

Diverse Functions Associate With Non-Coding Trans-species Polymorphisms in Humans

Keila Velazquez-Arcelay

Vanderbilt University

Mary Lauren Benton

Baylor University

John A. Capra ([✉ tony@capralab.org](mailto:tony@capralab.org))

University of California, San Francisco

Research Article

Keywords: trans-species polymorphisms, balancing selection, long-term balancing selection, non-coding variants, phenome-wide association study

Posted Date: May 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-559297/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Diverse functions associate with non-coding trans-species polymorphisms in**
2 **humans**

4 Keila Velazquez-Arcelay¹, Mary Lauren Benton², and John A. Capra^{1,3,4*}

6 ¹ Department of Biological Sciences, Vanderbilt University

7 ² Department of Computer Science, Baylor University

8 ³ Departments of Biomedical Informatics and Computer Science, Genetics Institute, and Center for Structural
9 Biology, Vanderbilt University

10 ⁴ Bakar Computational Health Sciences Institute and Department of Epidemiology and Biostatistics, University of
11 California, San Francisco

12 * Corresponding author: tony@capralab.org

14 **Abstract**

15 *Background:* Long-term balancing selection (LTBS) can maintain allelic variation at a locus over
16 millions of years and through speciation events. Variants shared between species, hereafter “trans-
17 species polymorphisms” (TSPs), often result from LTBS due to host-pathogen interactions. For
18 instance, the major histocompatibility complex (MHC) locus contains TSPs present across
19 primates. Several hundred candidate TSPs have been identified in humans and chimpanzees;
20 however, because many are in non-coding regions of the genome, the functions and adaptive roles
21 for most TSPs remain unknown.

22 *Results:* We integrated diverse genomic annotations, with a focus on non-coding regions, to
23 explore the functions of 125 previously identified regions containing multiple TSPs in humans and
24 chimpanzees. We analyzed genome-wide functional assays, expression quantitative trait loci
25 (eQTL), genome-wide association studies (GWAS), and phenome-wide association studies
26 (PheWAS). We identify functional annotations for 119 TSP regions, including 71 with evidence
27 of gene regulatory function from GTEx or genome-wide functional genomics data and 21 with
28 evidence of trait association from GWAS and PheWAS. TSPs in humans associate with many
29 immune system phenotypes, including response to pathogens, but we also find associations with a
30 range of other phenotypes, including body mass, alcohol intake, urate levels, chronotype, and risk-
31 taking behavior.

32 *Conclusions:* The diversity of traits associated with non-coding human TSPs further support
33 previous hypotheses that functions beyond the immune system are subject to LTBS. Furthermore,
34 several of these trait associations provide support and candidate genetic loci for previous
35 hypothesis about behavioral diversity in great ape populations, such as the importance of variation
36 in sleep cycles and risk sensitivity.

37
38 **Keywords:** trans-species polymorphisms; balancing selection; long-term balancing selection;
39 non-coding variants; phenome-wide association study

42 **Significance statement**
43 Most genetic variants present in human populations are young (<100,000 years old); however, a
44 few hundred are millions of years old with origins before the divergence of humans and
45 chimpanzees. Many of these trans-species polymorphisms (TSPs) were likely maintained by
46 balancing selection—evolutionary pressure to maintain genetic diversity at a locus. In spite of
47 their age, the functions driving this selection, especially for non-coding TSPs, are largely
48 unknown. We integrate genome-wide annotation strategies to demonstrate TSP associations with
49 immune system function, behavior (addition, cognition, risky behavior), uric acid metabolism,
50 and many other phenotypes. These results substantially expand our understanding of functions
51 for TSPs and support a role for balancing selection beyond the immune system.

52
53

54 **Background**

55 The interaction between populations and environments is dynamic. Over time, allele frequencies
56 in a population shift due to drift and adaptive responses to specific environmental pressures.
57 Most genetic variants are short-lived compared to the timescale of species. But on rare occasions
58 variants persistently segregate at intermediate frequencies for millions of years, sometimes
59 predating the most recent common ancestor (MRCA) between two sister species (Bitarello et al.,
60 2018; Cheng & DeGiorgio, 2019; DeGiorgio et al., 2014; Leffler et al., 2013; Siewert & Voight,
61 2017; Teixeira et al., 2015). These trans-species polymorphisms (TSPs) are likely a sign of
62 genomic regions under long-term balancing selection (LTBS). Over time, instances of LTBS
63 leave signatures in the genome that differentiate them from those under other forms of selection
64 (Bitarello et al., 2018; Key et al., 2014; Leffler et al., 2013; Siewert & Voight, 2017).

65 Several instances of likely LTBS have been observed in humans and other primates,
66 mostly within the major histocompatibility complex (MHC) or the ABO blood group locus. For
67 example, the MHC, or human leukocyte antigen (HLA) system in humans, is a family of varied
68 proteins expressed on the cell surface with essential functions in adaptive immune response and
69 regulation. Balancing selection on different components of the HLA region dates to the common
70 ancestor between chimpanzees and humans (Lawlor et al., 1988; Mayer et al., 1988) (Azevedo et
71 al., 2015). Similarly, the ABO gene has three alleles, and its variants lead to different blood cell
72 antigens, or lack of thereof, on the surface of the cell. Variation in this group could have a
73 benefit in the immune response to pathogens, and balanced polymorphisms at this locus date
74 back to the common ancestor of gorillas, orangutans, and humans (Ségurel et al., 2012). Several
75 other immune-related genes show LTBS between humans and other primates, e.g.: *TRIM5*, a
76 RING finger protein 88 (Battivelli et al., 2011; Cagliani et al., 2010; Ganser-Pornillos &
77 Pornillos, 2019), and *ZC3HAV1*, a zinc finger CCCH-type antiviral protein 1 (Cagliani et al.,
78 2012; De Filippo et al., 2016; Mao et al., 2013; Todorova et al., 2015). These genes have
79 important roles in host/pathogen response through inhibition of virus replication.

80 The high allelic variation maintained by balancing selection at a locus can also enable
81 adaptive selection in new environments. For example, some variants found under balancing

82 selection in African and ancestral human populations have experienced directional selection in
83 non-African populations (European and Asian), with one allele becoming predominant in the
84 population (De Filippo et al., 2016). This suggests the adaptive potential of the variation
85 maintained under balancing selection; however, in some cases the adaptive variants themselves
86 may have hitchhiked with those under LTBS.

87 Recent studies have developed methods to identify instances of balancing selection in
88 genome-wide data (Bitarello et al., 2018; Cheng & DeGiorgio, 2019; DeGiorgio et al., 2014;
89 Siewert & Voight, 2017, 2020). Some have focused on detecting LTBS using trans-species data,
90 while others have considered balancing selection over shorter timescales based on single-species
91 data. For example, De Giorgio (2014) developed likelihood-ratio tests (T_1 and T_2) based on
92 computing probabilities of polymorphism and substitution under LTBS based on inter-species
93 coalescent modeling. With this method they identified balancing selection on HLA regions, but
94 also in a gene that had no previous associations with balancing selection, *FANK1*, which is
95 involved in the suppression of apoptosis during/after the process of meiosis. They also found
96 enrichment in other functions: cell adhesion, membrane protein activity, and components of
97 membranes. A more recent study (Cheng & DeGiorgio, 2019) expanded the T_2 method to seek
98 trans-species balancing selection without direct consideration of trans-species polymorphism and
99 identified a handful of additional LTBS candidates. Bitarello et al. (2018) developed Non-central
100 Deviation (NCD) statistics that quantify the deviation from expectations under balancing
101 selection of the local site frequency spectrum (SFS) in windows in African and European 1000
102 Genomes populations. They identified thousands of candidates for balancing selection in
103 humans. They also showed varying directional selection in different populations, providing
104 evidence for the adaptive potential of regions under balancing selection. Siewert & Voight
105 (2017) developed β , a summary statistic for detecting genomic windows with clusters of
106 intermediate frequency alleles suggestive of balancing selection. They also recently updated the
107 β statistic to consider both polymorphism and substitution data (Siewert & Voight, 2020).
108 Among the highest scoring windows in these two analyses, they highlighted three genes
109 (*CADM2*, *WFS1*, and *ACSBG2*) with functions outside the immune system.

110 Trans-species polymorphisms, especially when more than one falls on a haplotype,
111 suggest the action of LTBS. Leffler et al. (2013) compared polymorphisms across the genome in
112 Yoruba individuals from the 1000 Genomes Project to those found in Western chimpanzees
113 sequenced by the PanMap Project. They identified more than 100 non-coding haplotypes with
114 multiple TSPs within 4 kb and in high LD suggesting the presence of LTBS. Sequencing errors
115 and high mutation rates can produce TSPs (Cheng & DeGiorgio, 2019; Gao et al., 2015) but it is
116 unlikely to observe haplotypes with two TSPs in close proximity by chance.

117 Despite the importance and prevalence of balancing selection, most the non-coding
118 haplotypes bearing potential signatures of LTBS, i.e., multiple TSPs, have not been functionally
119 characterized. Here, we focus on the non-coding TSPs identified by Leffler et al. (2013).
120 Determining the functional roles of these polymorphisms in human adaptation and health will
121 deepen our understanding of the dynamics of balancing and positive selection and their roles in

122 adaptation to new environments. We identify potential functions of the TSPs in humans by
123 applying several genome-wide functional annotations and association tests. Our results identify
124 diverse functions, including effects unrelated to the immune system, that may underlie LTBS on
125 the human and chimpanzee lineages.

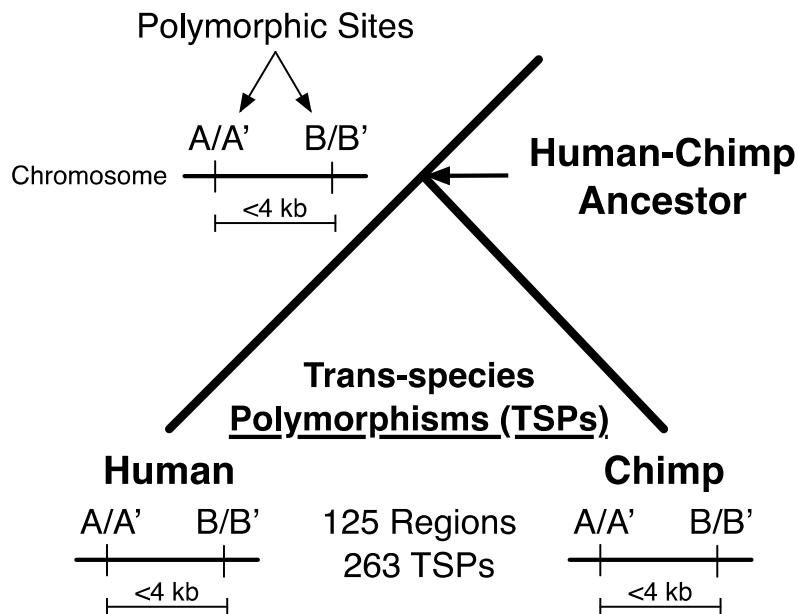
126

127 **Results**

128 *Human-chimpanzee TSPs*

129 We consider 125 human genomic regions containing multiple TSPs segregating in both humans
130 and chimpanzees in close proximity and high LD from Leffler et al. (2013). This set is composed
131 of 263 variants with strong evidence of identity-by-descent; i.e., the divergence between the
132 ancestral and derived alleles is deeper than the human-chimp speciation event. These TSPs were
133 identified based on the observation from coalescent theory (Ségurel et al., 2012) that pairs of
134 TSPs within 4 kb in the human genome are extremely unlikely to result from neutral processes,
135 and thus are strong candidates for LTBS (Figure 1). Hereafter, we refer to these as “TSP
136 regions”. While these criteria do not guarantee that all the TSP regions are the result of LTBS
137 (Gao et al., 2015), this set is likely strongly enriched for LTBS. We also quantify the evidence
138 for balancing selection on these loci provided by BetaScan2 and NCD, two additional recent
139 genome-wide scans (Bitarello et al., 2018; Siewert & Voight, 2020). In total, 80% (100/125) of
140 the TSP regions had additional evidence of balancing selection from at least one of BetaScan2 or
141 NCD (Methods; Supplementary Table 7).

