 1, 2A, 2B, 4-8, 10-14, and 17,

238    there were some differences in the remaining chromosomes. Chromosomes 20 and 21

239    in the Central/Eastern lineage had five and four times more regions affected by

240    selective sweeps than the Western lineage. Chromosomes 9, 19, and 22 showed two

241    times more regions. In the Western lineage chromosomes 16 and 15 exhibited four and

242    three times more regions under selection than chromosomes in the Central/Eastern

243    lineage. Chromosomes 3 and 18 showed two times more regions under selection in the

244    Western lineage compared to the Central/Eastern lineage. There was no evidence for

245    selection shared between both lineages on chromosomes 2B, 9, 10, 14-16, and 20, 21-

246    22.

247    **Functional annotation, enrichment, and cluster analysis of outlier**

248    **regions under selection**

249        We analyzed annotated outlier regions with complete or partially overlapping

250    genes and other genetic content (*e.g.*, non-coding genes, pseudogenes). Both lineages

251    had significantly more protein-coding and non-coding genes than randomly sampled

252    genome regions (one sample t-test, $p=0.0001$ & $p=0.0001$) (**Table 2**). Additionally, the

253    number of non-coding genes (ncRNA) was also significantly higher (one sample t-test,

254    $p=0.0001$ & $p=0.0016$), while the number of pseudogenes showed no significant

255    differences (one sample t-test, $p=0.3023$ & $p=0.6518$) (**Table 2**). Closer inspection of

13

256    the most enriched regions (**Table S2**) revealed these contained mostly protein-coding

257    genes and ncRNAs. The region with the highest significance value in the Western

258    lineage carries exclusively ncRNAs and one window did not contain any annotated

259    genetic features at all.

260

261    **Table 2. Genetic content of 200 kb windows under selection and ten randomly**
262    **sampled genome regions.**
263

| POPULATION | Real Regions | Random Regions (n=10) | | p-value[a] |
|---|---|---|---|---|
| **WEST** | | Average | StDev[b] | |
| Protein coding genes | 500 | 400.8 | 31.84 | 0.0001 |
| ncRNAs | 178 | 141.1 | 18.60 | 0.0001 |
| Pseudogenes | 5 | 6.2 | 2.90 | 0.3023 |
| **CENTRAL/EAST** | | | | |
| Protein coding genes | 536 | 401.1 | 30.65 | 0.0001 |
| ncRNAs | 155 | 136.7 | 24.00 | 0.0016 |
| Pseudogenes | 5 | 6.4 | 2.72 | 0.6518 |

264
265    [a]p-values were obtained using the "one sample t-test".
266    [b]StDev: Standard Deviation
267

268        We examined enriched gene ontology (GO) terms in the 'Biological Processes'

269    category (**Table S3**) and enriched KEGG pathways (**Table S4**) for genes under

270    selection in the Western lineage, the Central/Eastern lineage, or shared between the

271    two populations. Genes significantly enriched in the Western lineage are involved in

272    developmental processes (hair follicle development, embryonic development, pattern

273    specification, melanocyte differentiation), cellular and metabolic processes, and protein

274    localization and degradation. Enriched KEGG pathways in the Western lineage were

275    mainly related to diseases caused by pathogens or internal dysfunctions, branched-

276    chain amino acids (BCAAs) degradation, and neurological development. The

277    Central/Eastern lineage genes are enriched for innate immune system response,

278    cellular processes, and wound healing. Enriched KEGG pathways in the

279    Central/Eastern lineage are involved in several diseases affecting the heart muscle and

280    Amoebiasis. The shared dataset showed enrichment only in bone mineralization without

281    any KEGG pathway.

282         To minimize annotation redundancy and clarify the biological functions in each

283    lineage, we grouped genes into functional clusters based on similar biological meaning,

284    not physical distance [44]. **Fig 2b** and **Table S5** show functional enrichment clusters of

285    genes that were unique to the western group (purple), unique to the central/eastern

286    group (orange) and shared between the western and central/eastern lineage.

287         We grouped 610 candidate genes from the Western lineage into three clusters.

288    The cluster with the highest enrichment score (ES = 2.3) included four genes (PATR-

289    DOB, PATR-DMB, MAMU-DMA, HLA-DOA) functionally associated with the Major

290    Histocompatibility Complex (MHC) II (**Fig. 3a**). MHC II genes, located on chromosome

291    6, play a crucial role in the adaptive immune response by activating CD4 T cells to

292    respond to extracellular pathogens [45]. The second cluster (ES = 1.8) contains the

293    same four genes as the first cluster, plus gene HLA-DQA1. This cluster is defined by

294    additional gene functions and displays enrichment in additional disease pathways active

295    in diseases like Asthma, Graft-versus-host disease, Allograft rejection, type I diabetes

296    mellitus, and the intestinal immune network for IgA production [45]. The third cluster (ES

297    1.58) contains three genes containing the Ephrin receptor-binding domain. These three

298    genes (EFNA4, EFNA3, EFNA1) form a gene cluster on chromosome 1 from position

299    133,320,040 to 133,391,332. Depending on the context, Eph signaling pathways are

300    key determinants of neurological development, cell morphogenesis, tissue patterning,

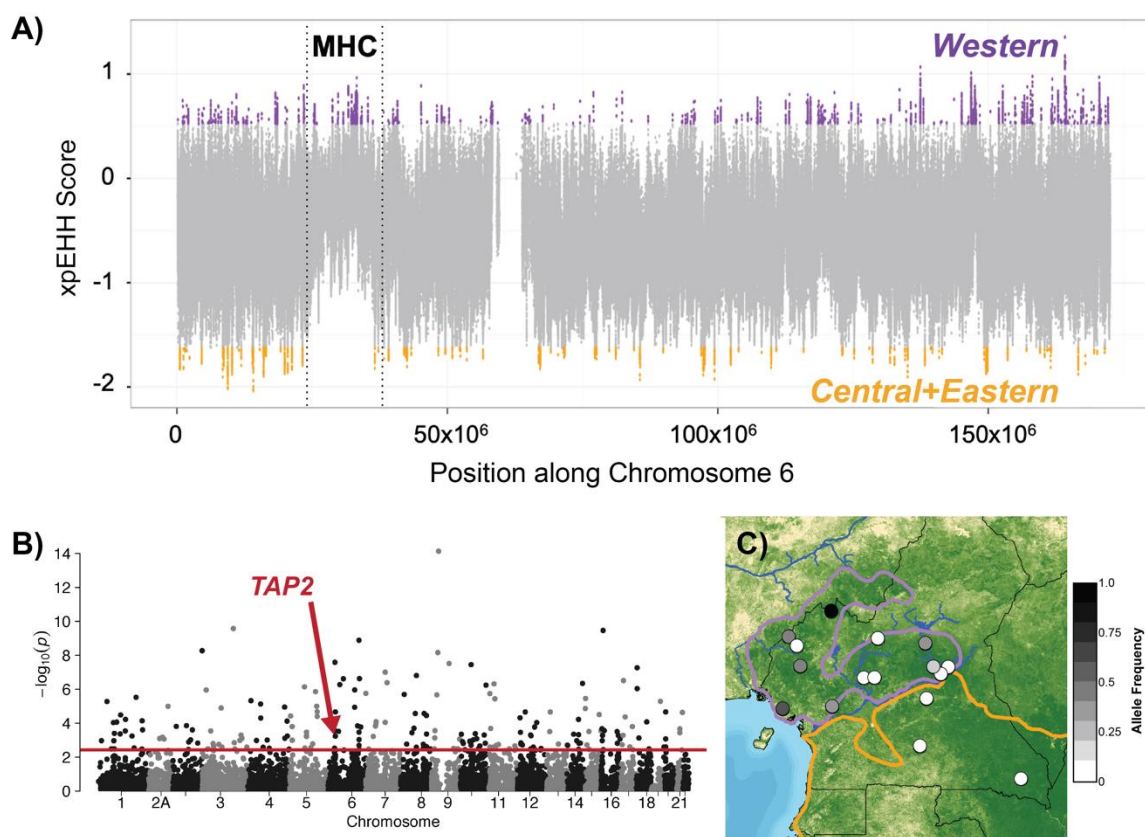301    angiogenesis, and neural plasticity [46, 47].

302



**Fig 3. Genome-wide variation of immune response genes under selection in chimpanzees.**
(A) XP-EHH analysis of SNPs on chromosome 6 from whole genome sequences of captive chimpanzees. Colored points represent SNPs within the 1% tail of the XP-EHH scores across the genome. The entire MHC region is notated, showing SNPs in MHC genes under selection in the Western lineage (*P. t. verus* and *P. t. ellioti*).
(B) Manhattan plot shows the genome-wide significance level (solid red line) for SNPs associated with Normalized Difference Vegetation Index (NDVI) - Brown with the *TAP2* SNP noted.
(C) Map of allele frequencies for the *TAP2* SNP superimposed onto NDVI and chimpanzee subspecies ranges in Cameroon.

303

304         In the Central/Eastern lineage, 593 genes were analyzed, forming three

305    functional clusters. The first (ES=1.8) showed enrichment of three genes (TFF3, TFF2,

306    TFF1) with a PD (or trefoil) domain. These three genes form a cluster on chromosome

307    21, but their functions are not understood: The peptides coded for in these segments

308    are in several tissues but are most abundant in the GI tract where they may stabilize the

309    mucosa and promote healing [48]. The second cluster (ES=1.7) contained five genes

310    belonging to the beta-defensin gene including DEFB126, DEFB127, DEFB129, and

311    DEFB132 are located on chromosome 20, and DEFB125 on chromosome 8. As

312    antimicrobial peptides they are important in the innate response, including resistance of

313    epithelial surfaces to microbial colonization and encapsulating viruses [49]. The third

314    cluster (ES=1.6) comprises genes INS, RLN3, and INSL6, all sharing an Insulin-like

315    domain.

316         Functional enrichment analysis of genes shared between both lineages revealed

317    only one cluster of six genes with an enrichment score of 2.2: IL1RL2 and IL18RAP

318    form a gene cluster on chr2A, CNTN6 is located on chr3, and ROBO3, ROBO4, and

319    HEPACAM form a gene cluster on chr11. These genes are all annotated with an

320    Immunoglobulin-like domain.

## 321 Wild chimpanzee SNP genotyping, population structure, and 322 selection analysis

### 323 Sequence analysis, filtering, SNP calling, and on-target read 324 assessment

325         We isolated DNA from fecal samples collected non-invasively from unhabituated

326    natural communities of chimpanzees sampled across Cameroon. For 192 of these

327    samples, we obtained 412,081,940 raw reads from single Illumina HiSeq PE125 lane –

328    an average of ~2.15 million reads per sample. In total, 275,443,720 of these reads

329    mapped to the chimpanzee reference genome; from these, we removed approximately

330    75 million reads and were left with 38,657,083 reads that mapped to our target sites

331    (**Fig. S5a**) – an average of 201,339 on-target reads per sample. The 9,986 targeted

332    sites had a mean read depth of 20x with one site showing as much as 166x coverage

333    (**Fig. S6**). After removing samples for missing data and relatedness, we were left with

334    two datasets ('10k' and '1k'). The '10k dataset' samples had significantly more on-target

335    reads per sample than the total dataset; an average of 328,863 on-target reads per

336    sample, representing ~16% of the total reads from these samples (**Figs. S5b** and **S5d**).

337    The '10k dataset' filtering process resulted in 7,878 SNPs and 112 samples, and all

338    samples from Boumba Bek (BB) and Campo Ma'an (CP) were removed. To retain more

339    geographic representation of samples from at least one of these sites, we created

340    another dataset ('1k dataset') by applying a more stringent site filter and the same

341    individual missingness filter above which resulted in 994 SNPs and 142 samples

342    (including two individuals from CP, but none from BB). The '1k dataset' samples had an

343    average of 268,773 on-target reads per sample, representing ~12% of the total reads

344    from these samples (**Figs. S5b** and **S5f**).

## Testing for isolation-by-distance and isolation-by-environment

346    We found that pairwise $F_{ST}$ values between sampling locations from different

347    habitats were significantly higher than pairwise $F_{ST}$ between sampling locations within

348    the same habitat for both the '10k' (one-tailed *t*-test, *p*-value = 2.2e-16; **Figs. 1c** and

349    **S7a**) and the '1k' dataset (one-tailed *t*-test, *p*-value = 2.1e-16; **Fig. S7b** and **S9a**).

350    Additionally, the geographic distance between sampling locations from different habitats

351    was significantly higher than between locations within the same habitat for all 19

352    sampling locations included in this study (one-tailed $t$-test, $p$-value = 6.324e-11).

