

Opposing selective forces operating on human-specific duplicated TCAF genes in Neanderthals and humans

PingHsun Hsieh

University of Washington School of Medicine

Vy Dang

University of Washington School of Medicine

Mitchell Vollger

University of Washington School of Medicine

Yafei Mao

University of Washington

Tzu-Hsueh Huang

University of Washington School of Medicine

Philip Dishuck

Univ of WA <https://orcid.org/0000-0003-2223-9787>

Carl Baker

University of Washington School of Medicine

Stuart Cantsilieris

University of Washington School of Medicine

Alexandra Lewis

University of Washington School of Medicine

Katherine Munson

University of Washington <https://orcid.org/0000-0001-8413-6498>

Melanie Scofield

University of British Columbia

AnneMarie Welch

University of Washington School of Medicine

Jason Underwood

Pacific Biosciences (United States)

Evan Eichler (✉ eee@gs.washington.edu)

University of Washington School of Medicine <https://orcid.org/0000-0002-8246-4014>

Keywords: Human-specific genes, thermal regulatory gene duplications, natural selection, archaic and modern humans

Posted Date: October 1st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-77798/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on August 25th, 2021. See the published version at <https://doi.org/10.1038/s41467-021-25435-4>.

1 **Opposing selective forces operating on human-specific duplicated *TCAF* genes in Neanderthals and**
2 **humans**

3 PingHsun Hsieh,^{1*} Vy Dang,^{1,†} Mitchell R. Vollger,¹ Yafei Mao,¹ Tzu-Hsueh Huang,¹ Philip C. Dishuck,¹
4 Carl Baker¹, Stuart Cantsilieris,^{1,‡} Alexandra P. Lewis,¹ Katherine M. Munson,¹ Melanie Sorensen,¹
5 AnneMarie E. Welch,^{1,‡} Jason G. Underwood,^{1,2} Evan E. Eichler^{1,3*}

6
7

8 *Co-corresponding authors

9

10 Correspondence: hsiehph@uw.edu

11 eee@gs.washington.edu

12 Affiliations:

13 1. Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA

14 2. Pacific Biosciences (PacBio) of California, Incorporated, Menlo Park, CA, USA

15 3. Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA

16 † Current address: Institute for Cell and Molecular Biology, University of Texas, Austin, TX, USA

17 ‡ Current address: Centre for Eye Research Australia, Department of Surgery (Ophthalmology),

18 University of Melbourne, Royal Victorian Eye and Ear Hospital, East Melbourne, VIC, Australia

19 ‡ Current address: Brain and Mitochondrial Research, Murdoch Children's Research Institute, Royal

20 Children's Hospital, Melbourne, VIC, Australia

21

22 Keywords: Human-specific genes, thermal regulatory gene duplications, natural selection, archaic and

23 modern humans

24

25 **SUMMARY**

26 TRP channel-associated factor 1/2 (TCAF1/TCAF2) proteins antagonistically regulate the cold-sensor
27 protein TRPM8 in multiple human tissues. Understanding their significance has been complicated given
28 the locus spans a gap-ridden region with complex segmental duplications in GRCh38. Using long-read
29 sequencing, we sequence-resolve the locus, annotate full-length *TCAF* models in human and nonhuman
30 primate genomes, and show substantial human-specific *TCAF* copy number variation. We identify two
31 human super haplogroups, H4 and H5, and establish that *TCAF* duplications originated ~1.7 million years
32 ago but diversified only in *Homo sapiens* by recurrent structural mutations that altered *TCAF* copy
33 number and regulation. Conversely, in all archaic-hominin samples the fixation for a specific H4
34 haplotype without duplication is likely due to positive selection. The significant, positive effect of H4 on
35 *TCAF2* expression in modern-day humans with candidate associations for hypothyroidism, nerve
36 compression, and diabetes suggests *TCAF* diversification among hominins potentially in response to cold
37 or dietary adaptations.

38

39 INTRODUCTION

40 Gene duplication contributes significantly to molecular evolution by providing the raw material for
41 genetic novelty and organismal adaptation^{1,2}. In the human lineage, recently duplicated regions
42 corresponding to segmental duplications (SDs) are known to give rise to new genes contributing to
43 synaptogenesis, neuronal migration, and neocortical expansion³⁻⁷ that distinguish human from other ape
44 species. Although relatively few in number, some of the largest genetic changes that differentiate humans
45 from archaic hominins involve gene-rich SDs^{8,9}. Among human-specific SD genes, the duplications of
46 *TCAF1/TCAF2* are particularly intriguing. These genes encode TRP channel-associated factors that bind
47 to the protein TRPM8 (transient receptor potential cation channel subfamily M member 8)—the ion
48 channel acting as thermal sensor in peripheral somatosensory neurons¹⁰ and are thought to be under
49 positive selection in Eurasian populations¹¹. Both TCAF1 and TCAF2 proteins interact directly with the
50 TRPM8 channel but have antagonistic effects in TRPM8 gating and trafficking into the plasma
51 membrane: while TCAF1 facilitates the TRPM8 channel opening and migration, the activity of TRPM8 is
52 completely suppressed by TCAF2¹². These results suggest that the relative abundance of TCAF proteins
53 and their competition in binding TRPM8 is likely critical in physiological regulation of the channel
54 activity¹².