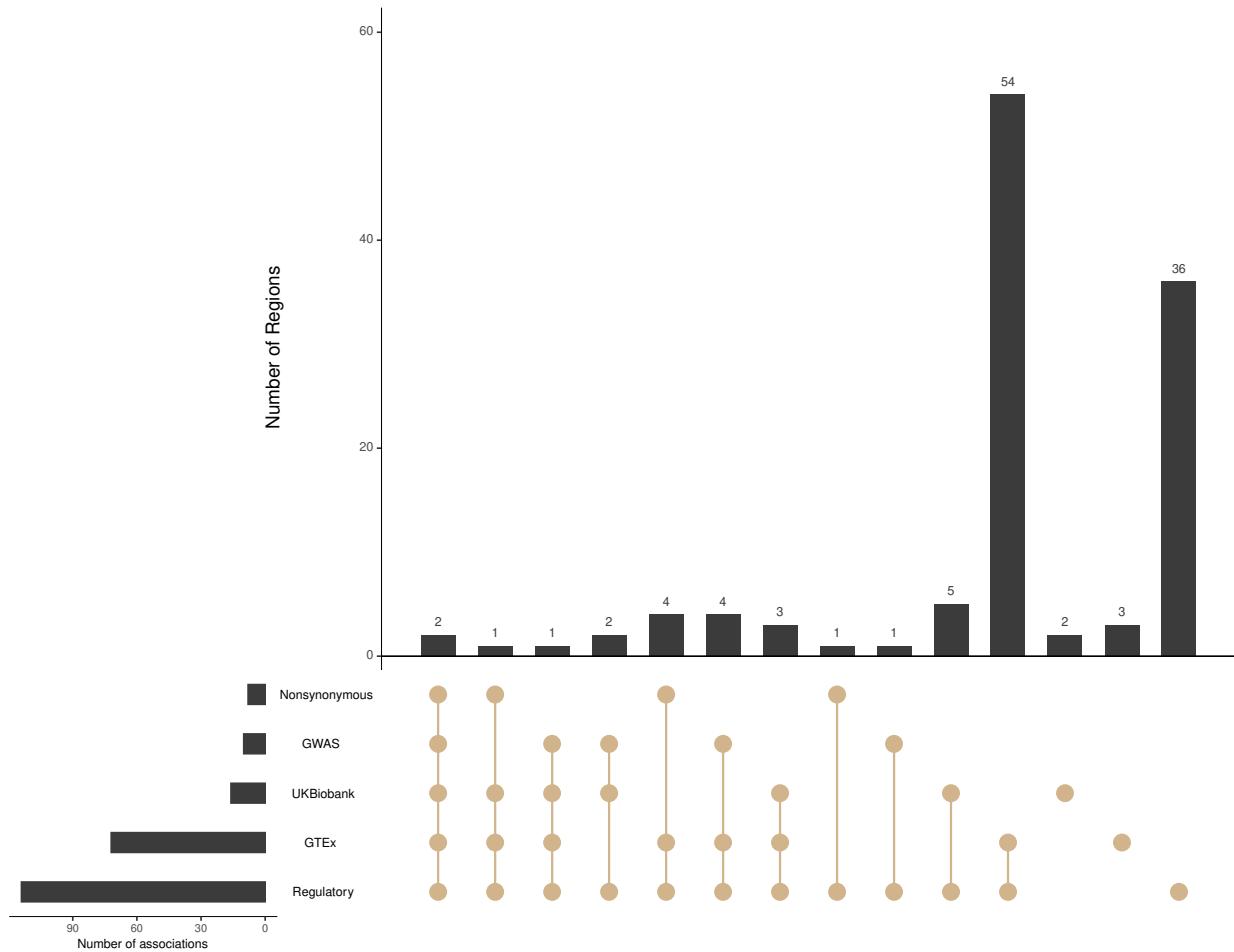
142 In the following, we analyze two sets of variants for the 125 TSP regions. First, we focus
143 on the 263 TSPs themselves. Second, to capture functions tagged by variants in high linkage
144 disequilibrium (LD) with TSPs, we also considered potential tag SNPs in high LD ($R^2 \geq 0.8$) with
145 TSPs in African, European, or East Asian populations from the 1000 Genomes Project. This LD-
146 expanded set includes 10,259 variants across the 125 TSP regions (Supplementary Figure 1). By
147 expanding to include LD, we capture additional associations, but may also introduce false
148 positives; thus, we report results on both sets throughout.



149
150 **Figure 1. Trans-species polymorphisms (TSPs) likely resulting from long-term balancing**
151 **selection (LTBS).** Schematic showing the criteria used by Leffler et al. (2013) to identify TSPs
152 likely maintained by LTBS. Each line represents a chromosome with polymorphisms segregating
153 in a population. A/A' are two alleles segregating in both humans and chimpanzees at one site
154 (i.e., a TSP), and B/B' are alleles segregating in both species at a nearby site. TSPs are very
155 unlikely to appear nearby (within 4 kb) without the action of balancing selection. We consider
156 125 TSP regions containing 263 TSPs. Within these regions, multiple functional scenarios are
157 possible. For example, one TSP may be under LTBS while the other is neutral, but maintained
158 due to tight linkage. Alternatively, the TSPs may have epistatic functions and both be under
159 selection. In addition to the 263 TSPs, we also considered functional associations with 9,996
160 variants in high LD ($r^2 > 0.8$) with a TSP at least one population from the 1000 Genomes Project
161 (Supplementary Figure 1).

162
163 *Trans-species polymorphisms overlap diverse functional annotations*
164 We intersected the TSPs with diverse lines of functional evidence from large-scale genomic
165 studies, including genome-wide functional genomics assays, eQTL, GWAS, and PheWAS. We
166 found at least one functional annotation for 95% of the TSP regions (119 out of 125) covering
167 130 TSPs and 4,807 LD SNPs (Figure 2). Here, we provide an overview of the overlap with
168 these annotations. In future sections, we provide details about each of these annotations. Variants
169 in 91% (114 out of 125) of regions overlap annotated gene regulatory regions. This includes 58
170 TSPs in 40 regions and 1334 LD variants in 112 regions. We also found 86 TSPs across 51
171 regions with evidence of being expression quantitative trait loci (eQTL) in 29 tissues. Including
172 the variants in LD with TSPs, 57.6% of the regions (72 out of 125) contain eQTLs ($p < E-5$ in

173 49 tissues. We found genome-wide significant associations with phenotypes in available
174 genome- or pheno- wide association studies for 19% of the regions (24 out of 125; 13 GWAS
175 and 16 geneAtlas PheWAS). Finally, 11.2% (14 out of 125) of these regions contain SNPs in
176 protein-coding regions; one TSP (rs12918619) produces a non-synonymous protein sequence
177 change, and 16 other non-synonymous variants in 7 regions are in high LD with the TSPs.
178 (These were not identified as coding in the previous study due to changes in protein annotations.)
179
180



181
182 **Figure 2. Functional annotations available for the expanded TSP regions.** Summary of the
183 annotations of each type available for TSP regions, including tagging variants in high LD with
184 TSPs. A total of 119 out of 125 TSP regions contain at least one line of functional evidence.
185 Multiple lines (two or more) are available for 78 regions.
186

187 *Evidence of gene regulatory function for TSPs*

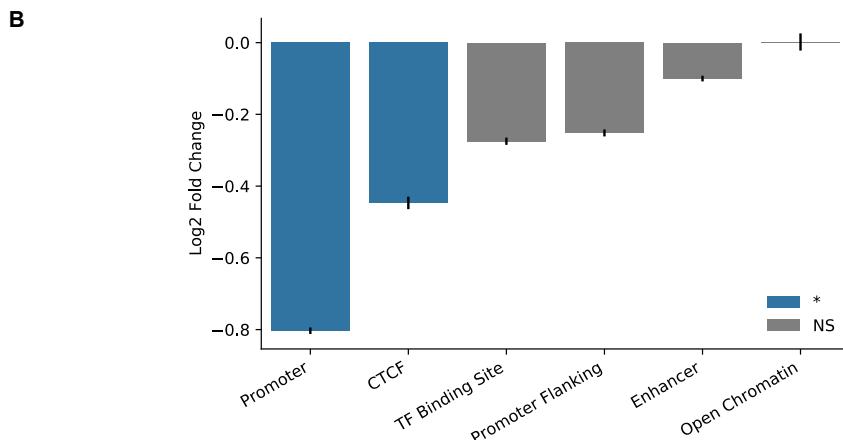
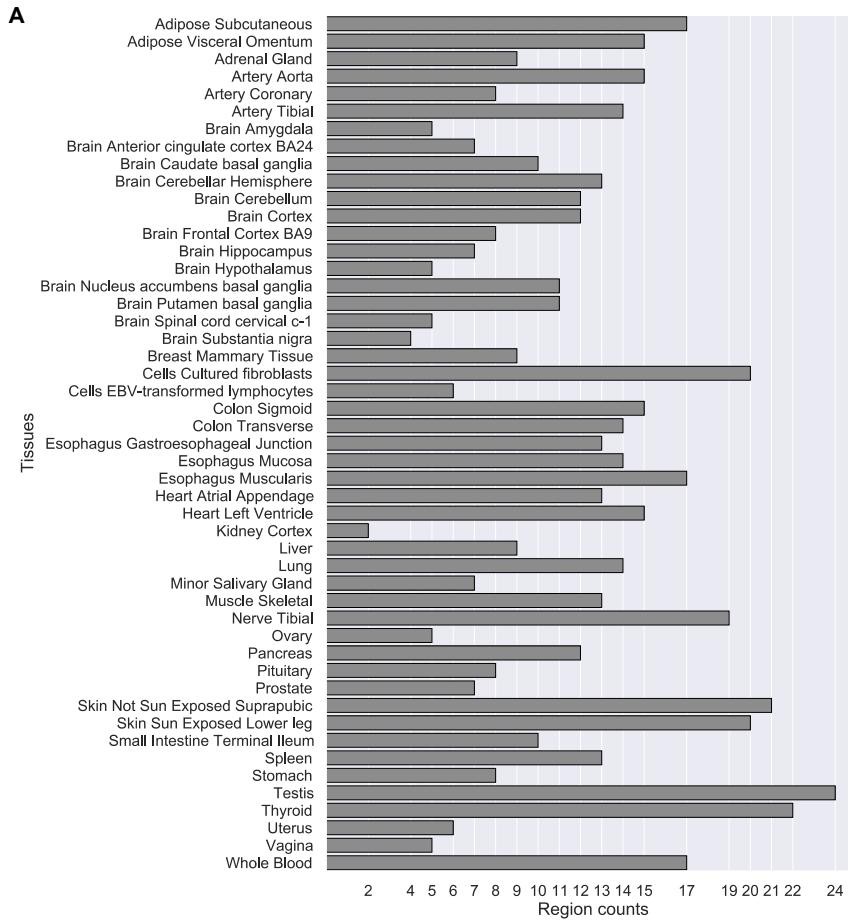
188 We hypothesized that many of the non-coding TSPs in our set perform gene regulatory
189 functions. To evaluate this possibility, we intersected the TSPs and variants in high LD with
190 maps of functional regulatory regions from the Ensembl regulatory build (Zerbino et al., 2015).
191 We found 58 TSPs with regulatory annotations in 40 TSP regions, and additionally 1334 LD

192 variants in 114 regions. These include variants in CTCF binding sites, open chromatin regions,
193 promoter flanking regions, enhancers, promoters, and known TF binding sites (Supplementary
194 Table 1).