353        We also performed a permutation test to account for the fact that population

354    structure across habitats can confound the detection of isolation-by-distance (IBD). This

355    categorized population pairs by geographic distance and randomized their habitat

356    origins, forming a null distribution of t-statistics. Using this distribution, we assessed if

357    $F_{ST}$ differed more between than within habitats/populations. For the '10k' and '1k'

358    datasets, we found that $F_{ST}$ was significantly higher between populations/habitats than

359    within them ($p$-value < 0.0001; **Figs. S8b** and **S9b**). This suggests that IBD alone

360    cannot fully explain the high $F_{ST}$ values between populations/habitats. We also ran the

361    permutation test for *P. t. ellioti* sampling locations alone and found that $F_{ST}$ is

362    significantly higher between *P. t. ellioti* (Rainforest) and *P. t. ellioti* (Ecotone) than within

363    them compared to the null distribution ($p$-value = 0.0002; **Fig. S10**). Taken together, the

364    results of the permutation tests suggest that habitat differences play a much stronger

365    role than geographic distance alone, although the signal is slightly stronger within *P. t.*

366    *ellioti* than in *P. t. troglodytes*. This may be attributed to the fact that *P. t. troglodytes* in

367    Cameroon occupies more uniform Congo Basin forested habitat south and east of the

368    Sanaga River. In contrast, *P. t. ellioti* occupies the comparatively diverse Gulf of Guinea

369    forest comprising lowland forest, montane forest, and the forest-savanna gradient north

370    of the Sanaga River.

371

## Population structure, hybridization, and demographic history

We next investigated population structure, hybridization, and demographic history. Principal Components Analysis (PCA) (**Figs. S11, S12, S13, S14** and **Table S10**), population clustering analysis results (**Figs. S15** and **S16**), and Analysis of Molecular Variance (AMOVA) (**Table S11**) consistently distinguished between *P. t. ellioti* and *P. t. troglodytes*. The results from wild chimpanzee samples were consistent with results from the genome analysis of captive individuals (**Figs. S17**, **S18**, **S19**, **S20** and **S21**), indicating that our SNP discovery approach from captive individuals is likely capturing pockets of genetic differentiation present in wild individuals. In addition, certain individuals showed hybrid ancestry, notably an F1 hybrid in the *P. t. ellioti* population and a potential backcrossed hybrid in *P. t. troglodytes*. The demographic history model indicates that *P. t. ellioti* and *P. t. troglodytes* split from one another around 478,000 years ago, with continuous but rare gene flow between them since splitting, underlining a complex demographic history characterized by significant admixture and evolutionary divergence within the region. **S1 Text** provides more detailed results from these analyses. Based on these results, we concluded that neither the IBD model nor simple allopatric divergence along the banks of the Sanaga River fully explains the separation of *P. t. ellioti* from *P. t. troglodytes*.

## Mapping wild chimpanzee genomic variation across habitats

These findings drew our attention to investigating how habitat variation corresponds with neutral and adaptive genetic differentiation among chimpanzees in Cameroon. Using a gradient forest model [50] and 31 environmental predictor variables

394    sourced from publicly available databases (See **Methods**, **SI Text,** and **Table S7**), we

395    quantified environmental associations with genomic loci, pinpointing key environmental

396    drivers and projecting genomic diversity spatially. We identified 581 SNPs with

397    significant environmental associations, representing 6% of all SNPs from wild

398    chimpanzees genotyped in this study. From these, 346 unique candidate genes within

399    10kb windows of these SNPs, matched outliers from captive chimpanzee genome

400    scans. When mapped to the study, these showed clear signals driven by a phylogenetic

401    split between *P. t. ellioti* and *P. t. troglodytes* across the Sanaga River <u>and</u> habitat

402    variation across Cameroon (**Fig. 4c**). Latitude (a proxy for geographic distance) had a

403    pronounced effect along PC1 (**Fig. 4d**). Isothermality and surface moisture also

404    contributed heavily to the model in differentiating between coastal and interior rainforest

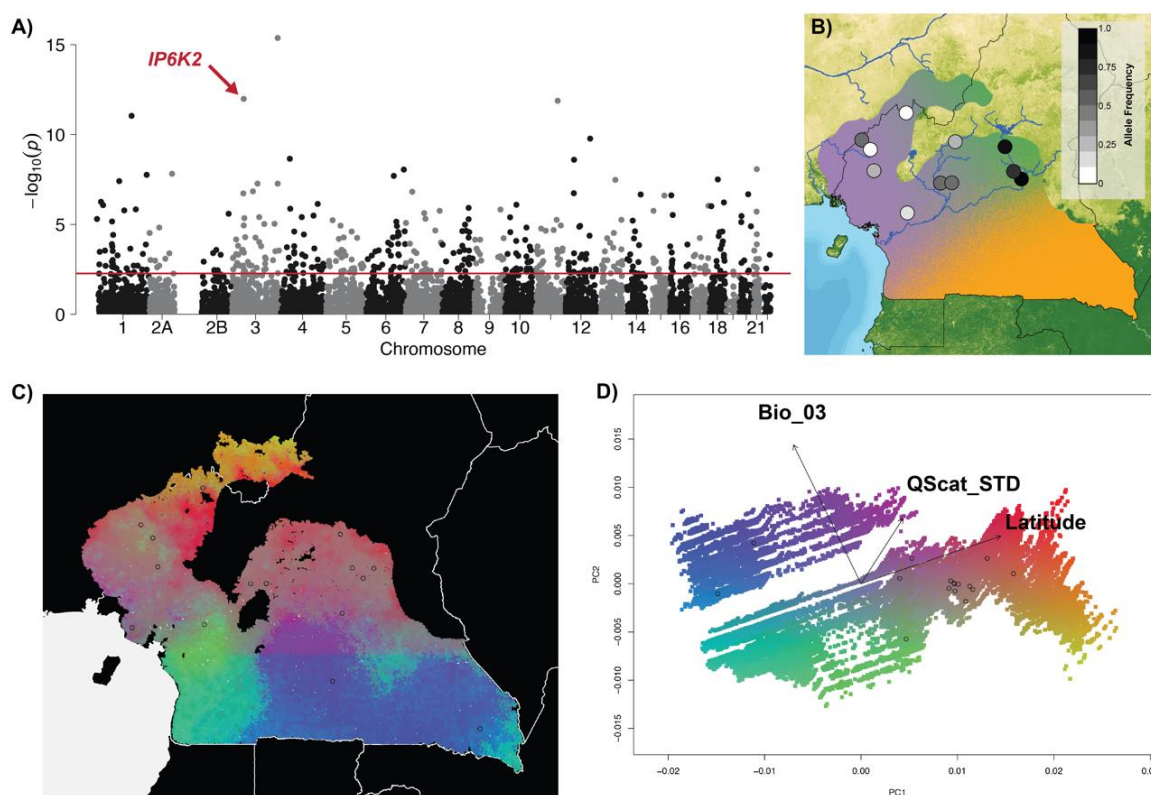405    habitats, as well as rainforest versus ecotone habitats (**Fig. 4c** and **4d**).

21

**Fig 4. Dietary gene under selection and gene-environment relationships**.
(A) Manhattan plot shows the genome-wide significance level (solid red line) for SNPs associated with Annual Mean Temperature (BIO1) with the *IP6K2* SNP noted.
(B) Spatialized allele frequencies for the *IP6K2* SNP showing differentiation between *P. t. ellioti* populations.
(C) Gradient forest-transformed climate variables show climate adaptation across the study area.
(D) Colors are based on a PCA of transformed climate variables.

406

407         Among all predictors tested in the model, latitude had the highest $R^2$ weighted

408 importance, likely reflecting the deep split between chimpanzee subspecies and/or

409 bioclimate turnover across the rainforest-savanna gradient. Precipitation during the dry

410 season and vegetation density were also important for predicting chimpanzee genomic

411 diversity (**Fig. S24**). The second most important axis of variation in the gradient forest

412 model primarily contributed to isothermality (bio3) and surface moisture (QScat_STD).

413 Thus, the variables contributing the most align with a rainforest/savanna ecotone split

414  (**Figs. 4c** and **4d**), consistent with previous studies of niche modeling [42].

## Detecting environmentally associated loci under selection in wild chimpanzee populations

417  We used Latent Factor Mixed Models (LFMM) to test for signals of selection on

418  individual SNPs in a manner that controls for confounding effects of population

419  structure. We identified 1,690 SNPs significantly associated with one of 31

420  environmental predictors (**Table S7**) after accounting for population structure (K=3). We

421  then identified 905 unique candidate genes within 10kb windows of the environmentally

422  associated outlier SNPs, all of which were outliers in the captive chimpanzee genomes

423  selection scan. Of the population groupings, we identified 695 associated with General

424  Temperature variables, 388 associated with Temperature Range, 66 associated with

425  Temperature Seasonality, 160 associated with Precipitation (Wet/Cold), 305 associated

426  with Precipitation (Dry/Warm), 456 associated with Surface Moisture, 448 associated

427  with Tree Cover, 355 associated with Vegetation Greenness, 325 associated with

428  Vegetation Brownness, and 428 associated with Topography. A simple Mantel test

429  revealed a significant correlation between pairs of environmental predictor variables and

430  shared outlier SNPs returned by LFMM between variable pairs (Mantel $r$ = 0.425, $p$ =

431  1.00e-6) demonstrating that independent LFMM performed as expected.

## Quantifying environmental relationships with candidate genes

433  We searched for functional enrichment signals in environment-associated genes

434  in two complementary ways. First, we compared the functions of genes associated with

435  environmental variation with genes that show no signs of positive selection. This

436  comparison helped us determine whether genes influenced by environmental factors

23

437 and potentially under selection differ functionally from those evolving under neutral

438 conditions. We used genes outside outlier regions from the captive chimpanzee whole-

439 genome analysis as a reference. We identified 47 biological processes enriched in

440 1,018 unique environmentally associated outlier genes from both the gradient forest and

441 LFMM models (**Table S12**). There were several enrichment clusters, notably two

442 processes functionally associated with immune response, one was related to the Major

443 Histocompatibility Complex (MHC) II – an important part of the adaptive immune system

444 - and eight processes associated with neurological development, including 60 unique

445 genes. We also found 48 enriched KEGG pathways in this subset of outliers (**Table**

446 **S13**). Key clusters included pathways in neurological development (56 genes), digestion

447 and metabolism (40 genes), and immune response (40 genes).

448 Our second analysis examined if genes influenced by environmental variation

449 showed functional enrichments compared to those under positive selection without a

450 clear environmental impact. This test uses a much smaller set of background genes

451 composed only of those assayed in wild chimpanzee SNP scan but were not

452 environmentally associated outliers in the LFMM and gradient forest models.

453 Unsurprisingly, the enrichment analysis using this more limited background set of genes

454 resulted in one significantly enriched biological process and KEGG pathway each, both

455 relating to neurological development, specifically axon guidance (**Table S14).**

456 Of the genes linked with immune response and MHC II, Transporter 2, ATP

457 binding cassette subfamily B member (*TAP2*) stands out. It contains a SNP that

458 significantly associated with Vegetation Brownness (NDVI BRN) ($-\log_{10}$ = 3.231762344,

459 p < 0.001) (**Fig. 3b**) and is associated with *General Temperature* variables in the LFMM

24

460 analysis of wild chimpanzees. The *TAP2* SNP in wild chimpanzees is nearly fixed in *P.*

461 *t. troglodytes* and is variable across *P. t. ellioti* habitats (**Fig. 3c**). Additionally, *TAP2* was

462 found to be under natural selection in the analysis of captive chimpanzee whole

463 genomes, and it is part of the enriched KEGG pathways under selection in *P. t. verus*

464 and *P. t. ellioti* (**Table S4**). *TAP2* is a component of the transporter associated with

465 antigen processing (TAP) complex, which plays a role in ensuring that MHC class I

466 (MHC-I) molecules are expressed on the cell surface [51]. TAP complex proteins,

467 including TAP2, are essential for viral peptide transport from the cytoplasm onto MHC-I

468 receptors within the endoplasmic reticulum [52]. In humans, several *TAP2* gene variants

469 are linked to an increased HIV-1 infection risk [53].