195 Overlap of a variant with a regulatory annotation does not necessarily imply a regulatory
196 function. To consider additional evidence of regulatory function, we examined eQTL in GTEx
197 from 50 tissues for overlap with TSPs. At least one eQTL was found for 51 of the TSP regions
198 (40%). Among these 51 regions, 29 TSPs are themselves eQTL; for the remainder, variants in
199 high LD with TSPs were eQTL. The eQTL were found across diverse tissues (Figure 3A), and
200 there was no enrichment for specific gene ontology (GO) terms among the set of genes
201 influenced by TSP eQTL. This suggests that the targets of balancing selection have functions in
202 gene regulation across diverse tissues beyond the immune system.

203 Next, we tested if the TSP regions are enriched for overlap with any specific types of
204 regulatory regions. We compared the observed overlap between TSP regions and each type of
205 regulatory annotation to the distribution of overlaps expected if TSP regions were randomly
206 distributed across the genome. We shuffled the TSP regions 1,000 times maintaining their length
207 and chromosome distributions and avoiding genome assembly gaps and ENCODE blacklist
208 regions and counted the number of overlaps observed with regulatory elements for each random
209 permutation. The TSPs overlap slightly fewer base pairs annotated with regulatory functions than
210 expected by chance, with significant ($P < 0.05$) depletion for promoter and CTCF sites (Figure
211 3B). Since variants in these regions are likely to influence gene regulation in many tissues (e.g.,
212 compared to enhancers which are often context-specific), this suggests that individual TSPs may
213 be less pleiotropic than expected by chance.

214



215

216 **Figure 3. TSPs are eQTLs in diverse tissues and are depleted for overlap with promoters**
 217 **and CTCF sites.** (A) The number of TSP regions that contain an eQTL for each GTEx tissue.
 218 Variation in TSP regions associates with gene expression in diverse tissues. The associated genes
 219 also have Gene Ontology (GO) annotations from diverse functional categories (Supplementary
 220 Figure 2). (B) TSP regions are significantly depleted for overlap with promoters and CTCF sites
 221 compared to length- and chromosome-matched non-coding regions from the genomic
 222 background. The error bars represent 95% confidence intervals.

223
224 *Genome-wide association studies link TSPs to traits*
225 Genome-wide association studies have identified thousands of associations between genetic
226 variants and human traits. We intersected the TSP regions with associations reported in the
227 GWAS Catalog (downloaded 2020/11), which is composed of 227,262 associations. Since TSPs
228 themselves were not always directly tested in GWAS studies, we also include genome-wide
229 significant ($p < 1E-8$) associations with the tag variants in high LD with TSPs. We found
230 significant associations for 29 different variants (Figure 4; Supplementary Table 2). Two main
231 functional categories were identified in the GWAS associations for these variants:
232 immunological functions and neurological/behavioral traits. The associations with immune traits
233 were expected given the results of previous balancing selection studies and the few well-
234 characterized instances of LTBS. We identified many variants in LD with TSPs as associated
235 with blood measurement phenotypes and diseases related to immune response failure
236 (Supplementary Tables 2, 3, 4). Some of these traits include sarcoidosis (chromosome 4 near
237 TSPs rs114383553 and rs17007061), ulcerative colitis and chronic inflammatory disease (chr2
238 near TSPs rs13426764/rs11694806), and Behçet's disease (chromosome 11 near TSPs rs2249268
239 and rs60427879).

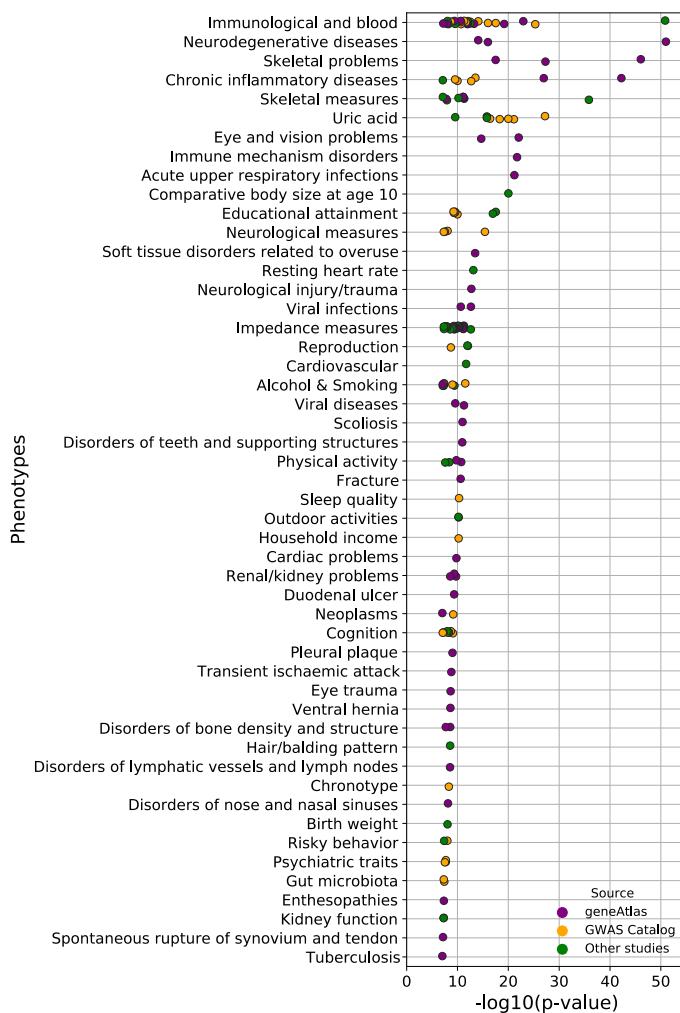
240 We also found many neurological/behavioral related associations among TSP region
241 variants. These traits include cognitive performance (chromosome 2 near rs13426764 and
242 rs11694806 and chromosome 3 near rs9869178/rs2118072), chronotype (chromosome 4 near
243 rs1887944 and rs7147645), addiction (alcohol use: chromosome 16 near rs9933768 and
244 rs57790054; smoking: chromosome 2 near rs13426764 and rs11694806), risky behavior
245 (automobile speeding propensity: chromosome 3:rs9869178/rs2118072), and experiencing mood
246 swings (chromosome 2 near rs13426764 and rs11694806). In addition to the immune response
247 and neurological categories, we observed associations in reproductive traits (polycystic ovary
248 syndrome, testosterone levels), urate levels, pancreatic cancer, hepatocyte growth factor levels,
249 and gut microbiota. We discuss several of these associations in more detail in following sections.
250

251 *Phenome-wide association studies link TSPs to additional diverse traits*
252 The growth of biobanks with linked genetic and phenotypic data has enabled the testing of the
253 association of genetic variants with diverse traits within a single cohort. This PheWAS approach
254 enables exploration of the functional and potentially pleiotropic effects of variants of interest
255 (Bush et al., 2016). Using published associations from the UK Biobank, we analyzed the
256 association of TSPs with 778 traits; all 125 of the TSP regions were tested. Overall, 21 TSPs in
257 16 regions had at least one genome-wide significant association ($P < 1E-8$, Figure 4;
258 Supplementary Table 3). Though testing different phenotypes than the GWAS, these associations
259 were qualitatively similar to the GWAS results, in that blood and immune system phenotypes
260 had many associations with TSPs, but the TSPs were also associated with a more diverse set of
261 phenotypes. We found associations in three major categories: immune response traits, body and
262 physical measurements, and neuropsychiatric traits related to addiction. We also observed

263 associations with many other phenotypes, for example: walking habits, heart disorders, renal and
264 kidney problems, pleural plaques, transient ischemic attack, cancerous tumor, and rupture of
265 synovium and tendon (Supplementary Tables 3 and 4).

266

267



268

269 **Figure 4. Genome- and phenotype-wide association studies link TSPs to diverse traits.**

270 Genome-wide significant ($P < 1E-8$) associations from the GWAS Catalog (yellow), a PheWAS
271 over the UK Biobank from the geneAtlas (purple), and other studies summarized in the GWAS
272 Atlas (green) (Watanabe et al., 2019). Each dot represents an association between a TSP region
273 and a trait. Many immune-related traits are associated with TSPs, but there are also associations
274 with a wider variety of phenotypes including osseous, neurological, and nervous system traits.
275 Five extreme associations with immunological and blood traits ($P < 1E-60$) were truncated for
276 this visualization. Since few TSPs themselves were directly tested in GWAS, we include GWAS
277 Catalog associations with tag variants in high LD ($r^2 > 0.8$) with TSPs. All associations are listed
278 in Supplementary Tables 2–4.

279

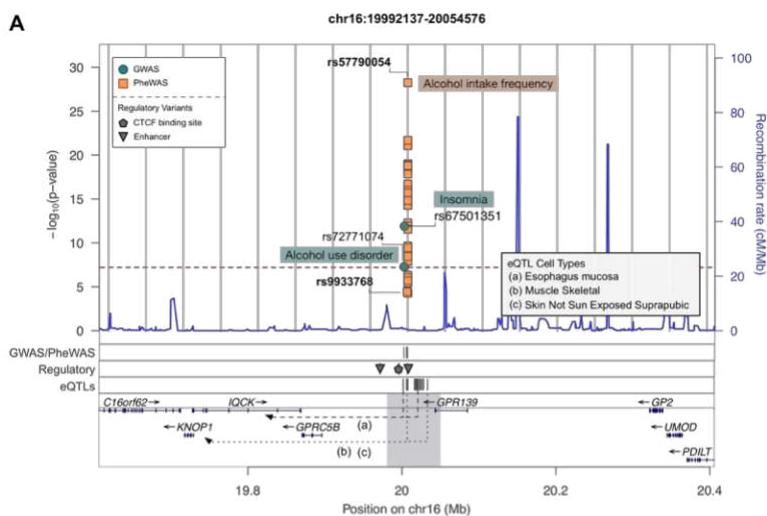
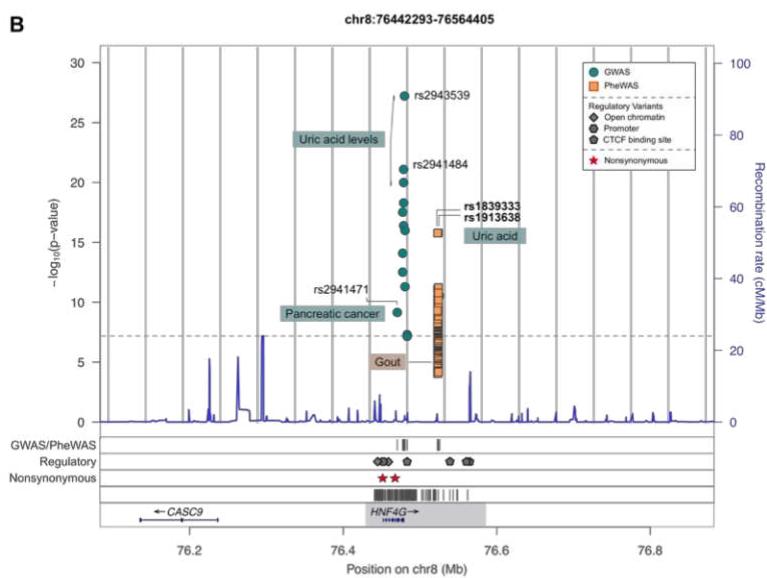
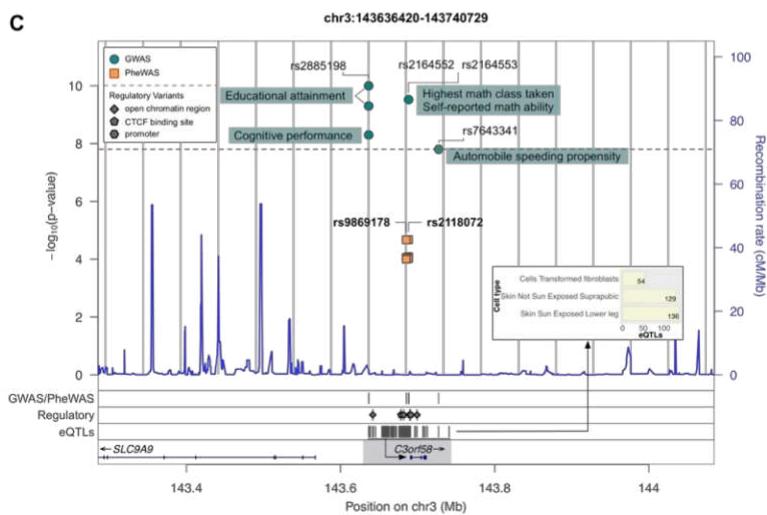
280 *Evidence of protein-coding function for TSPs*
281 The TSPs were originally filtered by Leffler et al. (2013) to be non-coding; however, three of the
282 263 TSPs are coding variants: two variants containing synonymous alleles in the genes dynein
283 (*DNHD1*) and in a transcript of *DNHD1*'s paralog *DNAH9*, and one non-synonymous variant in
284 *PKD1L2*. This discrepancy is due to changes in genomic annotations over time. For example,
285 *PKD1L2* is a calcium channel with potential roles in kidney function. This gene is a polymorphic
286 pseudogene in humans, with the reference genome encoding the pseudogene form; this likely
287 explains why this coding variant was not identified as coding in previous studies. With the LD-
288 expanded set, a total of 9 regions had non-synonymous codon changes (19 variants,
289 Supplementary Table 5). Some of the genes containing these variants include: Sperm Associated
290 Antigen 16 (*SPAG16*), which codes for two proteins that associate with the axoneme of sperm
291 cells; Hepatocyte Nuclear Factor 4 Gamma (*HNF4G*) codes for a receptor involved in DNA
292 binding transcription activity; Leucine Rich Colipase Like 1 (*LRCOL1*) which has enzyme
293 activation activity and is involved in digestion. All of the coding variants had CADD scores
294 suggestive of deleterious effects (scaled C-score > 20). Further work is needed to determine if
295 these coding variants in high LD with the TSPs influence selection.

296
297 *Illustrative examples of diverse functions associated with TSP regions*
298 Integrating the above data, we found 79 TSP regions with two or more lines of functional
299 evidence. This included twelve regions with annotations from at least four sources. To illustrate
300 the diverse functions associated with TSPs, we highlight four of these regions. All the example
301 regions were also detected by at least one of two other methods of finding genomic regions under
302 balancing selection (Supplementary Table 7) (Bitarello et al., 2018; Siewert & Voight, 2020).

303
304 *Body mass and alcohol intake.* A TSP (rs57790054) on 16p12.3 (hg38.chr16:19992138-
305 20043254) is strongly associated with several growth and body mass phenotypes as well as
306 alcohol intake frequency (Figure 5A; P < 4E-10 for all). Another variant in high LD in
307 Europeans (rs72771074, R²=0.89) with a TSP (rs57790054) in this locus was associated with
308 alcohol use disorder in a previous GWAS in a European cohort (P = 5E-8) (Sanchez-Roige et al.,
309 2019). The nearest gene, *GPR139*, encodes for a G-protein coupled receptor expressed in the
310 brain that is involved in alcohol drinking and withdrawal symptoms in rats (Kononoff et al.,
311 2018). This region contains several variants in LD with TSPs in regulatory regions, such as
312 CTCF binding sites (rs117293173, rs13338055, rs74011247, and rs79521770). One TSP
313 (rs57790054) is an eQTL for the gene *KNOP1* (aka *C16orf88*), and 26 other LD SNPs are eQTL
314 for both *KNOP1* and *IQCK*. These two genes have both been associated with obsessive
315 compulsive disorder, among other diseases (Mattheisen et al., 2015). These results suggest that
316 effects on growth and BMI or on addictive behaviors could be under LTBS. We note that there is
317 some evidence of ethanol consumption in chimpanzees, but it is unclear how widespread its
318 availability was over the past several million years (Hockings et al., 2015).
319

320 *Urate levels.* The TSPs (rs1839333, rs1913638) on 8q21.13 are both significantly associated ($P < 2.0e-18$) with uric acid levels in multiple GWAS in European and Asian ancestry populations (Figure 5B) (Kanai et al., 2018; Kötgen et al., 2013; Tin et al., 2019). These variants are also associated with a range of body mass traits in the UK Biobank. Another variant in this locus (rs2941471, $R^2=0.97$ and $R^2=0.82$ in East Asians and Europeans respectively) is associated with pancreatic cancer ($p=7E-10$). Though elevated uric acid in the blood is associated with many conditions, it is a marker for pancreatic cancer (Stotz et al., 2014). This locus also contains two non-synonymous variants in the gene *HNF4G*, a transcription factor expressed in the liver, kidney, and pancreas, in high LD with the TSPs (rs1805098 and rs2943549). The TSPs are also expression and splicing QTL for *HNF4G* ($P = 1.9E-4$ and $8.6E-6$, respectively). Variants in *HNF4G* are associated with several traits, including the development of hyperuricemia (Chen et al., 2017).

332

A**B****C****Figure 5. Illustrative examples of non-immune functions associated with TSPs.**

335 (A) A TSP in 16p12.3 is associated with body mass index (BMI) and alcohol intake. Regional
336 association plot showing statistically significant genome- and phenome-wide associations
337 (threshold $p \leq 1E-08$), regulatory annotations from Ensembl, and eQTLs from GTEx. One of the
338 TSPs (rs57790054, orange) is associated with alcohol intake and several growth and body mass
339 phenotypes in the UK Biobank. A variant in high LD (rs72771074, green) has been associated
340 with alcohol use disorder in a previous GWAS. The TSP is also strongly associated with growth
341 (comparative body size at age 10, 9.6e-21) and body mass index (3.5e-12). The TSPs are nearby
342 *GPR139*, a gene encoding a G-protein coupled receptor expressed in the brain, whose expression
343 levels influence alcohol drinking behavior in rats. The TSP region also contains several CTCF
344 binding sites. TSP are shown in bold text.

345 (B) TSPs in 8q21.11 are associated with urate levels. Regional association plot showing
346 statistically significant genome- and phenome-wide associations ($P \leq 1E-08$), eQTL, regulatory
347 and coding SNPs. LD SNPs in this region are associated with urate levels and pancreatic cancer.
348 A TSP (rs1839333) is also associated with gout, although the p-value did not meet our strict
349 threshold. TSP are shown in bold text. Figures created using LocusZoom (Pruim et al., 2011).

350 (C) TSP locus on 3q24. Regional association plot showing statistically significant genome- and
351 phenome-wide associations (threshold $p \leq E-08$), regulatory and eQTLs. This locus is
352 characterized by neurological traits involved in educational attainment, cognitive performance,
353 and risky behavior (automobile speeding propensity).

354

355 *Chronotype*. The TSPs (rs1887944, rs7147645) on 14q31.2 are both nominally associated ($P \leq$
356 9.5e-6) with chronotype in a study of nearly 500,000 people (Jones et al., 2019). This is the only
357 phenotype association with these variants. One variant in high LD (rs17119051) overlaps a
358 regulatory region, but there are no eQTL among the TSPs or other LD variants. Asynchrony in
359 chronotypes in a population has been proposed to be beneficial due to potential protection
360 against predators and other dangers during the vulnerable hours of sleep. This so-called sentinel
361 hypothesis developed out of the observation in Hadza hunter-gatherers of Tanzania that due to
362 variation in sleep phase timing all members of a group are rarely asleep simultaneously (Samson
363 et al., 2017). Thus, it is possible that LTBS maintained variation at this locus to promote
364 chronotype diversity to mitigate the risks of sleeping.

365

366 *Risky behavior and educational attainment*. TSPs (rs9869178 and rs2118072) in 3q24
367 (hg19.chr3:143636420-143740729) are associated with a range of risky behaviors and
368 educational attainment in both individual GWAS studies and the UK Biobank (Figure 5C). For
369 example, they are associated with automobile speeding propensity ($P = 4.4E-8$) (Linnér, 2019)
370 and with educational attainment ($P = 2.4E-7$) (Lee et al., 2018). The TSPs are also modestly
371 associated with variation in brain white matter microstructure (Anterior corona radiata mean
372 diusivities, $P = 1.96E-6$) (Zhao et al., 2019). Many of the variants in high LD with the TSPs in
373 this region have similar associations and overlap annotated regulatory regions: regulatory open
374 chromatin (rs10662845 and rs6766439), promoter (rs1898263, rs1992094, rs4431106,

375 rs7631704, rs7651567, and rs7653431), and CTCF binding sites (rs7628282, rs7650239,
376 rs7650332, rs9840519, rs9840971, rs9878070, and *rs9840157*). Furthermore, the TSPs (and 159
377 high LD variants) are significant eQTLs ($P \leq 1E-5$) for the gene *DIPK2A* (*C3orf58*) across four
378 GTEx tissues (small intestine terminal ileum, transformed fibroblasts, skin from the lower leg,
379 and suprapubic skin). *DIPK2A* has not been comprehensively functionally characterized, but it
380 contains a protein kinase domain and is broadly expressed, including in the developing and adult
381 brain. Deletion of this gene has been linked to autism, and its expression is responsive to
382 neuronal activity (Morrow et al., 2008). Associations with behavioral and cognitive traits must
383 be interpreted with caution as these traits are very challenging to quantify and strongly
384 influenced by social factors that may vary with other characteristics. Nonetheless, these
385 associations point to an influence of the TSPs on behaviors relevant to risk tolerance. Thus, it is
386 possible that maintaining a diversity of risk tolerance in human and chimpanzee populations has
387 been beneficial.

388

389 **Discussion**

390 In this study we aimed to characterize the function of genomic regions with two or more TSPs in
391 LD and close proximity. These variants have a deep ancestry in the common ancestor between
392 humans and chimpanzees, and have persisted together in the genomes of both species for
393 millions of years. Due to the maintenance of these polymorphisms over such long periods, they
394 are likely evolving under LTBS. The majority of the non-coding TSP regions previously
395 identified do not have known functions. To address this challenge, we identified functional
396 annotations for 114 out of 125 TSP regions with the help of newly developed genomic
397 annotation tools (Figure 2).

398 Our results suggest that non-coding TSPs likely maintained by LTBS have diverse
399 functions beyond enabling a flexible immune response to pathogens. This expands on several
400 recent studies of balancing selection over shorter timescales that have identified a small number
401 of regions with functions outside the immune system (Bitarello et al., 2018; Sato & Kawata,
402 2018; Siewert & Voight, 2017; Viscardi et al., 2018).

403 The associations we identify suggest possible behavioral, neurological, and other traits
404 that may have driven LTBS. In particular, our results provide support and candidate loci for
405 previous hypotheses about the need for neurological and behavioral diversity in populations. For
406 example, the chronotype association supports the sentinel hypothesis that variability in sleep
407 patterns is the result of natural selection acting to reduce vulnerability of groups while members
408 sleep at night (Samson et al., 2017). Similarly, selection has recently been shown to act on risk-
409 taking behavior in anole lizards (Lapiedra et al., 2018). Thus, our identification of an association
410 between a TSP and human risk-taking behavior (Figure 4C) suggests that LTBS may have
411 maintained genetic variants that contribute to variation in risk taking behavior in humans and
412 chimpanzees. A gene with evidence of eQTL in this region (*C3orf58*) encodes for a protein
413 kinase and has been associated with autism and other neurological disorders (Dudkiewicz et al.,
414 2013), adding further evidence of possible neurological drivers to LTBS.

415 Our results also raise the intriguing possibility that variants that modulate urate levels
416 have been under LTBS. Uricase, the enzyme that metabolizes uric acid into an easily excreted
417 water-soluble form in most mammals, has been lost in great apes. This gene was disabled by a
418 series of mutations that slowly decreased activity over primate evolution, increasing the levels of
419 uric acid in blood (Kratzer et al., 2014). It has been hypothesized that this loss of uricase activity
420 was driven by increase fructose in primate diets due to fruit eating (Johnson et al., 2009; Kratzer
421 et al., 2014). It has also been proposed that high levels of uric acid, a potent antioxidant, played
422 an important role in the evolution of intelligence, acting as antioxidant in the brain (Álvarez-
423 Lario & Macarrón-Vicente, 2010). However, as reflected in the associations with this locus,
424 elevated uric acid levels contribute to many common diseases in modern humans, including
425 chronic hypertension, cardiovascular disease, kidney and liver diseases, metabolic syndrome,
426 diabetes, and obesity (Gustafsson & Unwin, 2013). This suggests potential functional tradeoffs at
427 this locus. However, we emphasize that proving the environmental drivers of past selection is
428 challenging.