470       We identified the gene, Leucine rich repeat and Ig domain containing 2

471 (*LINGO2*), as having one of the strongest associations with the environmental predictor

472 variable, mean annual normalized vegetation index (NDVI), in the LFMM analysis (-$\log_{10}$

473 = 8.32330639, p = 0.00000000475) (**Fig. 5a**). A linear regression also revealed a

474 significant association between the allele frequencies of the *LINGO2* SNP and mean

475 annual NDVI ($R^2$ = 0.4826, p = 0.003495) (**Fig. 5b**). Low allele frequencies were found

476 in *P. t. ellioti* (Rainforest), with variable frequencies found in *P. t. ellioti* (Ecotone) (**Fig.**

477 **5c**). *LINGO2* is highly expressed in human brain tissue [54] and affects synapse

478 development and function [55]. *LINGO2* is also under positive selection in Lidia cattle

479 breed subpopulations and partially drives neurobehavioral phenotype variation among
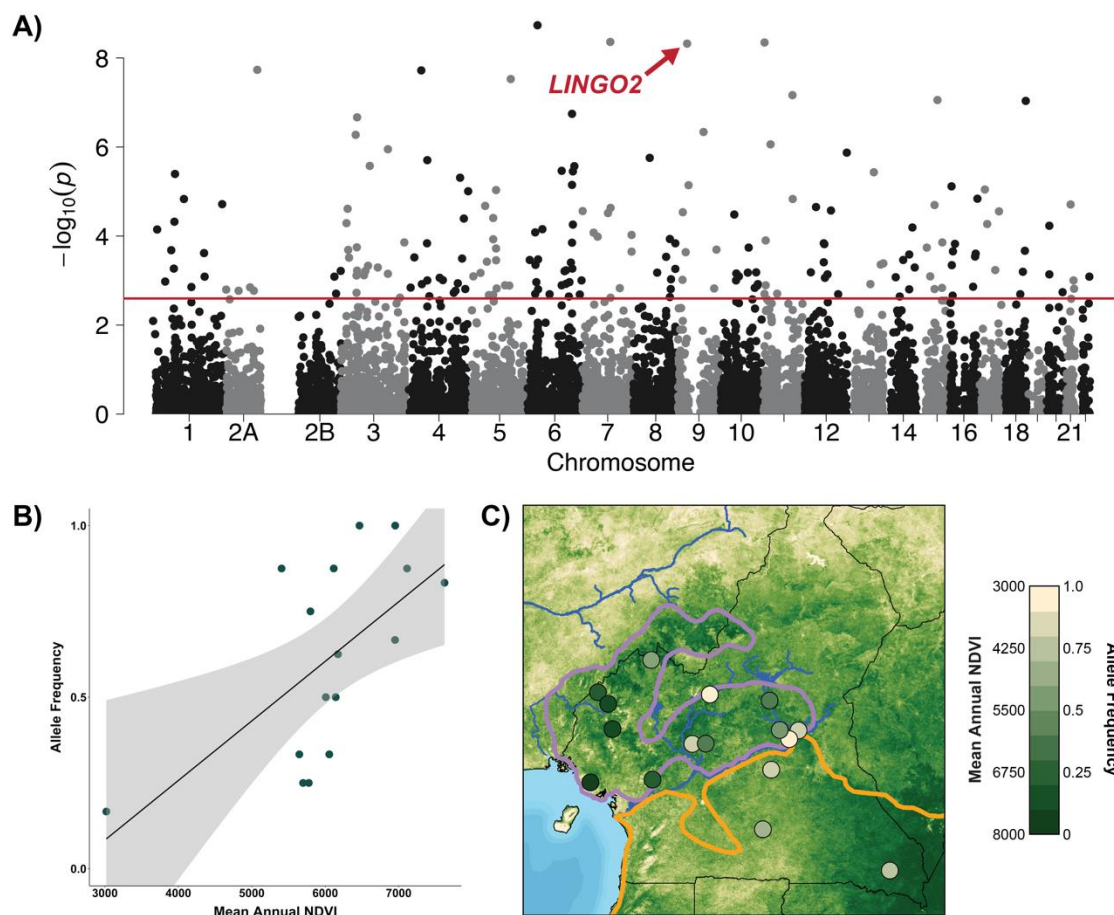
480 them [56].

25

**Fig 5. Genome-wide variation of neurological development genes under selection.**
(A) Manhattan plot shows the genome-wide significance level (solid red line) for SNPs association with mean annual normalized vegetation index (NDVI) with the *LINGO2* SNP noted.
(B) Correlation between *LINGO2* allele frequency and mean annual NDVI.
(C) Spatialization of allele frequencies for this SNP superimposed onto mean annual NDVI.

481

482        Of the 24 digestion and metabolism-related genes identified in the biological

483    processes and KEGG pathways, we further narrowed down our search by using

484    additional measures to quantify relationships of each of the genes with environmental

485    variables and were able to identify two genes with associated SNPs exhibiting

486    significant linear relationships directly with their associated environmental variables

487    across space, suggesting a potential role for diversifying selection across the

488    forest/savanna ecotone gradient. Acetyl-CoA acetyltransferase 2 (*ACAT2*) contains the

26

489     SNP at position 161,530,902 on chromosome 6 (**Fig. S26a**), which had the strongest

490     association of all SNPs with temperature seasonality according to results of the LFMM

491     analysis ($-\log_{10}$ = 3.594887672, p = 0.000254163) (**Fig. S26b**). Linear regression

492     revealed a strong and significant relationship between *ACAT2*'s outlier SNP and

493     temperature seasonality ($R^2$ = 0.5651, p = 0.0005) **(Fig. S26c)**. When plotting allele

494     frequencies of *ACAT2*'s outlier SNP, higher frequencies were observed in the ecotone's

495     northern sampling sites. Sampling sites within the range of the *P. t. troglodytes*

496     population had lower frequencies of the allele (**Fig. S26d**). The product of the *ACAT2*

497     gene is known to be involved in cholesterol and beta-oxidation lipid metabolism [57].

498         Another gene identified to have an environmentally associated outlier SNP was

499     Phospholipase C like 2 (*PLCL2*). *PLCL2* contains the SNP at position 17,268,745 on

500     chromosome 3 (**Fig. S27a**). We identified a strong relationship between this SNP and

501     the environmental variable *precipitation of the wettest month* through LFMM analysis ($-$

502     $\log_{10}$ = 3.36472742, p = 0.00043179) (**Fig. S27b**). Linear regression revealed a strong

503     and significant relationship between *PLCL2*'s outlier SNP and the environmental

504     variable *precipitation of the wettest month* ($R^2$ = 0.3422, p = 0.0102) **(Fig. S27c)**. We

505     observed higher allele frequencies of the *PLCL2* SNP in the *P. t. ellioti* ecotone

506     population, with the *P. t. ellioti* rainforest population having the lowest frequencies

507     across Cameroon (**Fig. S27d**). *PLCL2* is associated with obesity in mouse models.

508     Individuals lacking the allele were shown to have a leaner phenotype; were able to

509     resist induced obesity due to increased protection from glucose metabolism disorders

510     and insulin resistance; and exhibited higher energy expenditure [58].

511         Finally, a SNP in the Inositol hexakisphosphate kinase 2 (*IP6K2*) gene was

512 identified as a significant outlier differentiating the *P. t. ellioti* ecotone and rainforest

513 populations and associated with the environmental variable *Mean Annual Temperature*

514 through LFMM analysis (-$\log_{10}$ = 11.98296666, p < 0.0001) (**Fig. 4a**). *IP6K2*'s SNP was

515 significantly more divergent than neutral SNPs between the two *P. t. ellioti* populations

516 ($F_{ST}$ = 0.49, p < 0.0001) (**Fig. 4b**). In humans, the *IP6K2* gene is linked with

517 inflammatory bowel disease and cellular response to flavonoids, plant metabolites found

518 in fruits and vegetables [59]. The human KEGG pathway containing *IP6K2* is associated

519 with VACTERL/VATER syndrome, often associated with congenital heart disease and

520 chondrodysplasia [60, 61].

521

## Discussion

523     We presented genome-wide SNP genotyping from a representative sample of

524 112 wild chimpanzees from across Cameroon, along with eight newly sequenced

525 genomes of captive chimpanzees to enhance SNP discovery. We supplemented these

526 data by combining genetic analyses with environmental association scans to search for

527 evidence of environmentally-mediated selection. While prior studies have largely

528 concentrated on neutral evolution mechanisms across this contact zone between

529 chimpanzee lineages, our findings support a role for diversifying selection in the

530 divergence of chimpanzee subspecies across different environments. The proposed

531 population history of chimpanzees across this contact zone is well supported in this

532 study. *P. t. ellioti* and *P. t. troglodytes* last shared a common ancestor around 478,000

533 years ago, with occasional gene flow between them evidenced by an F1 hybrid in *P. t.*

534    *ellioti* and a potential backcrossed hybrid in *P. t. troglodytes*. These findings support

535    prior studies suggesting that this contact zone between subspecies best fits an

536    isolation-with-migration population model in which allopatric divergence and positive

537    selection contribute to the partitioning of genetic variation [62].

538        The evidence supporting a role for environmentally-mediated selection across

539    this contact zone is also compelling. We found 1,690 unique SNPs were associated with

540    at least one of 31 environmental predictors, indicating that prevailing environmental

541    conditions contribute to local adaptation in *P. t. elloti* and *P. t. troglodytes*, and to a

542    lesser extent, among populations within *P. t. ellioti*. These SNPs are distributed among

543    905 outlier genes enriched for 48 biological processes. Overall, the sets of genes with

544    highly divergent allele frequencies that separate *P. t. ellioti* from *P. t. troglodytes*

545    suggest a role for selection in pathways important in two main categories: immune

546    response and life history traits (neurological development, behavior, and dietary

547    function).

548        It is important to reiterate that all outliers identified in wild chimpanzees using

549    LFMM-based approaches were also identified as outliers in the haplotype homozygosity

550    selection scans of captive chimpanzee genomes. This two-tiered approach offers

551    heightened reliability of selection scans in wild populations while mitigating the

552    incidence of false positives in our final dataset. Moreover, the congruence of these

553    identified genomic regions between the two methods and two complementary datasets

554    suggests that these outliers are subject to positive selection and not merely an artifact

555    of demographic history or neutral population genetic structure.

556

## Signatures of selection from variable pathogen histories

557

558         The lack of natural SIVcpz infection in *P. t. ellioti* sparks interest because it

559 persists despite opportunities for transmission. SIVcpz*Ptt* virus infects *P. t. troglodytes*,

560 crossed the species barrier on at least four occasions: from chimpanzees to humans in

561 southern Cameroon, giving rise to the HIV-1 group M pandemic and to HIV-1 group N

562 [63, 64]. HIV-1 group O and P also arose from transmission from chimpanzees to

563 gorillas before subsequent transmission to humans [65, 66]. Thus, SIVcpz can cross

564 genus boundaries which makes its absence in *P. t. ellioti* particularly striking since this

565 subspecies still exchanges occasional migrants with *P. t. troglodytes*. Finally, the

566 presence of prey primate species that harbor endemic SIV strains also creates multiple

567 pathways for cross-species transmissions [34, 36, 67, 68]. Thus, we speculate that the

568 absence of an SIVcpz in *P. t. ellioti* must be at least partially explained by adaptations

569 that interrupt SIVcpz cell entry and/or boost immune response to clear SIVcpz infection.

570         Four processes functionally associated with the Major Histocompatibility

571 Complex (MHC) II on chromosome 6 play a crucial role in the adaptive immune

572 response. MHC II peptides stimulate CD4+ T cells that activate downstream immune

573 responses to intracellular pathogens, including viruses. In particular, the Th1/Th2 cell

574 differentiation pathway determines the type of helper cell a CD4+ T cell will become.

575 Naïve CD4+ T cells recognize an MHC class II molecule, activate, and divide to produce

576 clone effector CD4+ T cells specific for a particular antigen. CD4+ T cells can

577 differentiate into T helper type-1 (Th1), T helper type-2 (Th2), or other T helper types,

578 each with distinct cytokine-secretion phenotypes, production of distinct interferons, and

579 different downstream immune responses. This finding corresponds well with a growing

580  body of evidence that positive selection associated with pathogen defenses has

581  contributed to the genetic and phenotypic differentiation of chimpanzee subspecies,

582  especially *P. t. troglodytes* and *P. t schweinfurthii*, which is perhaps due to exposure to

583  different viruses [32].

584  This finding naturally called our attention to the absence of SIVcpz in *P. t. ellioti*

585  attributed to a lack of gene flow between *P. t. ellioti* and *P. t. troglodytes* [30]. Given that

586  gene flow occurs between *P. t. ellioti* and *P. t. troglodytes*, and that SIVcpz*ptt P. t.*

587  *troglodytes* is the source of multiple cross-species infections in both humans and

588  gorillas, it is logical to assume that SIVcpz should naturally infect *P. t. ellioti*. We

589  observed evidence of positive selection in *P. t. ellioti* due to highly differentiated SNPs

590  enriched in genic sites. Among these, *TAP2* (**Fig. 3c**) variants increase the risk of HIV

591  infection in humans [53], and may have a similar function in chimpanzees. Given the

592  low level of gene flow, and the absence of sequence data upstream or downstream of

593  the *TAP2* in our data, we cannot conclude whether this is evidence for recent adaptation

594  to SIVcpz*ptt* or evidence of ancient selection resulting from exposure to SIV-like viruses.