429 Some of the phenotype associations we discovered may reflect manifestations of
430 variation on traits in modern environments that could not be long-term drivers of balancing
431 selection. As an extreme example, influence on smoking behavior could not have been the cause
432 of LTBS given the relatively recent wide availability of nicotine. Though we note that there is
433 some evidence of ethanol consumption in chimpanzees (Hockings et al., 2015). Even if they
434 reflect modern environments, these associations provide hints about possible behavioral,
435 neurological, or other traits that may have driven LTBS. For instance, plant chemicals can hijack
436 reward systems in the brain that motivate repetition and learning (U.S. Department of Health &
437 Human Services, 2016). The same systems that influence these action and consequently
438 reproductive fitness potentially created a byproduct of excessive seeking of dopamine or other
439 reward chemicals.

440 There are several caveats to our work. First, factors other than LTBS, such as high
441 mutation rates and sequencing errors, can produce a TSP (Cheng & DeGiorgio, 2019; Gao et al.,
442 2015). However, the presence of two or more TSPs in the regions we considered strongly
443 suggest LTBS, and 80% of the regions had evidence of LTBS from an additional prediction
444 method. Nonetheless, candidate regions of interest for future study should be further analyzed for
445 possible confounders. Second, even with recent growth of genetic and phenotypic databases, our
446 knowledge of the functions of most regions of the genome is sparse. Thus, failure to observe a
447 functional association does not imply that a region does not have an important function. Third,
448 the genome- and phenome-wide association tools we used are limited to the samples that have
449 been analyzed; available data do not represent the full scope of human variation. Most of the
450 individuals analyzed in available genetic association studies are of European ancestry. Variant
451 functions and the ability to detect associations vary across human populations; however, we
452 anticipate that TSPs should have functional effects across populations, unless modern
453 environments have masked the pressure driving LTBS. Fourth, even in PheWAS, a limited
454 number of phenotypes have been quantified across individuals, and these studies are focused on

455 a subset of clinically relevant rather evolutionarily relevant traits. Fifth, in some analyses, we
456 considered annotations based on trait associations with variants in high LD ($r^2 > 0.8$) with TSPs.
457 This could potentially introduce false positives if the variant also tags a different causal non-TSP
458 variant. Given the long-term selection on TSPs, they are strong candidates for causal variants,
459 but functional studies are needed to confirm these statistical associations. Finally, our analyses
460 have focused on the human context; due to lack of functional data, it is not possible to explore
461 the function of TSPs in chimpanzees. Nonetheless, we feel that our integration of genome-scale
462 annotations and biobank data highlight the diversity of functions associated with LTBS.

463

464 **Conclusions**

465 In conclusion, we assign putative functions for TSP regions that likely persisted due to balancing
466 selection dating back to at least the common ancestor of humans and chimpanzees. These
467 annotations expand beyond immune functions to traits relevant to behavior, cognition, and body
468 shape. Notably, we also find that most regions with multiple TSPs overlap gene regulatory
469 annotations suggesting LTBS on gene expression levels. As methods improve for quantifying the
470 effects of variants on gene regulation in different tissues and how these relate to organism-level
471 phenotypes, we anticipate deeper mechanistic understanding of the functions and potential
472 evolutionary pressures on these regions.

473

474

475 **Methods**

476 *Trans-species polymorphisms*

477 The initial set of 125 regions containing 263 TSPs analyzed in this study was published by
478 Leffler et al. (2013). The set is composed of regions that: 1) contain at least two trans-species
479 polymorphisms—i.e., variants that are segregating in both 51 Yoruba individuals in the 1000
480 Genomes Pilot 1 and 10 chimpanzees from the PanMap project—within 4 kb of each other in
481 both species; and 2) are in high LD in humans and chimpanzees.

482 To increase our ability to identify annotations in each locus, the dataset was expanded to
483 include variants in high LD (threshold $R^2=0.8$) with each of the TSPs as is common in
484 association studies. We computed linkage disequilibrium with the TSP variants from 1000
485 Genomes Project Phase 3 data using rAggr a web tool developed by the University of Southern
486 California (<http://raggr.usc.edu>). We considered LD in African, East Asia, and European
487 populations. Variants with no reported RSID were excluded from the analysis. The dataset was
488 thus expanded by 9,996 SNPs in high LD with the TSPs for a total of 10,259 SNPs.

489

490 *LTBS prediction comparison*

491 We compared the TSP regions in this study with LTBS candidate regions from two different
492 methods developed to detect long-term balancing selection. BetaScan2 (Siewert & Voight, 2020)
493 is a recent statistic for detecting balancing selection based enrichment for variants in a region
494 with low variation in allele frequency and a deficit of substitutions. We identified overlaps

495 between the TSP regions and genomic regions detected by BetaScan2. Among the TSPs with
496 Beta scores, 56% (57/102) had values greater than the 2.0 standardized beta score threshold used
497 by the authors. Considering the entire TSP regions and LD, 83% (92/111) had at least one variant
498 with a score above this threshold. We also computed overlap with regions identified by the NCD
499 statistic (Bitarello et al., 2018). The overlap with the regions detected by NCD is 28% (35/125
500 regions). In total, 80% (100/125) of the TSP regions were supported by either the BetaScan2 or
501 NCD.

502

503 *Genome- and Phenome-wide associations*

504 The GWAS Catalog collects variant-trait associations from published genome-wide association
505 studies. The database is currently composed of more than 200,000 associations. We used the
506 GWAS Catalog version October 2020 to find functional associations for the LTBS variants. The
507 search was done using the BEDTools intersect function between the GWAS catalog and the LD-
508 expanded TSP dataset (Quinlan, 2014).

509 PheWAS is an analysis strategy built on top of medical records with information about
510 patient phenotypes and associated variants. The geneAtlas catalog
[511 \(<http://geneatlas.roslin.ed.ac.uk/phewas/>\)](http://geneatlas.roslin.ed.ac.uk/phewas/) takes advantage of the data provided by the UK
512 Biobank cohort, which contains medically relevant data from nearly 500,000 British individuals
513 of European ancestry. This database contains 3 million variants in 778 traits. We matched our set
514 of variants against the geneAtlas database to search for traits associated with LTBS.

515

516 *GTEX eQTL data*

517 To evaluate potential gene regulatory effects of TSPs in non-coding regions, we analyzed data
518 from GTEx, a project developed to quantify the consequence of genetic variation on expression
519 at the tissue level ([520 <https://www.gtexportal.org/home/>](https://www.gtexportal.org/home/)). The GTEx project v8 data have identified
521 eQTL across 50 tissues based on analyses of nearly 1,000 individuals to identify differential
522 expression through SNP variation. The intersection between the TSPs and LD SNPs and the
523 GTEx eQTL returned a large collection of TSP eQTL.

524

525 *Other functional categories*

526 We used ANNOVAR to identify variants overlapping protein coding regions. We used the
527 Ensembl Regulatory Build (Zerbino et al., 2015) to identify variants overlapping regions with
528 regulatory function. In all, 58 TSPs in 40 regions and 1334 LD SNPs in 114 regions overlapped
529 regulatory annotations.

530

531 *Enrichment for overlap with regulatory regions*

532 We used a permutation framework to calculate whether TSPs were more enriched for overlap
533 with regulatory regions than expected by chance (Benton et al., 2019). We quantified the number
534 of overlapping TSPs for each type of regulatory region (open chromatin, promoter, enhancer,
promoter-flanking, CTCF binding site, TF binding site). We then compared the observed TSP

535 overlap to a null distribution of expected overlap generated by randomly shuffling the regulatory
536 regions 1000 times across the genome. We maintain the original length and chromosome
537 distributions for shuffled regions and exclude all ENCODE blacklist and gap regions (Kundaje,
538 2013). We then computed an empirical p-value for the observed TSP overlap based on the
539 distribution of overlaps for the set of matched shuffled regions.

540

541 **Declarations**

542

543 **Availability of Data and Materials**

544 The data underlying this article are available in the article and in its online supplementary
545 material.

546

547 **Competing interests**

548 The authors declare that they have no competing interests.

549

550 **Funding**

551 This work was supported by the National Institutes of Health [grant R35GM127087; grant
552 T32LM012412], and the Burroughs-Wellcome Fund. The funders did not play any role in the
553 study design, collection, analysis and interpretation of data, or in writing the manuscript.

554

555 **Author's contributions**

556 Conceptualization: JAC; Methodology: KV, MLB, JAC; Investigation: KV, MLB, JAC; Writing
557 – Original Draft: KV, JAC; Writing – Review & Editing: KV, MLB, JAC; Funding Acquisition:
558 JAC; Resources: JAC; Supervision: JAC.

559

560 **Abbreviations**

561 LTBS: Long-term balancing selection

562 TSP: Trans-species polymorphisms

563 LD: Linkage disequilibrium

564 eQTL: Expression quantitative trait loci

565 GTEx: Genotype-Tissue Expression

566 GWAS: Genome-wide association study

567 PheWAS: Phenome-wide association study

568

569 **Acknowledgements**

570 We thank Evonne McArthur, David Rinker, and other members of the Capra Lab for helpful
571 comments on this work. This work was conducted in part using the resources of the Advanced
572 Computing Center for Research and Education at Vanderbilt University, Nashville, TN.

573

574

- 575 REFERENCES
- 576 Álvarez-Lario, B., & Macarrón-Vicente, J. (2010). Uric acid and evolution. *Rheumatology*,
577 49(11), 2010–2015. <https://doi.org/10.1093/rheumatology/keq204>
- 578 Azevedo, L., Serrano, C., Amorim, A., & Cooper, D. N. (2015). Trans-species polymorphism in
579 humans and the great apes is generally maintained by balancing selection that modulates the
580 host immune response. *Human Genomics*, 9(1). <https://doi.org/10.1186/s40246-015-0043-1>
- 581 Battivelli, E., Migraine, J., Lecossier, D., Yeni, P., Clavel, F., & Hance, A. J. (2011). Gag
582 Cytotoxic T Lymphocyte Escape Mutations Can Increase Sensitivity of HIV-1 to Human
583 TRIM5 , Linking Intrinsic and Acquired Immunity. *Journal of Virology*, 85(22), 11846–
584 11854. <https://doi.org/10.1128/jvi.05201-11>
- 585 Benton, M. L., Talipineni, S. C., Kostka, D., & Capra, J. A. (2019). Genome-wide enhancer
586 annotations differ significantly in genomic distribution, evolution, and function. *BMC*
587 *Genomics*, 20(1), 1–22. <https://doi.org/10.1186/s12864-019-5779-x>
- 588 Bitarello, B. D., De Filippo, C., Teixeira, J. C., Schmidt, J. M., Kleinert, P., Meyer, D., &
589 Andres, A. M. (2018). Signatures of long-term balancing selection in human genomes.
590 *Genome Biology and Evolution*, 10(3), 939–955. <https://doi.org/10.1093/gbe/evy054>
- 591 Bush, W. S., Oetjens, M. T., & Crawford, D. C. (2016). Unravelling the human genome-
592 phenotype relationship using genome-wide association studies. In *Nature Reviews Genetics*
593 (Vol. 17, Issue 3, pp. 129–145). Nature Publishing Group.
594 <https://doi.org/10.1038/nrg.2015.36>
- 595 Cagliani, R., Fumagalli, M., Biasin, M., Piacentini, L., Riva, S., Pozzoli, U., Bonaglia, M. C.,
596 Bresolin, N., Clerici, M., & Sironi, M. (2010). Long-term balancing selection maintains
597 trans-specific polymorphisms in the human TRIM5 gene. *Human Genetics*, 128(6), 577–
598 588. <https://doi.org/10.1007/s00439-010-0884-6>
- 599 Cagliani, R., Guerini, F. R., Fumagalli, M., Riva, S., Agliardi, C., Galimberti, D., Pozzoli, U.,
600 Goris, A., Dubois, B., Fenoglio, C., Forni, D., Sanna, S., Zara, I., Pitzalis, M.,
601 Zoledziewska, M., Cucca, F., Marini, F., Comi, G. P., Scarpini, E., ... Sironi, M. (2012). A
602 trans-specific polymorphism in ZC3HAV1 is maintained by long-standing balancing
603 selection and may confer susceptibility to multiple sclerosis. *Molecular Biology and*
604 *Evolution*, 29(6), 1599–1613. <https://doi.org/10.1093/molbev/mss002>
- 605 Chen, B. D., Chen, X. C., Pan, S., Yang, Y. N., He, C. H., Liu, F., Ma, X., Gai, M. T., & Ma, Y.
606 T. (2017). TT genotype of rs2941484 in the human HNF4G gene is associated with
607 hyperuricemia in Chinese Han men. *Oncotarget*, 8(16), 26918–26926.
608 <https://doi.org/10.18632/oncotarget.15851>
- 609 Cheng, X., & DeGiorgio, M. (2019). Detection of Shared Balancing Selection in the Absence of
610 Trans-Species Polymorphism. *Molecular Biology and Evolution*, 36(1), 177–199.
611 <https://doi.org/10.1093/molbev/msy202>
- 612 De Filippo, C., Key, F. M., Ghirotto, S., Benazzo, A., Meneu, J. R., Weihmann, A., Parra, G.,
613 Green, E. D., & Andrés, A. M. (2016). Recent Selection Changes in Human Genes under
614 Long-Term Balancing Selection. *Molecular Biology and Evolution*, 33(6), 1435–1447.
615 <https://doi.org/10.1093/molbev/msw023>
- 616 DeGiorgio, M., Lohmueller, K. E., & Nielsen, R. (2014). A Model-Based Approach for
617 Identifying Signatures of Ancient Balancing Selection in Genetic Data. *PLoS Genetics*,
618 10(8). <https://doi.org/10.1371/journal.pgen.1004561>
- 619 Dudkiewicz, M., Lenart, A., & Pawłowski, K. (2013). A Novel Predicted Calcium-Regulated
620 Kinase Family Implicated in Neurological Disorders. *PLoS ONE*, 8(6).