595  Evidence is mounting that chimpanzees have had a long and continuing relationship

596  with SIV-like viruses such that differences in viral exposures and immune responses

597  have likely been a central feature of the evolution of chimpanzees [32, 33, 69].

598  For instance, *P. t. troglodytes* and *P. t. schweinfurthii* also show evidence of

599  recent positive selection in genes involved in SIV/HIV cell entry and immune response

600  to SIV, biological pathways responsible for T-helper cell differentiation, including CD4

601  [33], and multiple genes that SIV/HIV use to infect and control host cells including

602  CCR3, CCR9 and CXCR6 [32]. There is also compelling evidence of past selective

603 sweeps leading to reduced diversity in the MHC II repertoire of *P. t. verus* that has been

604 attributed to past infections with SIV or SIV-like viruses [70]. Although we cannot

605 speculate further given the nature of the data in this study, our findings add to the

606 mounting evidence that chimpanzees have experienced long-lasting host-virus

607 relationships with SIV-like viruses and that these relationships have been a critical

608 process underpinning their evolution. More detailed investigations are needed on

609 whether the positive selection in *P. t. ellioti* is due to past, recent, or ongoing infection

610 with SIVcpz and/or related viruses.

## Signatures of selection across variable habitats

612 Cameroon is also a uniquely positioned 'natural laboratory' to examine the

613 relative contributions of neutral evolutionary forces versus natural selection in the

614 evolution of many animals, including chimpanzees. In addition to being home to the

615 Sanaga River, a well-known biogeographic boundary for many species, the country is

616 exceptionally ecologically diverse. We speculate that the area is important for

617 understanding how habitat variation and behavioral diversity may impact chimpanzee

618 evolution. The Congo Basin Forest in the south, the Gulf of Guinea Forest in the west,

619 and the Sahelian habitats in the north of Cameroon all converge and interdigitate to

620 form a unique ecotone habitat composed of open woodland, savanna, and riparian

621 forest [71]. This ecotone is a known engine of diversification for many species [72-78].

622 Differences between *P. t. ellioti* and *P. t. troglodytes* have been linked with habitat

623 variation across Cameroon [62], which suggests a possible role of local adaptation in

624 their genetic differentiation. Finally, there is a further genetic distinction within *P. t. ellioti*

625 itself, with one gene pool associated with the mountainous rainforest in western

32

626    Cameroon and the other with the ecotone in central Cameroon [24] (**Fig. 1b**). Each

627    gene pool has a unique ecological niche [42, 79] with marked differences in key

628    socioecological variables, including sex-specific differences in community structure [80]

629    and dietary preferences [81].

630        We observed compelling evidence for positive selection that distinguished *P. t.*

631    *troglodtyes* from *P. t. ellioti* across this contact zone. We also found evidence of

632    diversifying selection that distinguished *P. t. ellioti* populations that occupy different

633    niches [42], which adds strength to our previous findings that both allopatric divergence

634    due to genetic drift and environmentally-mediated local adaptation contribute to

635    sustaining the prolonged separation of these two subspecies across this narrow contact

636    zone between them. In particular, we found 246 genes involved in cellular, metabolic,

637    and developmental processes were associated with one or more of the 31

638    environmental predictor variables. Genes with the most divergent allele frequencies

639    were associated with latitudinal variation and separate *P. t. ellioti* from *P. t. troglodytes*.

640    We detected an additional more subtle signal of positive selection among *P. t. ellioti*

641    chimpanzees located in western Cameroon's mountainous, forested regions compared

642    to the population inhabiting central Cameroon's drier ecotone forests. Environmentally

643    driven pressures between habitats shape adaptive variation, especially between

644    rainforest and ecotone habitats. Chimpanzees in these two different habitat types were

645    previously identified as distinctive ecological populations occupying unique niches [42,

646    79]. The two ecological populations also exhibit distinctive differences in diet and

647    nesting preferences [81, 82], key elements of chimpanzee cultural diversity.

648     Thus, multiple phenotypic axes appear to be under environmentally mediated

649     selection that can be linked to habitat variation and variation in chimpanzee

650     socioecology. Genes under environmentally mediated selection associated with

651     neurological development (e.g., *LINGO2*) could be shaped by selective pressures linked

652     to the development of cultural traits in diverse habitats [37], while those associated with

653     diet and metabolism are likely shaped by pressures related to fruit availability and

654     seasonality [81]. One of the most compelling findings of our study is the identification of

655     24 genes related to digestion and metabolism with the strongest signals of

656     environmentally-mediated selection, including *ACAT2*, *PLCL2*, and *IP6K2*, which

657     present promising avenues for future research.

658     This study adds to the emerging evidence that neutral evolutionary forces alone

659     cannot explain the prolonged persistence of the separation of *P. t. ellioti* from *P. t.*

660     *troglodytes* across the narrow contact zone between them. Local adaptation to

661     prevailing conditions has led to divergence in sets of genes important in immune

662     response, neurological development, behavior, and dietary function. Together, these

663     findings suggest that local adaptation, notably to varying pathogen pressure and

664     different habitat types, has shaped chimpanzee subspecies differentiation in Cameroon,

665     and likely across their broad range. Future studies exploring how these genetic

666     differences map to phenotypic differences in wild populations are needed to better

667     understand precisely which traits — particularly those associated with pathogen

668     defense, diet, social organization, and other aspects of chimpanzee cultural diversity —

669     provide for local adaptation and divergence among chimpanzee populations across this

670     region.

671

# Materials and Methods

## Captive chimpanzee genomes

### Sequencing and read mapping

675     Our captive chimpanzee genomic dataset includes 24 previously sequenced [16]

676  and 8 new chimpanzee genomes from Cameroon, representing all four subspecies: 4 *P.*

677  *t. verus*, 10 *P. t. ellioti*, 12 *P. t. troglodytes*, and 6 *P. t. schweinfurthii* (**Table S1**). These

678  eight genomes were sequenced using established methods [16] and deposited in

679  GenBank. Details on the samples, the estimated origins of the captive chimpanzees

680  [83], and GenBank accession numbers are in **Fig S1** and **Table S1**. We mapped raw

681  sequencing reads against the chimpanzee reference genome Pan_troglodytes-2.1.4

682  (*panTro4*; https://www.ncbi.nlm.nih.gov/assembly/GCF_000001515.6/) using BWA-

683  MEM v0.7.12 [84] with default parameters. After removing PCR duplicates using

684  PICARD v1.119 (https://broadinstitute.github.io/picard/index.html), we called variants

685  using FREEBAYES v0.9.20 [85]. After filtering, 12,754,225 high-quality bi-allelic SNPs

686  on the autosomes were retained.

### Genome scans for signals of selection

688     We divided SNP datasets into a 'Western lineage' (*P. t. verus* & *P. t. ellioti*; *n*=15)

689  and 'Central/Eastern lineage' (*P. t. troglodytes* & *P. t. schweinfurthii*; *n*=17). We applied

690  two selection scan methods, cross-population extended haplotype homozygosity (XP-

691  EHH) [86] and integrated haplotype score (iHS) [87] to detect sweeps using hapbin

692  v1.2.0 [88]. Since both tests require haplotypes, we phased the whole-genome SNP

35

693  datasets (12,754,225 SNPs) with SHAPEIT v2.r837 [89] following established methods

694  [18]. Genetic maps for *panTro4* were provided by de Manuel & Kuhlwilm *et al.* [18] and

695  Auton & Feldel-Alon *et al.* [90]. iHS calculations used SNPs with a minor allele

696  frequency (MAF) over 5%. We determined the ancestral state of each allele using the 6-

697  primate EPO alignment (ftp://ftp.ensembl.org/pub/release-80/fasta/ancestral_alleles/)

698  [91, 92]. After phasing, ancestral allele assignment, and MAF filtering we used

699  4,577,055 SNPs in the Western lineage and 6,475,338 SNPs in the Central/Eastern

700  lineage for iHS. XP-EHH scores compared both lineages using 12,450,633 SNPs, and

701  the results were normalized across the genome.

## Defining population-informative neutral SNPs

703  Using normalized XP-EHH and iHS values, we identified SNPs expected to

704  follow neutral evolution that met the following criteria: (*i*) a *p*-value of > 0.05, (*ii*) absent

705  from the top 1% genomic regions under selection (see *Defining genomic 'outlier' regions*

706  *in captive chimpanzees*), (*iii*) be located >10kb from a gene, and (*iv*) be in linkage

707  equilibrium. Using these parameters, we defined 147,700 neutral SNPs reflecting

708  chimpanzee population structure. We annotated these using the Variant Effect Predictor

709  (VEP) v82 and the UpDownDistance plugin [93].

## Defining genomic 'outlier' regions

711  To understand the amount of selection on the genome, we considered numbers

712  of base pairs under selection (magnitude) and the size of regions affected (genome

713  space). We employed iHS [87] and XP-EHH [86] to detect signatures of local positive

714  selection. Both assume that long-range haplotypes remain unaffected by recombination,

715  signifying natural selection even with small datasets [43, 94]. They are also

36

716 complementary: while iHS detects partial sweeps, XP-EHH identifies near-fixation

717 events. Following Pickrell *et al.* [43], chromosomes were split into 100kb non-

718 overlapping windows, and the fraction of SNPs with |iHS| > 2 and the maximum XP-

719 EHH was used as a test statistic. We analyzed the fraction of SNPs with |iHS| > 2 and

720 the maximum XP-EHH per window. We turned these into empirical *p*-values by binning

721 windows by SNP count, with iHS dropping windows with < 100 SNPs. Each window's

722 statistic value was compared against others in its bin to determine an empirical p-value.

723 All bins were then sorted by this *p*-value. The top 1% of each test statistic was noted.

724 'Outlier regions' were windows in this 1% (*p*-value < 0.01). Adjacent windows were

725 merged, retaining the smallest *p*-value.

## Characterizing genomic regions under selection

727 As XP-EHH and iHS are complementary, we analyzed the 1% tail of each test

728 merging adjacent windows. Windows were extended by 50kb on either side and

729 annotated for gene content using Ensembl's BioMart [95], including protein-coding

730 genes, pseudogenes, and RNA coding genes. Genes within outlier regions were

731 considered candidate genes. We tested whether these genomic regions carry certain

732 types of gene content more often than expected by chance by randomly selecting

733 regions equivalent in length and annotating them as described above. We repeated this

734 process 10 times for the Western and Central/Eastern populations, respectively. We

735 counted the number of protein coding genes, non-coding RNAs, and pseudogenes in

736 the real and randomized datasets. We then calculated mean and standard deviation for

737 each and performed a one-sample t-test to determine significance.

## Functional annotation and enrichment analysis of whole-genome datasets

We used DAVID Bioinformatics Resources v6.8 [44] to annotate candidate genes and perform an enrichment analysis with default functional annotations (GO terms, KEGG pathways, protein domains). We concentrated on the 'functional annotation clustering' function using the highest classification stringency and adjusted the enrichment thresholds for EASE to 0.05, reducing non-significant term inclusion. This clustering reduces redundancy by grouping similar annotations. Clusters received a Group Enrichment Score based on their *p*-value, ranking their biological importance. High scores likely mean lower *p*-values for annotation members [44]. We omitted windows found in both Western and Central/Eastern lineages, analyzing them separately. We set a background gene population as the entire chimpanzee genome, as recommended for genome-wide studies [44].

# Wild chimpanzee SNP genotyping, population history, and selection analysis

## Fecal sample collection, DNA extraction, and quantification

We sampled wild, non-habituated chimpanzee populations using non-invasive methods during a series of field studies from 2003 to 2015 spanning remote forested regions of Cameroon (**Fig. 1b**). Sampling occurred in protected and unprotected areas, as detailed in **Table S6.** Chimpanzee fecal samples were collected and stored following established protocols [24]. All samples were transported from Cameroon to the United States in full compliance with the Convention of International Trade in Endangered

760 Species of Wild Fauna and Flora (CITES), the Centers for Disease Control (CDC)

761 export and import regulations, and with approval from the Government of Cameroon.

762 Following established protocols [24], DNA was extracted from fecal samples with

763 the QIAamp DNA Stool Mini Kit (Qiagen). Due to the low proportion of endogenous

764 DNA in fecal gDNA extracts [96, 97], samples were sometimes extracted up to six times

765 to ensure enough chimpanzee DNA for sequencing processes. The concentration of

766 endogenous DNA was measured via quantitative real-time PCR using the Quantifiler™

767 Human DNA Quantification Kit (Applied Biosystems) and methods from prior studies

768 [24].

769 **SNP ascertainment, library preparation, DNA capture enrichment, and**
770 **sequencing**

771 We genotyped 9,986 SNPs of wild chimpanzees from Cameroon, chosen from

772 the larger set of 12,450,633 SNPs identified in the captive chimpanzee genome dataset.