- 621 https://doi.org/10.1371/journal.pone.0066427
622 Ganser-Pornillos, B. K., & Pornillos, O. (2019). Restriction of HIV-1 and other retroviruses by
623 TRIM5. In *Nature Reviews Microbiology* (Vol. 17, Issue 9, pp. 546–556). Nature
624 Publishing Group. https://doi.org/10.1038/s41579-019-0225-2
625 Gao, Z., Przeworski, M., & Sella, G. (2015). Footprints of ancient-balanced polymorphisms in
626 genetic variation data from closely related species. *Evolution*, 69(2).
627 https://doi.org/10.1111/evo.12567
628 Gustafsson, D., & Unwin, R. (2013). The pathophysiology of hyperuricaemia and its possible
629 relationship to cardiovascular disease, morbidity and mortality. In *BMC Nephrology* (Vol.
630 14, Issue 1). https://doi.org/10.1186/1471-2369-14-164
631 Hockings, K. J., Bryson-Morrison, N., Carvalho, S., Fujisawa, M., Humle, T., McGrew, W. C.,
632 Nakamura, M., Ohashi, G., Yamanashi, Y., Yamakoshi, G., & Matsuzawa, T. (2015). Tools
633 to tipple: Ethanol ingestion by wild chimpanzees using leaf-sponges. *Royal Society Open
634 Science*, 2(6). https://doi.org/10.1098/rsos.150150
635 Johnson, R. J., Sautin, Y. Y., Oliver, W. J., Roncal, C., Mu, W., Gabriela Sanchez-Lozada, L.,
636 Rodriguez-Iturbe, B., Nakagawa, T., & Benner, S. A. (2009). Lessons from comparative
637 physiology: Could uric acid represent a physiologic alarm signal gone awry in western
638 society? In *Journal of Comparative Physiology B: Biochemical, Systemic, and
639 Environmental Physiology* (Vol. 179, Issue 1, pp. 67–76). Springer Verlag.
640 https://doi.org/10.1007/s00360-008-0291-7
641 Jones, S. E., Lane, J. M., Wood, A. R., van Hees, V. T., Tyrrell, J., Beaumont, R. N., Jeffries, A.
642 R., Dashti, H. S., Hillsdon, M., Ruth, K. S., Tuke, M. A., Yaghootkar, H., Sharp, S. A., Jie,
643 Y., Thompson, W. D., Harrison, J. W., Dawes, A., Byrne, E. M., Tiemeier, H., ... Weedon,
644 M. N. (2019). Genome-wide association analyses of chronotype in 697,828 individuals
645 provides insights into circadian rhythms. *Nature Communications*, 10(343).
646 https://doi.org/10.1038/s41467-018-08259-7
647 Kanai, M., Akiyama, M., Takahashi, A., Matoba, N., Momozawa, Y., Ikeda, M., Iwata, N.,
648 Ikegawa, S., Hirata, M., Matsuda, K., Kubo, M., Okada, Y., & Kamatani, Y. (2018).
649 Genetic analysis of quantitative traits in the Japanese population links cell types to complex
650 human diseases. *Nature Genetics*, 50(3), 390–400. https://doi.org/10.1038/s41588-018-
651 0047-6
652 Key, F. M., Teixeira, J. C., de Filippo, C., & Andrés, A. M. (2014). Advantageous diversity
653 maintained by balancing selection in humans. *Current Opinion in Genetics and
654 Development*, 29, 45–51. https://doi.org/10.1016/j.gde.2014.08.001
655 Kononoff, J., Kallupi, M., Kimbrough, A., Conlisk, D., de Guglielmo, G., & George, O. (2018).
656 Systemic and intra-habenular activation of the orphan G protein-coupled receptor GPR139
657 decreases compulsive-like alcohol drinking and hyperalgesia in alcohol-dependent rats.
658 *ENeuro*, 5(3). https://doi.org/10.1523/ENEURO.0153-18.2018
659 Köttgen, A., Albrecht, E., Teumer, A., Vitart, V., Krumsiek, J., Hundertmark, C., Pistis, G.,
660 Ruggiero, D., O'Seaghda, C. M., Haller, T., Yang, Q., Tanaka, T., Johnson, A. D., Kutalik,
661 Z., Smith, A. V., Shi, J., Struchalin, M., Middelberg, R. P. S., Brown, M. J., ... Gieger, C.
662 (2013). Genome-wide association analyses identify 18 new loci associated with serum urate
663 concentrations. *Nature Genetics*, 45(2), 145–154. https://doi.org/10.1038/ng.2500
664 Kratzer, J. T., Lanaspa, M. A., Murphy, M. N., Cicerchi, C., Graves, C. L., Tipton, P. A.,
665 Ortlund, E. A., Johnson, R. J., & Gaucher, E. A. (2014). Evolutionary history and metabolic
666 insights of ancient mammalian uricases. *Proceedings of the National Academy of Sciences*

- 667 *of the United States of America*, 111(10), 3763–3768.
668 <https://doi.org/10.1073/pnas.1320393111>
- 669 Kundaje, A. (2013). *A comprehensive collection of signal artifact blacklist regions in the human*
670 *genome*.
671 ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/supplementary/integration_data_jan2011/b
672 *yFreeze/jan2011/blacklists/hg19-blacklist-README.pdf*,
673 <https://sites.google.com/site/anshulkundaje/projects/blacklists>
- 674 Lapiendra, O., Schoener, T. W., Leal, M., Losos, J. B., & Kolbe, J. J. (2018). Predator-driven
675 natural selection on risk-taking behavior in anole lizards. *Science*, 360(6392), 1017–1020.
676 <https://doi.org/10.1126/science.aap9289>
- 677 Lawlor, D. A., Ward, F. E., Ennis, P. D., Jackson, A. P., & Parham, P. (1988). HLA-A and B
678 polymorphisms predate the divergence of humans and chimpanzees. *Nature*, 335(6187),
679 268–271. <https://doi.org/10.1038/335268a0>
- 680 Lee, J. J., Wedow, R., Okbay, A., Kong, E., Maghzian, O., Zacher, M., Nguyen-Viet, T. A.,
681 Bowers, P., Sidorenko, J., Karlsson Linnér, R., Fontana, M. A., Kundu, T., Lee, C., Li, H.,
682 Li, R., Royer, R., Timshel, P. N., Walters, R. K., Willoughby, E. A., ... Turley, P. (2018).
683 Gene discovery and polygenic prediction from a genome-wide association study of
684 educational attainment in 1.1 million individuals. *Nature Genetics*, 50(8), 1112–1121.
685 <https://doi.org/10.1038/s41588-018-0147-3>
- 686 Leffler, E. M., Gao, Z., Pfeifer, S., Ségurel, L., Auton, A., Venn, O., Bowden, R., Bontrop, R.,
687 Wall, J. D., Sella, G., Donnelly, P., McVean, G., & Przeworski, M. (2013). Multiple
688 instances of ancient balancing selection shared between humans and chimpanzees. *Science*,
689 340(6127), 1578–1582. <https://doi.org/10.1126/science.1234070>
- 690 Linnér, R. K. (2019). Genome-wide association analyses of risk tolerance and risky behaviors in
691 over 1 million individuals identify hundreds of loci and shared genetic influences. *Nature*
692 *Genetics*, 51(2), 245–257. <https://doi.org/10.1038/s41588-018-0309-3>
- 693 Mao, R., Nie, H., Cai, D., Zhang, J., Liu, H., Yan, R., Cuconati, A., Block, T. M., Guo, J. T., &
694 Guo, H. (2013). Inhibition of Hepatitis B Virus Replication by the Host Zinc Finger
695 Antiviral Protein. *PLoS Pathogens*, 9(7). <https://doi.org/10.1371/journal.ppat.1003494>
- 696 Mattheisen, M., Samuels, J. F., Wang, Y., Greenberg, B. D., Fyer, A. J., Mccracken, J. T.,
697 Geller, D. A., Murphy, D. L., Knowles, J. A., Grados, M. A., Riddle, M. A., Rasmussen, S.
698 A., McLaughlin, N. C., Nurmi, E. L., Askland, K. D., Qin, H. D., Cullen, B. A., Piacentini,
699 J., Pauls, D. L., ... Nestadt, G. (2015). Genome-wide association study in obsessive-
700 compulsive disorder: Results from the OCGAS. *Molecular Psychiatry*, 20(3).
701 <https://doi.org/10.1038/mp.2014.43>
- 702 Mayer, W. E., Jonker, M., Klein, D., Ivanyi, P., van Seventer, G., & Klein, J. (1988). Nucleotide
703 sequences of chimpanzee MHC class I alleles: evidence for trans-species mode of
704 evolution. *The EMBO Journal*, 7(9), 2765–2774. <https://doi.org/10.1002/j.1460-2075.1988.tb03131.x>
- 706 Morrow, E. M., Yoo, S. Y., Flavell, S. W., Kim, T. K., Lin, Y., Hill, R. S., Mukaddes, N. M.,
707 Balkhy, S., Gascon, G., Hashmi, A., Al-Saad, S., Ware, J., Joseph, R. M., Greenblatt, R.,
708 Gleason, D., Ertelt, J. A., Apse, K. A., Bodell, A., Partlow, J. N., ... Walsh, C. A. (2008).
709 Identifying autism loci and genes by tracing recent shared ancestry. *Science*, 321(5886),
710 218–223. <https://doi.org/10.1126/science.1157657>
- 711 Pruim, R. J., Welch, R. P., Sanna, S., Teslovich, T. M., Chines, P. S., Gliedt, T. P., Boehnke, M.,
712 Abecasis, G. R., Willer, C. J., & Frishman, D. (2011). LocusZoom: Regional visualization

- 713 of genome-wide association scan results. *Bioinformatics*, 27(13), 2336–2337.
714 <https://doi.org/10.1093/bioinformatics/btq419>
- 715 Quinlan, A. R. (2014). BEDTools: The Swiss-Army tool for genome feature analysis. *Current*
716 *Protocols in Bioinformatics*, 47(1), 11–12. <https://doi.org/10.1002/0471250953.bi1112s47>
- 717 Samson, D. R., Crittenden, A. N., Mabulla, I. A., Mabulla, A. Z. P., & Nunn, C. L. (2017).
718 Chronotype variation drives night-time sentinel-like behaviour in hunter-gatherers.
719 *Proceedings of the Royal Society B: Biological Sciences*, 284(20170967).
720 <https://doi.org/10.1098/rspb.2017.0967>
- 721 Sanchez-Roige, S., Palmer, A. A., Fontanillas, P., Elson, S. L., Adams, M. J., Howard, D. M.,
722 Edenberg, H. J., Davies, G., Crist, R. C., Deary, I. J., McIntosh, A. M., Clarke, T. K., Elson,
723 L., Fontanillas, P., Furlotte, N. A., Hinds, D. A., Huber, K. E., Kleinman, A., Litterman, N.
724 K., ... Wilson, C. H. (2019). Genome-wide association study meta-analysis of the alcohol
725 use disorders identification test (AUDIT) in two population-based cohorts. *American*
726 *Journal of Psychiatry*, 176(2), 107–118. <https://doi.org/10.1176/appi.ajp.2018.18040369>
- 727 Sato, D. X., & Kawata, M. (2018). Positive and balancing selection on SLC18A1 gene
728 associated with psychiatric disorders and human-unique personality traits. *Evolution*
729 Letters, 2(5), 499–510. <https://doi.org/10.1002/evl3.81>
- 730 Ségurel, L., Thompson, E. E., Flutre, T., Lovstad, J., Venkat, A., Margulis, S. W., Moyse, J.,
731 Ross, S., Gamble, K., Sella, G., Ober, C., & Przeworski, M. (2012). The ABO blood group
732 is a trans-species polymorphism in primates. *Proceedings of the National Academy of*
733 *Sciences of the United States of America*, 109(45), 18493–18498.
734 <https://doi.org/10.1073/pnas.1210603109>
- 735 Siewert, K. M., & Voight, B. F. (2017). Detecting long-term balancing selection using allele
736 frequency correlation. *Molecular Biology and Evolution*, 34(11), 2996–3005.
737 <https://doi.org/10.1093/molbev/msx209>
- 738 Siewert, K. M., & Voight, B. F. (2020). BetaScan2: Standardized Statistics to Detect Balancing
739 Selection Utilizing Substitution Data. *Genome Biology and Evolution*, 12(2), 3873–3877.
740 <https://doi.org/10.1093/gbe/eva013>
- 741 Stotz, M., Szkandera, J., Seidel, J., Stojakovic, T., Samonigg, H., Reitz, D., Gary, T., Kornprat,
742 P., Schaberl-Moser, R., Hoefler, G., Gerger, A., & Pichler, M. (2014). Evaluation of uric
743 acid as a prognostic blood-based marker in a large cohort of pancreatic cancer patients.
744 *PLoS ONE*, 9(8). <https://doi.org/10.1371/journal.pone.0104730>
- 745 Teixeira, J. C., De Filippo, C., Weihmann, A., Meneu, J. R., Racimo, F., Dannemann, M.,
746 Nickel, B., Fischer, A., Halbwax, M., Andre, C., Atencia, R., Meyer, M., Parra, G., Pääbo,
747 S., & Andrés, A. M. (2015). Long-term balancing selection in LAD1 maintains a missense
748 trans-species polymorphism in humans, chimpanzees, and bonobos. *Molecular Biology and*
749 *Evolution*, 32(5), 1186–1196. <https://doi.org/10.1093/molbev/msv007>
- 750 Tin, A., Marten, J., Halperin Kuhns, V. L., Li, Y., Wuttke, M., Kirsten, H., Sieber, K. B., Qiu,
751 C., Gorski, M., Yu, Z., Giri, A., Sveinbjornsson, G., Li, M., Chu, A. Y., Hoppmann, A.,
752 O'Connor, L. J., Prins, B., Nutile, T., Noce, D., ... Köttgen, A. (2019). Target genes,
753 variants, tissues and transcriptional pathways influencing human serum urate levels. *Nature*
754 *Genetics*, 51(10), 1459–1474. <https://doi.org/10.1038/s41588-019-0504-x>
- 755 Todorova, T., Bock, F. J., & Chang, P. (2015). Poly(ADP-ribose) polymerase-13 and RNA
756 regulation in immunity and cancer. In *Trends in Molecular Medicine* (Vol. 21, Issue 6, pp.
757 373–384). Elsevier Ltd. <https://doi.org/10.1016/j.molmed.2015.03.002>
- 758 U.S. Department of Health & Human Services. (2016). Chapter 2: The neurobiology of

- 759 substance use, misuse, and addiction. In *Facing Addiction in America: The Surgeon
760 General's Report on Alcohol, Drugs, and Health*.
- 761 Viscardi, L. H., Paixão-Côrtes, V. R., Comas, D., Salzano, F. M., Rovaris, D., Bau, C. D.,
762 Amorim, C. E. G., & Bortolini, M. C. (2018). Searching for ancient balanced
763 polymorphisms shared between Neanderthals and modern humans. *Genetics and Molecular
764 Biology*, 41(1), 67–81. <https://doi.org/10.1590/1678-4685-gmb-2017-0308>
- 765 Watanabe, K., Stringer, S., Frei, O., Umićević Mirkov, M., de Leeuw, C., Polderman, T. J. C.,
766 van der Sluis, S., Andreassen, O. A., Neale, B. M., & Posthuma, D. (2019). A global
767 overview of pleiotropy and genetic architecture in complex traits. *Nature Genetics*, 51(9).
768 <https://doi.org/10.1038/s41588-019-0481-0>
- 769 Zerbino, D. R., Wilder, S. P., Johnson, N., Juettemann, T., & Flieck, P. R. (2015). The Ensembl
770 Regulatory Build. *Genome Biology*, 16(1). <https://doi.org/10.1186/s13059-015-0621-5>
- 771 Zhao, B., Zhang, J., Ibrahim, J. G., Luo, T., Santelli, R. C., Li, Y., Li, T., Shan, Y., Zhu, Z.,
772 Zhou, F., Liao, H., Nichols, T. E., & Zhu, H. (2019). Large-scale GWAS reveals genetic
773 architecture of brain white matter microstructure and genetic overlap with cognitive and
774 mental health traits (n = 17,706). *Molecular Psychiatry*. <https://doi.org/10.1038/s41380-019-0569-z>
- 775
- 776

Figures

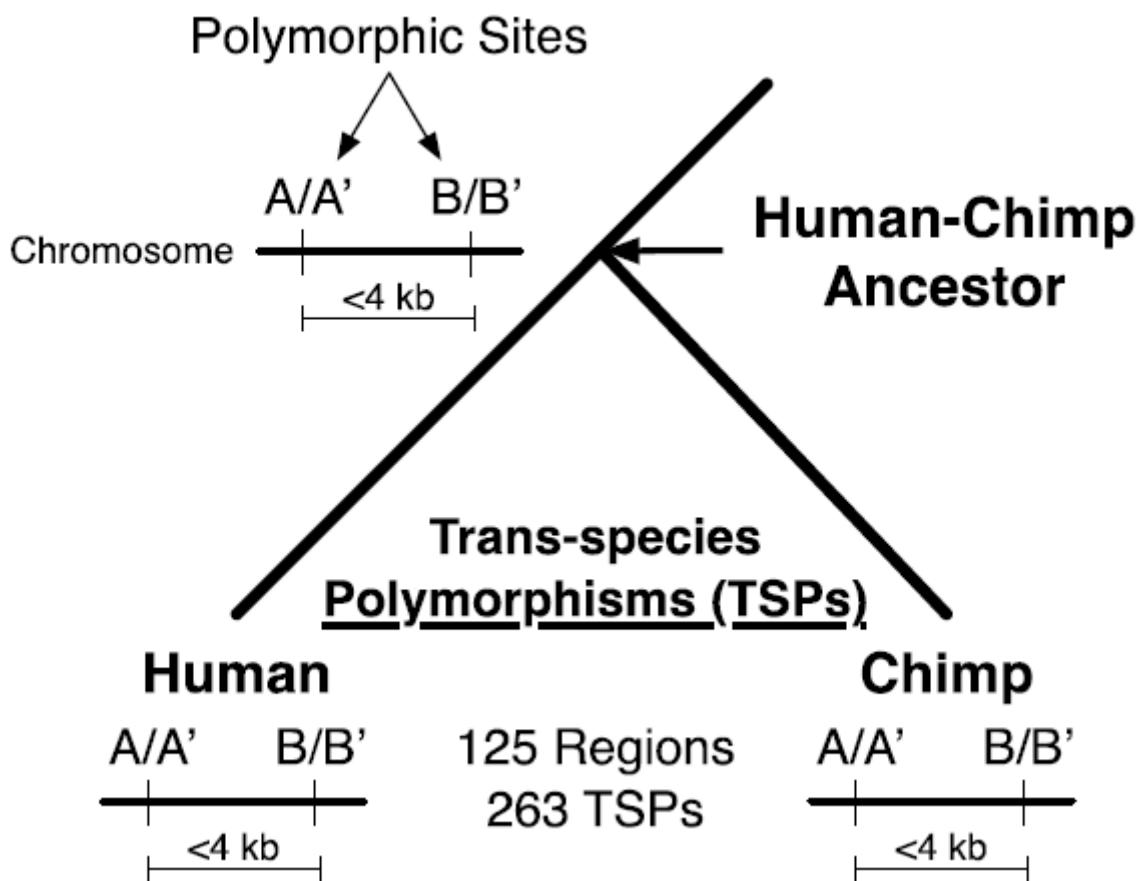


Figure 1

Trans-species polymorphisms (TSPs) likely resulting from long-term balancing selection (LTBS). Schematic showing the criteria used by Leffler et al. (2013) to identify TSPs likely maintained by LTBS. Each line represents a chromosome with polymorphisms segregating in a population. A/A' are two alleles segregating in both humans and chimpanzees at one site (i.e., a TSP), and B/B' are alleles segregating in both species at a nearby site. TSPs are very unlikely to appear nearby (within 4 kb) without the action of balancing selection. We consider 125 TSP regions containing 263 TSPs. Within these regions, multiple functional scenarios are possible. For example, one TSP may be under LTBS while the other is neutral, but maintained due to tight linkage. Alternatively, the TSPs may have epistatic functions and both be under selection. In addition to the 263 TSPs, we also considered functional associations with 9,996 variants in high LD ($r^2 > 0.8$) with a TSP at least one population from the 1000 Genomes Project (Supplementary Figure 1).