773 This selection comprised: (i) population informative neutral SNPs (n=3,492) randomly

774 selected 147,700 neutral SNPs defined in the whole-genome dataset from above; (ii)

775 'outlier' SNPs (n=6,494) identified through iHS or xpEHH tests as being in the top 1%

776 for selection signals and within or 10k bp up- or down- stream of a known gene; and, (iii)

777 SNPs in genes involved in immune response, disease resistance, and dietary

778 adaptation in humans (n=20) [98, 99]. For each targeted SNP, we designed two 80

779 nucleotide biotinylated RNA probes, overlapping by 20bp, to create 100bp windows

780 around each SNP using the panTro4 chimpanzee reference genome. After rigorous

781 filtration using the Arbor Biosciences BLAST pipeline, we finalized a bait-set of 19,974

782 probes, assigning one or two probes to each SNP based on the outcome of the

783 stringent filtering process.

39

784          gDNA samples were prepared in clean facilities at Arbor Biosciences to prevent

785    contamination. DNA was quantified, sonicated, and size-selected for around 300nt

786    fragment lengths. Samples were converted to sequencing libraries via adapter ligation.

787    They were index-amplified based on the DNA input amount. Up to 2µg of each library

788    was then enriched using the myBaits system v3. Different enrichment and amplification

789    protocols were applied depending on the DNA quantity in the starting extract. Libraries

790    were combined for equal representation, sequenced on an Illumina HiSeq PE125 lane

791    at HudsonAlpha, and protocols were consistent with studies on degraded or low

792    endogenous DNA samples [96, 97].

## SNP calling and on-target read assessment

794         We filtered sequence reads with the Illumina CASAVA-1.8 FASTQ Filter

795    (http://cancan.cshl.edu/labmembers/gordon/fastq_illumina_filter/) and mapped them to

796    the chimpanzee panTro4 genome using BWA-MEM. After removing PCR duplicates

797    with PICARD, we called variants using FREEBAYES. We evaluated DNA capture

798    enrichment and sequencing for the raw and filtered sequence reads using SAMtools

799    v1.3.1 (101) and VCFtools v0.1.15 [100] following established methods (97) for all 192

800    individuals. Using VCFtools, we filtered variant calls based on quality and coverage.

801    SNPs with <5x coverage or quality score <30 were recorded as missing data [97].

802    Positions with >60% missing data, a minor allele frequency below 5%, or individuals

803    with >70% missing data were removed. This resulted in 7,878 SNPs from 112 samples,

804    termed the '10k dataset'. Due to removing all Boumba Bek (BB) and Campo Ma'an (CP)

805    samples, a second '1k dataset' was made with stricter site filtering, yielding 994 SNPs

806    and 142 samples, which included two from CP but none from BB. We also removed

807    closely related and duplicate samples using the R package related [101] using the

808    triadic likelihood method [102], resulting in 85 individuals in the '10k dataset' and 108

809    individuals in the '1k dataset.'

810 **Testing for isolation-by-environment and inferring population**

811 **structure**

812    **S1 Text** provides full details on methods to test IDB, IBE and to infer population

813    structure, hybridization, and demographic history. In brief, we examined IDB versus IBE

814    using pairwise $F_{ST}$ values between sampling locations using Arlequin v3.5 [103], while

815    geographic distances were determined with the geosphere package in R [104], focusing

816    on areas with more than one individual. Population structure was inferred by PCA and

817    ADMIXTURE analysis [105], DISTRUCT v1.1 [106], CLUMPP v1.1.2 [107], with

818    geographic without pre-assigned population labels using the SNPrelate package in R,

819    focusing on 'neutral' SNPs identified from captive chimpanzee genomes. We mapped

820    genetic clusters using TESS [108] and Ad-Mixer v1.0 [109], accounting for IDB. We

821    calculated observed and expected heterozygosity using the adegenet package in R

822    [110], identified potential hybrids using NEWHYBRIDS v1.0 [111] as implemented in the

823    R packages *hybriddetective* [112] and *parallelnewhybrid* [113], and investigated

824    demographic history using δaδi [114] to model asymmetric migration patterns between

825    *P. t. ellioti* and *P. t. troglodytes.*

826    Environmental data layers (**Table S7**) were compiled and analyzed to assess

827    habitat suitability and IBE for chimpanzees in Cameroon and Nigeria. These layers,

828    sourced from publicly available databases, included diverse variables such as

829    topography, hydrography, climate, vegetation, moisture content, and tree cover. After

830     standardizing these layers to a 30-arcsecond resolution and converting them to the

831     WGS84 coordinate system, the dataset underwent cross-correlation analysis to pinpoint

832     environmental factors significantly influencing chimpanzee distribution (**Table S8**).

## Mapping genomic variation across habitats

834         Using the R package *gradientForest* [50], we calculated associations between

835     allele frequencies and environmental variation across suitable habitat. This extended

836     random forest model identifies links between response variables (e.g., SNP allele

837     frequencies) and spatial environmental factors [115] by iteratively processing datasets,

838     assessing outliers and predictor significance. Gradient forests further apply regression

839     to multiple responses, revealing genomic variation from environmental shifts. This can

840     pinpoint areas of high intraspecific variation, subspecies transitions, or barriers

841     separating genomic variation related to the environment [116]. Following established

842     methods, we refined the environmental dataset (**Table S7**) to reduce noise and applied

843     the gradient forest model [76].

844         We ran gradient forests on 7,878 SNP allele frequencies used as a response

845     dataset, and 17 environmental variables as the predictors in our final model, including

846     measures of temperature, precipitation, vegetation, surface moisture, and geographic

847     features at the sampling locations. We ran 100 trees in our model, noting SNPs

848     significantly associated with any environmental variable ($R^2 > 0$) and the average

849     regression of all associated SNPs. To assess model performance, we randomized the

850     environmental data and ran 200 permutations of the model, creating a distribution of $R^2$

851  and significant SNP associations. We then ran 200 permutations of the actual model,

852  comparing these distributions. (**Fig. S25**).

## Detecting environmentally associated loci under selection

854  In order to understand the degree to which environmentally driven natural

855  selection may cause chimpanzees to be locally adapted to different habitats, we used

856  latent factor mixed models implemented in the program LFMM v1.5 [117]. LFMM

857  quantifies statistical associations between allele frequencies and environmental

858  variables, accounts for underlying population genetic structure, and detects loci with

859  stronger environmental correlations than population structure. We ran five MCMC

860  replicates for all environmental variable with 25,000 burn-in steps, 100,000 iterations,

861  and a latent factor of $K$=3 from *a priori* knowledge of wild chimpanzee population

862  structure in Cameroon (**Figs. S11, S12, S13, S14 S15, S16, S20** and **S21**). We

863  calculated median z-scores across runs and used them to calculate the genomic

864  inflation factor ($\lambda$) and adjusted *p*-values. To correct for multiple testing, we applied a

865  conservative false discovery rate (FDR) of 0.1 using the Benjamini-Hochberg algorithm.

866  We identified unique candidate SNPs linked to at least one environmental variable,

867  presented via Manhattan plots using the *qqman* package [118].

868  Using outlier analysis with LFMM, we grouped highly correlated environmental

869  variables together to create 'environmental groupings' (**Table S9**). These included:

870  General Temperature (n=7), Temperature Range (n=2), Temperature Seasonality (n=2),

871  Precipitation – Wet/Cold (n=4), Precipitation – Dry/Warm (n=4), Tree Cover (n=2),

872  Vegetation Brownness (n=2), Vegetation Greenness (n=3), Surface Moisture Content

873  (n=2), and Topography (n=3).To determine the degree to which types of environmental

874     variation may drive selection of different genomic regions, we identified panels of unique

875     candidates from each 'environmental grouping.' We also analyzed the impact of

876     multicollinearity of our environmental predictors on the LFMM results by correlating the

877     degree of association between pairs of environmental predictors and the number of

878     shared outlier SNPs, using a Mantel test in R. (**Table S8**).

879     **Enriched gene ontologies and KEGG pathways**

880         We identified candidate genes near candidate SNPs positions with Ensembl's

881     BioMart tool [119, 120], including both complete and partial genes within these

882     windows. We used the DAVID database [44] for annotation and enrichment analysis of

883     candidate gene lists focusing on the 'Biological Processes' category of the 'Gene

884     Ontology' database [121] and KEGG pathways [60] using a $p$-value threshold of 0.05

885     and two different background populations of genes to control for potential bias since the

886     SNPs assayed in wild chimpanzees were selected from a subset of those identified

887     using whole-genome data from captive individuals. The first background we used was

888     composed of the population of genes found outside regions under selection identified in

889     the whole-genome sequencing data (**Fig 2a**). This resulted in a broad view of

890     environmentally mediated selection in wild chimpanzees by including only genes in

891     putatively neutral regions of the genome. The second background population of genes

892     consisted of all genes assayed in wild chimpanzees, excluding environmental outliers.

893

894

# Acknowledgments

# References

909

910   1.    Coyne JA, Orr HA. Speciation. Sunderland, Mass.: Sinauer Associates; 2004. xiii, 545, 2 p. of
911         plates p.
912   2.    Nielsen R. Molecular Signatures of Natural Selection. Annu Rev Genet. 2005;39(1):197-218. doi:
913         10.1146/annurev.genet.39.073003.112420.
914   3.    Moritz C, Patton JL, Schneider CJ, Smith TB. Diversification of rainforest faunas: An integrated
915         molecular approach. Annu Rev Ecol Syst.  2000;31(1):533-63. doi:
916         doi:10.1146/annurev.ecolsys.31.1.533.
917   4.    Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, et al. Positive natural
918         selection in the human lineage. Science. 2006;312(5780):1614-20. doi: 10.1126/science.1124309
919   5.    Hancock AM, Witonsky DB, Gordon AS, Eshel G, Pritchard JK, Coop G, Di Rienzo A.
920         Adaptations to climate in candidate genes for common metabolic disorders. PLoS Genet.
921         2008;4(2):e32. doi: 10.1371/journal.pgen.0040032.
922   6.    Manel S, Holderegger R. Ten years of landscape genetics. Trends Ecol Evol. 2013;28(10):614-
923         21. doi: 10.1016/j.tree.2013.05.012.
924   7.    Novembre J, Galvani AP, Slatkin M. The Geographic Spread of the CCR5 Δ32 HIV-Resistance
925         Allele. PLoS Biol. 2005;3(11):1954-62. doi: 10.1371/journal.pbio.0030339.
926   8.    Hedrick PW. Population genetics of malaria resistance in humans. Heredity. 2011;107(4):283-
927         304. doi: 10.1038/hdy.2011.16.
928   9.    Tishkoff SA, Reed FA, Ranciaro A, Voight BF, Babbitt CC, Silverman JS, et al. Convergent
929         adaptation of human lactase persistence in Africa and Europe. Nature Genet. 2007;39(1):31-40.
930         doi: 10.1038/ng1946.
931   10.   Fumagalli M, Moltke I, Grarup N, Racimo F, Bjerregaard P, Jørgensen ME, et al. Greenlandic
932         Inuit show genetic signatures of diet and climate adaptation. Science. 2015;349(6254):1343-7.
933         doi: doi:10.1126/science.aab2319.
934   11.   Jarvis JP, Scheinfeldt LB, Soi S, Lambert C, Omberg L, Ferwerda B, et al. Patterns of Ancestry,
935         Signatures of natural selection, and genetic association with stature in Western African Pygmies.
936         PLoS Genet. 2012;8(4):e1002641. doi: 10.1371/journal.pgen.1002641.
937   12.   Hancock AM, Witonsky DB, Ehler E, Alkorta-Aranburu G, Beall C, Gebremedhin A, et al. Human
938         adaptations to diet, subsistence, and ecoregion are due to subtle shifts in allele frequency. Proc
939         Natl Acad Sci USA. 2010;107(supplement_2):8924-30. doi: doi:10.1073/pnas.0914625107.
940   13.   Bhandari S, Zhang X, Cui C, Yangla, Liu L, Ouzhuluobu, et al. Sherpas share genetic variations
941         with Tibetans for high-altitude adaptation. Mol Genet Genomic Med. 2017;5(1):76-84. doi:
942         10.1002/mgg3.264.
943   14.   Hanaoka M, Droma Y, Basnyat B, Ito M, Kobayashi N, Katsuyama Y, et al. Genetic variants in
944         *EPAS1* contribute to adaptation to high-altitude hypoxia in Sherpas. PLoS ONE.
945         2012;7(12):e50566. doi: 10.1371/journal.pone.0050566.
946   15.   Scheinfeldt LB, Soi S, Thompson S, Ranciaro A, Woldemeskel D, Beggs W, et al. Genetic
947         adaptation to high altitude in the Ethiopian highlands. Genome Biol. 2012;13(1):R1. doi:
948         10.1186/gb-2012-13-1-r1.
949   16.   Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, et al. Great ape
950         genetic diversity and population history. Nature. 2013;499(7459):471-5. doi: 10.1038/nature12228

951   17.   Xue YL, Prado-Martinez J, Sudmant PH, Narasimhan V, Ayub Q, Szpak M, et al. Mountain gorilla
952         genomes reveal the impact of long-term population decline and inbreeding. Science.
953         2015;348(6231):242-5. doi: 10.1126/science.aaa3952. PubMed PMID: WOS:000352613700048.
954   18.   de Manuel M, Kuhlwilm M, Frandsen P, Sousa VC, Desai T, Prado-Martinez J, et al. Chimpanzee
955         genomic diversity reveals ancient admixture with bonobos. Science. 2016;354(6311):477-81. doi:
956         10.1126/science.aag2602.

957   19.   Cagan A, Theunert C, Laayouni H, Santpere G, Pybus M, Casals F, et al. Natural selection in the
958         great apes. Mol Biol Evol. 2016;33(12):3268-83. doi: 10.1093/molbev/msw215. PubMed PMID:
959         WOS:000387925300021.
960   20.   Nater A, Mattle-Greminger MP, Nurcahyo A, Nowak MG, de Manuel M, Desai T, et al.
961         Morphometric, behavioral, and genomic evidence for a new orangutan species. Current Biol.
962         2017;27(22):3487-98.e10. doi: 10.1016/j.cub.2017.09.047.
963   21.   Anthony NM, Johnson-Bawe M, Jeffery K, Clifford SL, Abernethy KA, Tutin CE, et al. The role of
964         Pleistocene refugia and rivers in shaping gorilla genetic diversity in central Africa. Proc Natl Acad
965         Sci USA. 2007;104(51):20432-6. Epub 2007/12/14. doi: 10.1073/pnas.0704816105.
966   22.   Gonder MK, Oates JF, Disotell TR, Forstner MR, Morales JC, Melnick DJ. A new west African
967         chimpanzee subspecies? Nature. 1997;388(6640):337. doi: 10.1038/41005.
968   23.   Gonder MK, Disotell TR, Oates JF. New genetic evidence on the evolution of chimpanzee
969         populations, and implications for taxonomy. Int J Primatol.  2006;27(4):1103-27. doi:
970         10.1007/s10764-006-9063-y.
971   24.   Mitchell MW, Locatelli S, Ghobrial L, Pokempner AA, Sesink Clee PR, Abwe EE, et al. The
972         population genetics of wild chimpanzees in Cameroon and Nigeria suggests a positive role for
973         selection in the evolution of chimpanzee subspecies. BMC Evol Biol. 2015;15:3. doi:
974         10.1186/s12862-014-0276-y.
975   25.   Mattle-Greminger MP, Bilgin Sonay T, Nater A, Pybus M, Desai T, de Valles G, et al. Genomes
976         reveal marked differences in the adaptive evolution between orangutan species. Genome Biol.
977         2018;19(1):193. doi: 10.1186/s13059-018-1562-6.
978   26.   Rodrigues MF, Kern AD, Ralph PL. Shared evolutionary processes shape landscapes of genomic
979         variation in the great apes. Genet. 2024;226(4):iyae006. doi: 10.1093/genetics/iyae006.
980   27.   Fontsere C, Kuhlwilm M, Morcillo-Suarez C, Alvarez-Estape M, Lester JD, Gratton P, et al.
981         Population dynamics and genetic connectivity in recent chimpanzee history. Cell Genom.
982         2022;2(6). doi: 10.1016/j.xgen.2022.100133.
983   28.   Sharp PM, Plenderleith LJ, Hahn BH. Ape Origins of Human Malaria. Annu Rev Microbiol.
984         2020;74(1):39-63. doi: 10.1146/annurev-micro-020518-115628.
985   29.   Leendertz SAJ, Wich SA, Ancrenaz M, Bergl RA, Gonder MK, Humle T, Leendertz FH. Ebola in
986         great apes – current knowledge, possibilities for vaccination, and implications for conservation
987         and human health. Mammal Rev. 2016:n/a-n/a. doi: 10.1111/mam.12082.
988   30.   Locatelli S, McKean KA, Sesink Clee PR, Gonder MK. The Evolution of Resistance to Simian
989         Immunodeficiency Virus (SIV): A Review. Int J Primatol. 2014;35(2):349-75. doi: 10.1007/s10764-
990         014-9763-7.
991   31.   Bibollet-Ruche F, Russell RM, Liu W, Stewart-Jones GBE, Sherrill-Mix S, Li Y, et al. CD4
992         receptor diversity in chimpanzees protects against SIV infection. Proc Natl Acad Sci USA.
993         2019;116(8):3229-38. doi: 10.1073/pnas.1821197116.
994   32.   Schmidt JM, de Manuel M, Marques-Bonet T, Castellano S, Andrés AM. The impact of genetic
995         adaptation on chimpanzee subspecies differentiation. PLoS Genet. 2019;15(11):e1008485. doi:
996         10.1371/journal.pgen.1008485.
997   33.   Pawar H, Ostridge HJ, Schmidt JM, Andrés AM. Genetic adaptations to SIV across chimpanzee
998         populations. PLoS Genet. 2022;18(8):e1010337. doi: 10.1371/journal.pgen.1010337.
999   34.   Locatelli S, Harrigan RJ, Sesink Clee PR, Mitchell MW, McKean KA, Smith TB, Gonder MK. Why
1000        are Nigeria-Cameroon Chimpanzees (*Pan troglodytes ellioti*) free of SIVcpz infection? PLoS
1001        ONE. 2016;11(8):e0160788. doi: 10.1371/journal.pone.0160788.
1002  35.   Prince AM, Brotman B, Lee D-H, Andrus L, Valinsky J, Marx P. Lack of evidence for HIV Type 1-
1003        related SIVcpz infection in captive and wild chimpanzees (*Pan troglodytes verus*) in West Africa.
1004        AIDS Res Hum Retroviruses. 2002;18(9):657-60. doi: 10.1089/088922202760019356. PubMed
1005        PMID: 12079561.
1006  36.   Leendertz SAJ, Locatelli S, Boesch C, Kücherer C, Formenty P, Liegeois F, et al. No evidence for
1007        transmission of SIVwrc from western red colobus monkeys (*Piliocolobus badius badius*) to wild
1008        west african chimpanzees (*Pan troglodytes verus*) despite high exposure through hunting. BMC
1009        Microbiol. 2011;11(1):24. doi: 10.1186/1471-2180-11-24.
1010  37.   Kalan AK, Kulik L, Arandjelovic M, Boesch C, Haas F, Dieguez P, et al. Environmental variability
1011        supports chimpanzee behavioural diversity. Nat Commun. 2020;11(1):4451. doi: 10.1038/s41467-
1012        020-18176-3.

38.  Wessling EG, Deschner T, Mundry R, Pruetz JD, Wittig RM, Kühl HS. Seasonal Variation in Physiology Challenges the Notion of Chimpanzees (*Pan troglodytes verus*) as a Forest-Adapted Species. Front Ecol Evol. 2018;6. doi: 10.3389/fevo.2018.00060.

39.  Boesch C, Kalan AK, Mundry R, Arandjelovic M, Pika S, Dieguez P, et al. Chimpanzee ethnography reveals unexpected cultural diversity. Nat Hum Behav. 2020;4(9):910-6. doi: 10.1038/s41562-020-0890-1.

40.  Kühl HS, Boesch C, Kulik L, Haas F, Arandjelovic M, Dieguez P, et al. Human impact erodes chimpanzee behavioral diversity. Science. 2019;363(6434):1453-5. doi: doi:10.1126/science.aau4532.

41.  Stumpf R. Chimpanzees and Bonobos, Diveristy within and between species. In: Campbell CJ, Fuentes A, MacKinnon KC, Panger M, Bearder SK, editors. Primates in Perspective. New York: Oxford University Press; 2007. p. 321-44.

42.  Sesink Clee PR, Abwe EE, Ambahe RD, Anthony NM, Fotso R, Locatelli S, et al. Chimpanzee genetic structure in Cameroon and Nigeria is associated with habitat variation that may be lost under climate change. BMC Evolutionary Biology. 2015;15:2. doi: 10.1186/s12862-014-0275-z.

43.  Pickrell JK, Coop G, Novembre J, Kudaravalli S, Li JZ, Absher D, et al. Signals of recent positive selection in a worldwide sample of human populations. Genome Res. 2009;19(5):826-37. doi: 10.1101/gr.087577.108.

44.  Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. Nat Protoc. 2009;4:44. doi: 10.1038/nprot.2008.211.

45.  Villadangos JA. Presentation of antigens by MHC class II molecules: getting the most out of them. Mol Immunol. 2001;38(5):329-46. doi: 10.1016/s0161-5890(01)00069-4.

46.  Pasquale EB. EPH receptor signalling casts a wide net on cell behaviour. Nat Rev Mol Cell Biol. 2005;6(6):462-75. doi: 10.1038/nrm1662. PubMed PMID: WOS:000229629100013.

47.  Kania A, Klein R. Mechanisms of ephrin-Eph signalling in development, physiology and disease. Nat Rev Mol Cell Biol. 2016;17(4):240-56. doi: 10.1038/nrm.2015.16.

48.  Aihara E, Engevik KA, Montrose MH. Trefoil Factor Peptides and Gastrointestinal Function. Annu Rev Physiol. 2017; 79: 357–380.

49.  Weinberg A, Jin G, Sieg S, McCormick T. The Yin and Yang of Human Beta-Defensins in Health and Disease. Front Immunol. 2012;3(294). doi: 10.3389/fimmu.2012.00294.

50.  Ellis N, Smith SJ, Pitcher CR. Gradient forests: calculating importance gradients on physical predictors. Ecol. 2012;93(1):156-68. doi: 10.1890/11-0252.1.

51.  Hansen TH, Bouvier M. MHC class I antigen presentation: learning from viral evasion strategies. Nat Rev Immunol. 2009;9(7):503-13. doi: 10.1038/nri2575.

52.  Luo Y, Jacobs EY, Greco TM, Mohammed KD, Tong T, Keegan S, et al. HIV–host interactome revealed directly from infected cells. Nat Microbiol. 2016;1(7):16068. doi: 10.1038/nmicrobiol.2016.68.

53.  Abitew AM, Sobti RC, Sharma VL, Wanchu A. Analysis of transporter associated with antigen presentation (TAP) genes polymorphisms with HIV-1 infection. Mol Cel Biochem. 2020;464(1):65-71. doi: 10.1007/s11010-019-03649-x.

54.  Fagerberg L, Hallström BM, Oksvold P, Kampf C, Djureinovic D, Odeberg J, et al. Analysis of the Human Tissue-specific Expression by Genome-wide Integration of Transcriptomics and Antibody-based Proteomics *. Mol Cell Proteomics. 2014;13(2):397-406. doi: 10.1074/mcp.M113.035600.

55.  Williams SM, An JY, Edson J, Watts M, Murigneux V, Whitehouse AJO, et al. An integrative analysis of non-coding regulatory DNA variations associated with autism spectrum disorder. Mol Psychiatry. 2019;24(11):1707-19. doi: 10.1038/s41380-018-0049-x.

56.  Eusebi PG, Cortés O, Carleos C, Dunner S, Cañon J. Detection of selection signatures for agonistic behaviour in cattle. J Anim Breed Genet. 2018;135(3):170-7. doi: 10.1111/jbg.12325.

57.  Parini P, Davis M, Lada AT, Erickson SK, Wright TL, Gustafsson U, et al. ACAT2 Is Localized to Hepatocytes and Is the Major Cholesterol-Esterifying Enzyme in Human Liver. Circulation. 2004;110(14):2017-23. doi: 10.1161/01.CIR.0000143163.76212.0B.

58.  Oue K, Zhang J, Harada-Hada K, Asano S, Yamawaki Y, Hayashiuchi M, et al. Phospholipase C-related Catalytically Inactive Protein Is a New Modulator of Thermogenesis Promoted by β-Adrenergic Receptors in Brown Adipocytes. J Biol Chem. 2016;291(8):4185-96. doi: 10.1074/jbc.M115.705723.

59. Jostins L, Ripke S, Weersma RK, Duerr RH, McGovern DP, Hui KY, et al. Host–microbe interactions have shaped the genetic architecture of inflammatory bowel disease. Nature. 2012;491(7422):119-24. doi: 10.1038/nature11582.

60. Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic Acids Res. 2000;28(1):27-30. doi: 10.1093/nar/28.1.27.