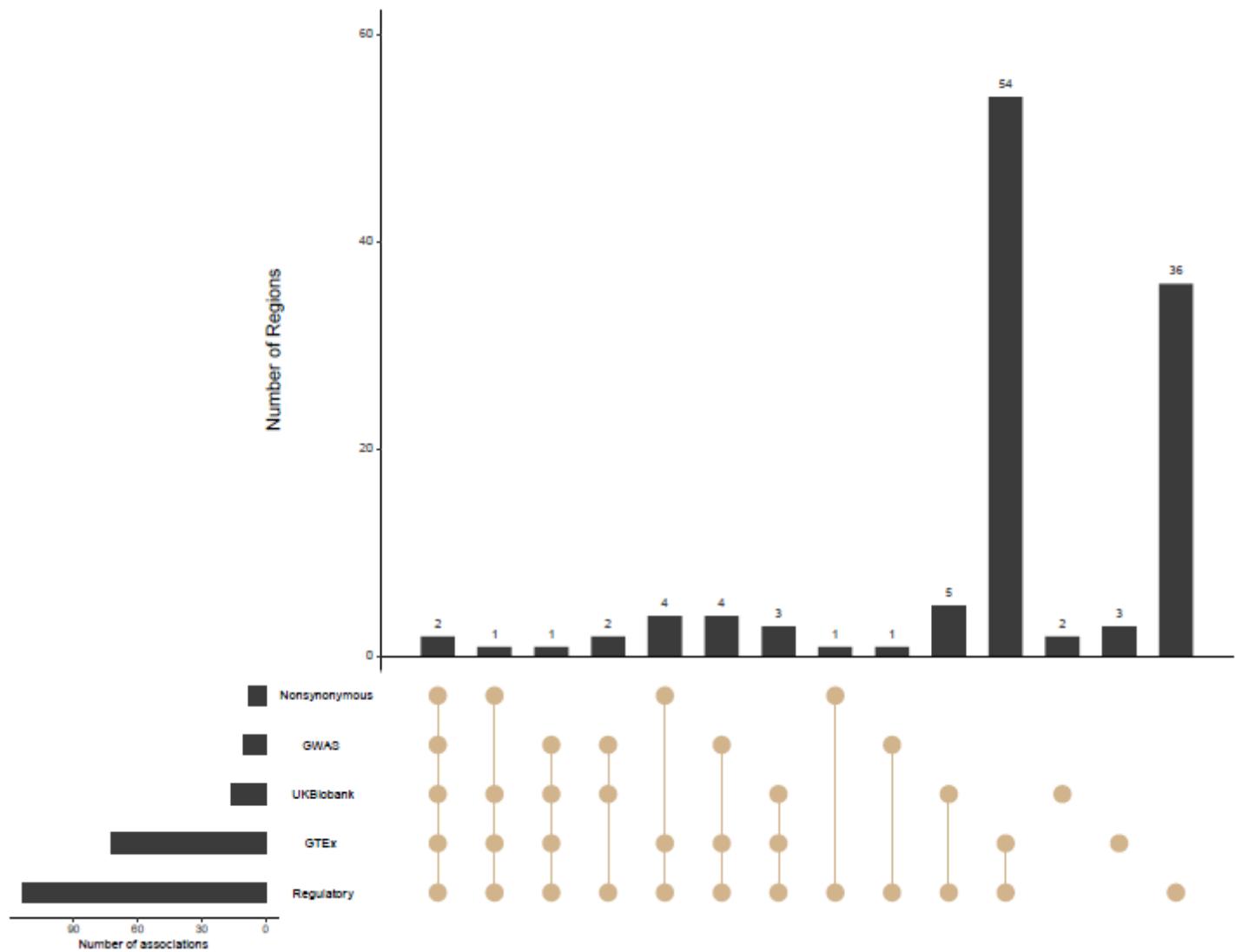


Figure 2

Functional annotations available for the expanded TSP regions. Summary of the annotations of each type available for TSP regions, including tagging variants in high LD with TSPs. A total of 119 out of 125 TSP regions contain at least one line of functional evidence. Multiple lines (two or more) are available for 78 regions.

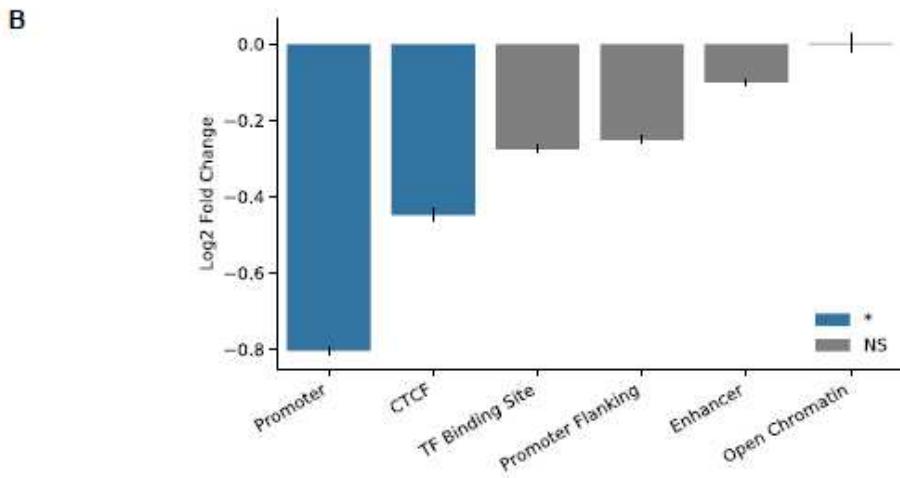
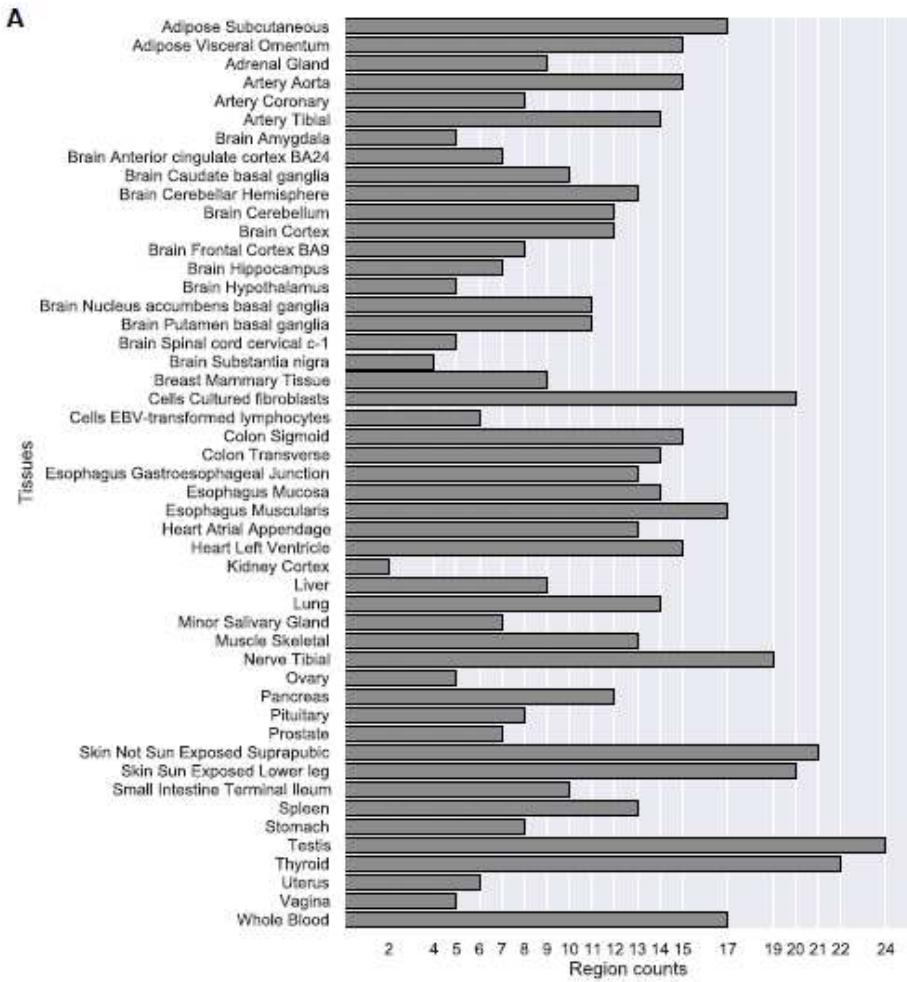


Figure 3

TSPs are eQTLs in diverse tissues and are depleted for overlap with promoters and CTCF sites. (A) The number of TSP regions that contain an eQTL for each GTEx tissue. Variation in TSP regions associates with gene expression in diverse tissues. The associated genes also have Gene Ontology (GO) annotations from diverse functional categories (Supplementary Figure 2). (B) TSP regions are significantly depleted

for overlap with promoters and CTCF sites compared to length- and chromosome-matched non-coding regions from the genomic background. The error bars represent 95% confidence intervals.



Figure 4

Genome- and phenotype-wide association studies link TSPs to diverse traits. Genome-wide significant ($P < 1E-8$) associations from the GWAS Catalog (yellow), a PheWAS over the UK Biobank from the geneAtlas (purple), and other studies summarized in the GWAS Atlas (green) (Watanabe et al., 2019). Each dot represents an association between a TSP region and a trait. Many immune-related traits are associated with TSPs, but there are also associations with a wider variety of phenotypes including osseous, neurological, and nervous system traits. Five extreme associations with immunological and blood traits ($P < 1E-60$) were truncated for this visualization. Since few TSPs themselves were directly tested in GWAS, we include GWAS Catalog associations with tag variants in high LD ($r^2 > 0.8$) with TSPs. All associations are listed in Supplementary Tables 2–4.

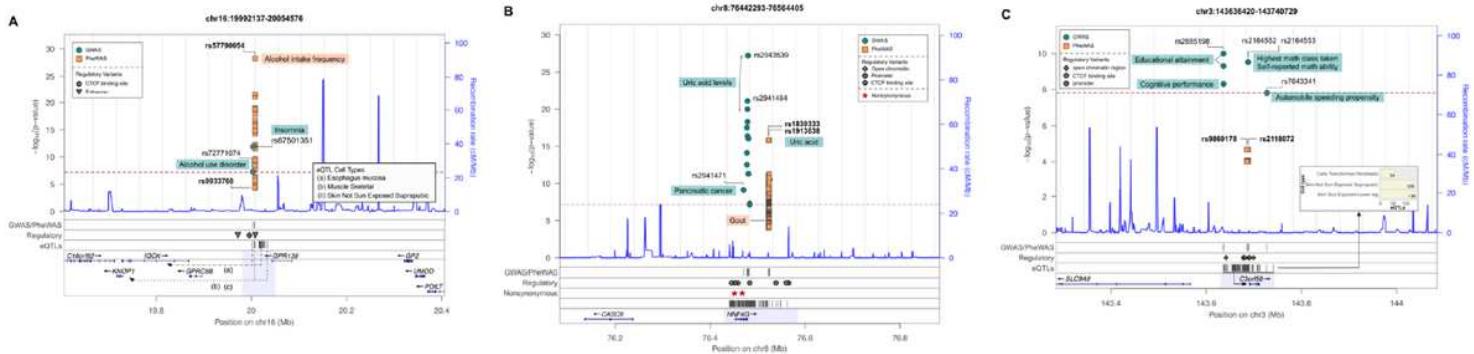


Figure 5

Illustrative examples of non-immune functions associated with TSPs. (A) A TSP in 16p12.3 is associated with body mass index (BMI) and alcohol intake. Regional association plot showing statistically significant genome- and phenotype-wide associations ($p \leq 1E-08$), regulatory annotations from Ensembl, and eQTLs from GTEx. One of the TSPs (rs57790054, orange) is associated with alcohol intake and several growth and body mass phenotypes in the UK Biobank. A variant in high LD (rs72771074, green) has been associated with alcohol use disorder in a previous GWAS. The TSP is also strongly associated with growth (comparative body size at age 10, $9.6e-21$) and body mass index ($3.5e-12$). The TSPs are nearby GPR139, a gene encoding a G-protein coupled receptor expressed in the brain, whose expression levels influence alcohol drinking behavior in rats. The TSP region also contains several CTCF binding sites. TSP are shown in bold text. (B) TSPs in 8q21.11 are associated with urate levels. Regional association plot showing statistically significant genome- and phenotype-wide associations ($P \leq 1E-08$), eQTL, regulatory and coding SNPs. LD SNPs in this region are associated with urate levels and pancreatic cancer. A TSP (rs1839333) is also associated with gout, although the p-value did not meet our strict threshold. TSP are shown in bold text. Figures created using LocusZoom (Pruim et al., 2011). (C) TSP locus on 3q24. Regional association plot showing statistically significant genome- and phenotype-wide associations ($p \leq E-08$), regulatory and eQTLs. This locus is characterized by neurological traits involved in educational attainment, cognitive performance, and risky behavior (automobile speeding propensity).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplement.pdf](#)
- [supplementarytables.xlsx](#)