61. Kanehisa M. Toward understanding the origin and evolution of cellular organisms. Protein Sci. 2019;28(11):1947-51. doi: 10.1002/pro.3715.

62. Mitchell MW, Locatelli S, Sesink Clee PR, Thomassen HA, Gonder MK. Environmental variation *and* rivers govern the structure of chimpanzee genetic diversity in a biodiversity hotspot. BMC Evol Biol. 2015;15:1. doi: 10.1186/s12862-014-0274-0.

63. Gao F, Bailes E, Robertson DL, Chen Y, Rodenburg CM, Michael SF, et al. Origin of HIV-1 in the chimpanzee *Pan troglodytes troglodytes*. Nature. 1999;397(6718):436-41. doi: 10.1038/17130.

64. Keele BF, Van Heuverswyn F, Li Y, Bailes E, Takehisa J, Santiago ML, et al. Chimpanzee reservoirs of pandemic and nonpandemic HIV-1. Science. 2006;313(5786):523-6. doi: 10.1126/science.1126531

65. Van Heuverswyn F, Li Y, Neel C, Bailes E, Keele BF, Liu W, et al. SIV infection in wild gorillas. Nature. 2006;444(7116):164-. doi: 10.1038/444164a.

66. D'arc M, Ayouba A, Esteban A, Learn GH, Boué V, Liegeois F, et al. Origin of the HIV-1 group O epidemic in western lowland gorillas. Proc Natl Acad Sci USA. 2015;112(11):E1343-E52. doi: doi:10.1073/pnas.1502022112.

67. Bailes E, Gao F, Bibollet-Ruche F, Courgnaud V, Peeters M, Marx PA, et al. Hybrid origin of SIV in chimpanzees. Science. 2003;300(5626):1713-. doi: doi:10.1126/science.1080657.

68. Sharp PM, Shaw GM, Hahn BH. Simian Immunodeficiency virus infection of chimpanzees. J Virol. 2005;79(7):3891-902. doi: 10.1128/JVI.79.7.3891-3902.2005.

69. Rudicell RS, Holland Jones J, Wroblewski EE, Learn GH, Li Y, Robertson JD, et al. Impact of Simian Immunodeficiency Virus infection on chimpanzee population dynamics. PLoS Pathog. 2010;6(9):e1001116. doi: 10.1371/journal.ppat.1001116

70. de Groot NG, Otting N, Doxiadis GGM, Balla-Jhagjhoorsingh SS, Heeney JL, van Rood JJ, et al. Evidence for an ancient selective sweep in the MHC class I gene repertoire of chimpanzees. Proc Natl Acad Sci USA. 2002;99(18):11748-53. doi: 10.1073/pnas.182420799.

71. Maisels F. Mbam Djerem National Park, Cameroon: at the forest's edge. Canopee. 2005;27:2-6.

72. Smith TB, Wayne RK, Girman DJ, Bruford MW. A role for ecotones in generating rainforest biodiversity. Science. 1997;276(5320):1855-7. doi: 10.1126/science.276.5320.1855.

73. Simard F, Ayala D, Kamdem G, Pombi M, Etouna J, Ose K, et al. Ecological niche partitioning between *Anopheles gambiae* molecular forms in Cameroon: the ecological side of speciation. BMC Ecol. 2009;9(1):17. doi:10.1186/1472-6785-9-17.

74. Smith TB, Thomassen HA, Freedman AH, Sehgal RNM, Buermann W, Saatchi S, et al. Patterns of divergence in the olive sunbird *Cyanomitra olivacea* (Aves: Nectariniidae) across the African rainforest-savanna ecotone. Biol J Linnean Soc. 2011;103(4):821-35. doi: 10.1111/j.1095-8312.2011.01674.x.

75. Freedman AH, Thomassen HA, Buermann W, Smith TB. Genomic signals of diversification along ecological gradients in a tropical lizard. Mol Ecol. 2010;19(17):3773-88. doi: 10.1111/j.1365-294X.2010.04684.x.

76. Zhen Y, Harrigan RJ, Ruegg KC, Anderson EC, Ng TC, Lao SRN, et al. Genomic divergence across ecological gradients in the Central African rainforest songbird (*Andropadus virens*). Mol Ecol. 2017;26(19):4966-77. doi: 10.1111/mec.14270.

77. Morgan K, Mboumba J-F, Ntie S, Mickala P, Miller CA, Zhen Y, et al. Precipitation and vegetation shape patterns of genomic and craniometric variation in the central African rodent *Praomys misonnei*. Proc Biol Sci B. 2020;287(1930):20200449. doi: doi:10.1098/rspb.2020.0449.

78. Freedman AH, Harrigan RJ, Zhen Y, Hamilton AM, Smith TB. Evidence for ecotone speciation across an African rainforest-savanna gradient. Mol Ecol. 2023;32(9):2287–300. doi: 10.1111/mec.16867.

79. Abwe EE, Morgan BJ, Tchiengue B, Kentatchime F, Doudja R, Ketchen ME, et al. Habitat differentiation among three Nigeria–Cameroon chimpanzee (*Pan troglodytes ellioti*) populations. Ecol Evol. 2019;9(3):1489-500. doi: 10.1002/ece3.4871.

80. Mitchell MW, Locatelli S, Abwe EE, Ghobrial L, Gonder MK. Male-driven differences in chimpanzee (*Pan troglodytes*) population genetic structure across three habitats in Cameroon and Nigeria. Int J Primatol. 2018;39(4):581-601. doi: 10.1007/s10764-018-0053-7.

81. Abwe EE, Morgan BJ, Doudja R, Kentatchime F, Mba F, Dadjo A, et al. Dietary ecology of the Nigeria–Cameroon chimpanzee (*Pan troglodytes ellioti*). Int J Primatol. 2020;41(1):81-104. doi: 10.1007/s10764-020-00138-7.

82. Abwe EE. Linking behavioral diversity with genetic and ecological variation in the Nigeria-Cameroon chimpanzee (*Pan troglodytes ellioti*): Drexel University; 2018.

83. Ghobrial L, Lankester F, Kiyang JA, Akih AE, de Vries S, Fotso R, et al. Tracing the origins of rescued chimpanzees reveals widespread chimpanzee hunting in Cameroon. BMC Ecol. 2010;10(1):2. doi: 10.1186/1472-6785-10-2.

84. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv:1303.3997v1 [q-bio.GN]. 2013. doi: 10.48550/arXiv.1303.3997

85. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. ArXiv e-prints [Internet]. 2012 July 01, 2012. Available from: https://ui.adsabs.harvard.edu/#abs/2012arXiv1207.3907G.

86. Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, et al. Genome-wide detection and characterization of positive selection in human populations. Nature. 2007;449:913. doi: 10.1038/nature06250

87. Voight BF, Kudaravalli S, Wen X, Pritchard JK. A Map of Recent Positive Selection in the Human Genome. PLoS Biol. 2006;4(3):e72. doi: 10.1371/journal.pbio.0040072.

88. Maclean CA, Hong NPC, Prendergast JGD. hapbin: An Efficient Program for Performing Haplotype-Based Scans for Positive Selection in Large Genomic Datasets. Mol Biol Evol. 2015;32(11):3027-9. doi: 10.1093/molbev/msv172.

89. Delaneau O, Marchini J, Zagury J-F. A linear complexity phasing method for thousands of genomes. Nat Methods. 2012;9:179. doi: 10.1038/nmeth.1785

90. Auton A, Fledel-Alon A, Pfeifer S, Venn O, Ségurel L, Street T, et al. A fine-scale chimpanzee genetic map from population sequencing. Science. 2012;336(6078):193-8. doi: 10.1126/science.1216872.

91. Paten B, Herrero J, Beal K, Fitzgerald S, Birney E. Enredo and Pecan: Genome-wide mammalian consistency-based multiple alignment with paralogs. Genome Res. 2008;18(11):1814-28. doi: 10.1101/gr.076554.108.

92. Paten B, Herrero J, Fitzgerald S, Beal K, Flicek P, Holmes I, Birney E. Genome-wide nucleotide-level mammalian ancestor reconstruction. Genome Res. 2008;18(11):1829-43. doi: 10.1101/gr.076521.108.

93. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The Ensembl Variant Effect Predictor. Genome Biol. 2016;17:14. doi: 10.1186/s13059-016-0974-4.

94. Huff CD, Harpending HC, Rogers AR. Detecting positive selection from genome scans of linkage disequilibrium. BMC Genomics. 2010;11:9. doi: 10.1186/1471-2164-11-8. PubMed PMID: WOS:000274642200001.

95. Smedley D, Haider S, Durinck S, Pandini L, Provero P, Allen J, et al. The BioMart community portal: an innovative alternative to large, centralized data repositories. Nucleic Acids Res. 2015;43(W1):W589-W98. doi: 10.1093/nar/gkv350.

96. Perry GH, Marioni JC, Melsted P, Gilad Y. Genomic-scale capture and sequencing of endogenous DNA from feces. Mol Ecol. 2010;19(24):5332-44. doi: 10.1111/j.1365-294X.2010.04888.x.

97. Hernandez-Rodriguez J, Arandjelovic M, Lester J, Filippo C, Weihmann A, Meyer M, et al. The impact of endogenous content, replicates and pooling on genome capture from faecal samples. Mol Ecol Res. 2018;18(2):319-33. doi: doi:10.1111/1755-0998.12728.

98. Fan S, Hansen MEB, Lo Y, Tishkoff SA. Going global by adapting local: A review of recent human adaptation. Science. 2016;354(6308):54-9. doi: 10.1126/science.aaf5098.

99. Hancock AM, Witonsky DB, Ehler E, Alkorta-Aranburu G, Beall C, Gebremedhin A, et al. Human adaptations to diet, subsistence, and ecoregion are due to subtle shifts in allele frequency. Proc Natl Acad Sci USA. 2010;107(Supplement 2):8924-30. doi: 10.1073/pnas.0914625107.

100. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. Bioinformatics. 2011;27(15):2156-8. doi: 10.1093/bioinformatics/btr330.

101. Pew J, Muir PH, Wang J, Frasier TR. related: an R package for analysing pairwise relatedness from codominant molecular markers. Mol Ecol Res. 2015;15(3):557-61. doi: 10.1111/1755-0998.12323.

102. Wang J. Triadic IBD coefficients and applications to estimating pairwise relatedness. Genetical Res. 2007;89(3):135-53. doi: 10.1017/S0016672307008798.

103. Excoffier L, Lischer HEL. Arlequin suite ver 3.5: A new series of programs to perform population genetics analyses under Linux and Windows. Mol Ecol Resources. 2010;10(3):564-7. doi: 10.1111/j.1755-0998.2010.02847.x.

104. Hijmans RJ. geosphere: Spherical Trigonometry. R package version 15-5. 2016. Available from: https://CRAN.R-project.org/package=geosphere.

105. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. Genome Res. 2009;19(9):1655-64. doi: 10.1101/gr.094052.109.

106. Rosenberg NA. DISTRUCT: A program for the graphical display of population structure. Mol Ecol Notes. 2004;4(1):2. doi: 10.1046/j.1471-8286.2003.00566.x.

107. Jakobsson M, Rosenberg NA. CLUMPP: A cluster matching and permutation program for dealing with label and switching and multimodality in analysis of population structure. Bioinformatics. 2007;23(14):6. doi: 10.1093/bioinformatics/btm233.

108. Chen C, Durand E, Forbes F, François O. Bayesian clustering algorithms ascertaining spatial population structure: A new computer program and a comparison study. Mol Ecol Notes. 2007;7(5):747-56. doi: 10.1111/j.1471-8286.2007.01769.x.

109. Mitchell MW, Rowe B, Sesink Clee PR, Gonder MK. TESS Ad-Mixer: A novel program for visualizing TESS Q matrices. Conserv Genet Res. 2013;5(4):1075-8. doi: 10.1007/s12686-013-9987-4.

110. Jombart T, Ahmed I. adegenet 1.3-1: new tools for the analysis of genome-wide SNP data. Bioinformatics. 2011;27(21):3070-1. doi: 10.1093/bioinformatics/btr521.

111. Anderson EC, Thompson EA. A model-based method for identifying species hybrids using multilocus genetic data. Genetics. 2002;160(3):1217-29.

112. Wringe BF, Stanley RRE, Jeffery NW, Anderson EC, Bradbury IR. hybriddetective: A workflow and package to facilitate the detection of hybridization using genomic data in R. Mol Ecol Res. 2017;17(6):e275-e84. doi: doi:10.1111/1755-0998.12704.

113. Wringe BF, Stanley RRE, Jeffery NW, Anderson EC, Bradbury IR. parallelnewhybrid: an R package for the parallelization of hybrid detection using newhybrids. Mol Ecol Res. 2017;17(1):91-5. doi: doi:10.1111/1755-0998.12597.

114. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. Inferring the Joint Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. PLoS Genet. 2009;5(10):e1000695. doi: 10.1371/journal.pgen.1000695.

115. Fitzpatrick MC, Keller SR. Ecological genomics meets community-level modelling of biodiversity: mapping the genomic landscape of current and future environmental adaptation. Ecol Let. 2015;18(1):1-16. doi: https://doi.org/10.1111/ele.12376.

116. Láruson ÁJ, Fitzpatrick MC, Keller SR, Haller BC, Lotterhos KE. Seeing the forest for the trees: Assessing genetic offset predictions from gradient forest. Evol Appl. 2022;15(3):403-16. doi: 10.1111/eva.13354.

117. Frichot E, Schoville SD, Bouchard G, François O. Testing for associations between loci and environmental gradients using latent factor mixed models. Mol Bio Evol. 2013;30(7):1687-99. doi: 10.1093/molbev/mst063.

118. Turner SD. qqman: an R package for visualizing GWAS results using Q-Q and manhattan plots. J Open Source Software. 2018;3(25):731. doi: https://doi.org/10.21105/joss.00731.

119. Kinsella RJ, Kähäri A, Haider S, Zamora J, Proctor G, Spudich G, et al. Ensembl BioMarts: a hub for data retrieval across taxonomic space. Database. 2011;2011:bar030. doi: 10.1093/database/bar030.

120. Zerbino DR, Achuthan P, Akanni W, Amode M R, Barrell D, Bhai J, et al. Ensembl 2018. Nucleic Acids Res. 2018;46(D1):D754-D61. doi: 10.1093/nar/gkx1098.

121. Consortium GO. The Gene Ontology (GO) database and informatics resource. Nucleic Acids Research. 2004;32(suppl_1):D258-D61. doi: 10.1093/nar/gkh036.

# Supporting Information

**S1 File. Supplemental Results.**

**S2. File Extended Methods and Materials.**

**S1 Fig. Origins of chimpanzees of Cameroon included in this study.** Sample locations and proportions of estimated ancestry were estimated in previous studies [83, 122].

**Table S1. Captive chimpanzee genomes included in this study**

**S2 Fig. Heterozygosity estimates of captive chimpanzee genomes.**

(A) Individual heterozygosity.

(B) Subspecies heterozygosity.

**S3 Fig. Population structure of captive chimpanzee genomes**.

(A) PCA of LD pruned SNP data set consisting of 1,113,142 SNPs.

(B) sNMF individual ancestry analysis of the LD pruned data set in a range of K values.

(C) PCA of Neutral SNP data set consisting of 147,000 SNPs.

(D) sNMF individual ancestry analysis of the neutral SNP data set in a range of K values.

**S4 Fig. Cross Entropy Results.**
Value of the cross-entropy criterion as a function of the number of ancestral populations in sNMF for (A) the LD pruned SNP panel and (B) the neutral SNP panel.

**Table S2. Top 10 regions under selection including their genetic content.**

**Table S3. Enriched GO terms in the "Biological Processes" domain.**

**Table S4. Enriched KEGG pathways.**

**Table S5. Functional enrichment clustering.**

**Table S6. Number of wild chimpanzee samples collected and used in this study.**

**Table S7. Environmental predictor variables used characterize chimpanzee habitats.**

**Table S8. Pearson correlation table of environmental variables.**

**Table S9. Environmental variable groupings.**

**S5 Fig. Illumina reads of wild chimpanzee samples evaluated into 4 categories.**
Raw reads – yellow, mapped reads (sequence reads that mapped to the panTro4

1266  reference genome) – green, mapped and deduplicated reads (PCR duplicates were
1267  removed) – light blue, and on-target reads mapped to our sites – dark blue.

1268  (A) Proportion of read types for all 192 sequenced chimpanzee samples.

1269  (B) Number of read types for 85 samples included in the '10k' dataset.

1270  (C) Number of read types for samples removed due to missingness (>30% missing
1271  sites).

1272  (D) Proportion of read types for 85 samples included in the 'complete' dataset.

1273  (E) Number of read types for duplicate samples removed following relatedness analysis.

1274  (F) Number of read types for 23 individuals included in the '1k' dataset, but not the '10k'
1275  dataset.

1276  **S6 Fig. Frequency of mean sequencing coverage of wild chimpanzee samples for**
1277  **each site.**

1278  Mean coverage across all sites was 20.2 reads/site. The red vertical line represents the
1279  minimum coverage needed to accurately call SNPs (5x coverage).

1280  **S7 Fig. Pairwise $F_{ST}$ between sites shows population structure. (A) '10k' dataset.**
1281  **(B) '1k' dataset.**

1282  **S8 Fig. Isolation-by-distance for '10k' dataset.**

1283  (A) Correlation between 'linearized $F_{ST}$' and geographic distance (km) generated using
1284  the '10k' dataset. Solid circles represent pairs of sampling locations from the same
1285  habitat. Dual-colored diamonds represent pairs of sampling locations from different
1286  habitats.

1287  (B) Null distribution of t-statistics from 10,000 permutations same- or different
1288  habitat/population pairs in four bins of geographic distance. The red dotted line shows
1289  the t-statistic value for actual data.

1290  **S9 Fig. Isolation-by-distance for '1k' dataset.**

1291  (A) Correlation between 'linearized $F_{ST}$' and geographic distance (km) generated using
1292  the '1k' dataset. Solid circles represent pairs of sampling locations from the same
1293  habitat. Dual-colored diamonds represent pairs of sampling locations from different
1294  habitats.

1295  (B) Null distribution of $t$-statistics from 10,000 permutations same- or different
1296  habitat/population pairs in four bins of geographic distance. The red dotted line shows
1297  the $t$-statistic value for actual data.

1298  **S10 Fig. Isolation-by-distance for only *P. t. ellioti* populations.**

1299 (A) Correlation between 'linearized $F_{ST}$' and geographic distance (km) generated using
1300 the '10k' dataset. Solid circles represent pairs of sampling locations from the same
1301 habitat. Dual-colored diamonds represent pairs of sampling locations from different
1302 habitats.

1303 (B) Null distribution of $t$-statistics from 10,000 permutations same- or different
1304 habitat/population pairs in four bins of geographic distance. The red dotted line shows
1305 the $t$-statistic value for actual data.

1306 **S11 Fig. PCA of all SNPs from the '10k' dataset.**

1307 The first two principal components recapitulate known population structure of
1308 chimpanzees in Cameroon. They show 3 clear populations, and one *P. t. ellioti*
1309 (Ecotone) individual (CMMD06) clustering with *P. t. troglodytes*, as well as multiple *P. t.*
1310 *ellioti* (Rainforest) clustering together with *P. t. ellioti* (Ecotone) and vice versa.

1311 **S12 Fig. PCA of only neutral SNPs from the '10k' dataset.**

1312 **S13 Fig. PCA of all SNPs from the '1k' dataset**

1313 **S14 Fig. PCA of only neutral SNPs from the '1k' dataset**

1314 **Table S10. The results of the Tracy-Widom test for all SNPs from the '10k' SNP**
1315 **dataset**.

1316 **S15 Fig. ADMIXTURE bar plots for $K$=2-3.**

1317 These three populations correspond to known population structure. However, at $K$=2,
1318 there is a signal of possible historic gene flow of *P. t. troglodytes* into *P. t. ellioti*
1319 (Ecotone). Moreover, the is one individual (CMMD06 – also identified in the PCA) as
1320 being a potential *ellioti/troglodytes* hybrid. At $K$=3, we see evidence of three additional
1321 individuals that may be Rainforest/Ecotone hybrids, as well as evidence of mixing
1322 between the populations.

1323 **S16 Fig. The cross-validation error results of ADMIXTURE analysis of wild**
1324 **chimpanzees ('10k' dataset).**

1325 **S17 Fig. PCA results of the merged captive and wild datasets.**

1326 **S18 Fig. ADMIXTURE bar plots for $K$=2-5 for merged captive and wild datasets.**

1327 **S19 Fig. The cross-validation error results of ADMIXTURE analysis of merged**
1328 **captive and wild chimpanzees ('10k' dataset).**

1329 **S20 Fig. Cluster analysis and spatial interpolation of population structure.**

1330 (A) TESS bar plots showing individual proportions of ancestry of wild chimpanzees.

1331 (B) Spatial interpolation of the $Q$ matrix for $K$=3 generated using TESS and Ad-Mixer.

1332 **S21 Fig. Estimating $K_{MAX}$ from TESS analysis.**

54

1333    D$K$ values estimated for $K$=1-5 across 10 replicate runs.

1334    **Table S11. Analysis of Molecular Variance (AMOVA).**

1335    **S22 Fig. Mean observed heterozygosity.**

1336    (A) Heterozygosity of all loci for all individuals grouped by population.
1337    (B) Heterozygosity for all individuals.
1338    There were no significant differences between heterozygosity for each population.

1339    **S23 Fig. Posterior plots of model performance.**

1340    (A) The observed Joint SFS for *P. t. troglodytes* and *P. t. ellioti.*

1341    (B) the simulated Joint SFS for *P. t. troglodytes* and *P. t. ellioti* under the most likely
1342    asymmetric migration scenario obtained from *δαδι.*

1343    (C) The residuals between the modeled and observed Joint SFS.

1344    (D) A 1D histogram of the residual values between the model and the observed data

1345    **S24 Fig. $R^2$ weighted importance of the environmental predictor variables to the**
1346    **Gradient Forest model of gene-environment relationships.**

1347    **S25 Fig. Results of randomized gradient forest models (n=200), as compared to**
1348    **results from the observed data (n=200).**

1349    An average 588 of 7,878 SNPs demonstrated a positive $R^2$ with at least one
1350    environmental variable, with an average $R^2$ = 0.155 in the Observed data distribution
1351    (n=200, represented by the red histograms in A) and B) above). A significantly different
1352    average was obtained when randomizing the associations between the genomic data
1353    (SNPs) and environmental predictors for both total SNPs with a positive $R^2$ (average
1354    total = 504, t = 5.011(unequal variances df = 202.28), p < 0.0001), as well as for the
1355    average $R^2$ (average = 0.152, t = 2.806 (unequal variance df = 261.97), p = 0.0054) of
1356    the randomized gradient forests runs (n=200).

1357    **Table S12. Enriched GO terms in the 'Biological Processes' domain for**
1358    **environmentally associated outliers (LFMM and gradient forest) in wild**
1359    **chimpanzees**.

1360    **Table S13. Enriched KEGG pathways for environmentally associated outliers**
1361    **(LFMM and gradient forest) in wild chimpanzees.**

1362    **Table S14. Enriched GO terms in the 'Biological Processes' domain for**
1363    **environmentally associated outliers (LFMM and gradient forest) in wild**
1364    **chimpanzees.**

1365    **S26 Fig. Evidence of selective pressures on acetyl-CoA acetyltransferase 2**
1366    **(*ACAT2*).**

1367    (A) Map of the *ACAT2* gene on chromosome 6 with brown star representing the SNP
1368    identified through outlier analysis between *P. t. troglodytes* and *P. t. ellioti*.

(B) Manhattan plot showing the significance (as the negative $\log_{10}$ p-value) of SNP associations with the environmental variable temperature seasonality. Grey colors distinguish different chromosomes. The red line represents the threshold for significant association (p = 0.05). The SNP contained in the ACAT2 gene is highlighted by the red arrow.

(C) Correlation between the allele frequency of the SNP contained in the ACAT2 gene and temperature seasonality values at each corresponding sampling location ($R^2$ = 0.5615, p = 0.0005).

(D) Allele frequencies of the SNP contained in the ACAT2 gene across Cameroon. Sampling sites are represented by circles that are shaded according to the frequency of the allele within the population. SNP frequencies are plotted against temperature seasonality across the region.

**S27 Fig. Evidence of selective pressures on phospholipase C like 2 (*PLCL2*).**

(A) Map of the *PLCL2* gene on chromosome 3 with brown star representing the SNP identified through outlier analysis between *P. t. troglodytes* and *P. t. ellioti*.

(B) Manhattan plot showing the significance (as the negative $\log_{10}$ p-value) of SNP associations with the environmental variable precipitation of wettest month. Grey colors distinguish different chromosomes. The red line represents the threshold for significant association (p = 0.05). The SNP contained in the PLCL2 gene is highlighted by the red arrow.

(C) Correlation between the allele frequency of the SNP contained in the PLCL2 gene and precipitation of wettest month values at each corresponding sampling location ($R^2$ = 0.3422, p = 0.0102).

(D). Allele frequencies of the SNP contained in the PLCL2 gene across Cameroon. Sampling sites are represented by circles that are shaded according to the frequency of the allele within the population. SNP frequencies are plotted against precipitation of wettest month across the region.