55 The *TCAF* family originated from an ancient gene duplication event at the basal of mammalian
56 phylogeny and stayed as single-copy genes throughout much of their evolution^{8,12}. Within the human
57 lineage, subsequent duplications over the last few million years have changed the copy number of both
58 *TCAF1* and *TCAF2*. In the human reference genome GRCh38, *TCAF1* and *TCAF2* are embedded and
59 span within a complex region of large, highly identical SDs (>99.5%) consisting of >250 thousand base
60 pairs (kbp) sequence and an annotated gap at chromosome 7q35. Studying these SDs is particularly
61 problematic because of the high sequence identity and missing sequences in reference genomes¹³⁻¹⁵.
62 Using large-insert bacterial artificial chromosome (BAC) clones from a haploid hydatidiform mole cell
63 line, we recently assembled a sequence-resolved haplotype spanning over this region and characterized
64 three sets of SDs present in both direct and indirect orientations. Copy number estimations from read-

65 depth analyses suggest that *TCAF* duplications are completely missing in archaic hominins, such as
66 Neanderthal and Denisovan, and nonhuman great apes⁸. In addition, while their sequence analysis
67 suggested that a single full-length *TCAF1* and *TCAF2* exist at the locus, respectively, additional
68 *TCAF1/TCAF2* copies appear to be truncated or incomplete⁸.

69 In the present study, we systematically explore the haplotype structure of the *TCAF* locus in order
70 to study its diversity, annotate the genes, and infer its evolutionary history in the context of selection. We
71 generate 15 sequence-resolved haplotypes from both human and nonhuman primate BAC libraries as well
72 as full-length non-chimeric (FLNC) transcript data from seven tissues using the Pacific Biosciences
73 (PacBio) single-molecule, real-time (SMRT) long-read sequencing technology. We leverage a large
74 collection of over 1,100 publicly available high-coverage Illumina genomes from modern, and archaic
75 humans as well as nonhuman primate samples and document the global diversity of *TCAF* SDs. Finally,
76 we integrate both phylogenetic and population genetic inference approaches to reconstruct the evolution
77 of the *TCAF* SDs in primates and identify putative signals of balancing selection in Native American,
78 Melanesian, and Siberian populations as well as positive selection in both Neanderthal and Denisovan.
79 We provide evidence for differential expressed quantitative trait loci (eQTLs) between two *TCAF* super
80 haplogroups in thyroid, whole blood, tibial nerve, spleen, and lung tissues. In addition, our association
81 analysis using the UK Biobank data suggest that the *TCAF* super haplogroups are likely associated with
82 hypothyroidism, diabetes, and nerve-related traits in some human populations. The study provides one of
83 the most complete genetic investigations of human-specific SDs shedding potential new insights into
84 structural adaptations important in thermal regulation.

85

86 **RESULTS**

87 ***Diversity of TCAF copy number in human and nonhuman primates.*** We applied a read-depth-based
88 genotyper¹⁶ to reevaluate copy number diversity of *TCAF* SDs (duplication segments A, B, and C;
89 **Figure 1A**) in a collection of recently published diverse human and nonhuman samples¹⁷⁻²³. The set
90 includes high-coverage Illumina genomes from a global panel of 1,102 human samples, four archaic

91 hominins, and 71 nonhuman great apes. Consistent with previous studies ⁸, all archaic hominin and
92 nonhuman primate samples are fixed for diploid copy number (CN) 2 (**Figure 1B**). Among contemporary
93 modern human samples, the overall diploid copy number estimates of the three *TCAF* SDs range from
94 two to eight copies and are highly correlated among each other (**Figure S1**), indicating that these three
95 SDs likely appear as a cassette in the evolution of this locus. In contrast to archaic humans where there is
96 a prediction of a single copy of this locus, we estimate that ~98% of modern humans carry more than two
97 copies (**Figures 1B, S1**), suggesting that most of the diversity emerged during the divergence of
98 hominins.

99 Among humans, we observed a wide range of variation across different geographic locations in
100 contemporary human populations (**Figure 1B-C**). While the median copy numbers are consistently CN4
101 (with the exception of Native Americans [median = 3]), African samples on average carry more copies
102 (mean CN = 4.6) than other populations (**Figure 1B**). We tested for copy number stratification by
103 applying both the V_{ST} statistic ²⁴ and CN differentiation test ²⁵. While there is little evidence for population
104 differentiation among pairs of super populations, a reduction in diploid CN among Native American
105 samples distinguishes this group from all others ($V_{ST} > 0.21$, Bonferroni p value for the CN differentiation
106 test < 0.0014 ; **Table S1**).

107 The observation that all diploid CN2 samples (i.e., without duplications) are present solely among
108 non-African populations raises the question as to whether it is due to archaic introgression. To address
109 this, we compared single-nucleotide variation at 20 kbp unique diploid sequences flanking the *TCAF* SD
110 locus (chr7:143,501,000-143,521,000, chr7:143,875,000-143,895,000 [GRCh38]) between archaic and
111 modern human genomes. In general, the archaic flanking haplotypes are similar to those present in
112 modern human samples, including the AFR samples (**Figure S2**). Focusing on sites that are fixed in
113 derived alleles in either the Neanderthal or Denisovan samples, the numbers of derived allele counts in
114 the CN2 samples are comparable to those drawn randomly from African samples ($p = 0.276$; 1,000 non-
115 parametric bootstrap samples). These results suggest that non-African CN2 haplotypes unlikely arose as a
116 result of introgression from Neanderthal or Denisovan.

117

118 *Characterization of TCAF haplotype diversity using long-read BAC-based assemblies.* These near-
119 perfect duplications have hampered not only genome assembly but also the characterization of underlying
120 genetic diversity as well as our understanding of the *TCAF* evolution. An investigation of four recently
121 published long-read whole-genome assemblies, for example, showed that all four failed to create
122 contiguous sequences for both haplotypes spanning over the *TCAF* SDs (**Figure S3, Table S2**). To
123 resolve this, we selected large-insert BAC clones and used PacBio long-read sequencing to generate 15
124 high-quality sequence-resolved haplotypes from eight human samples and three nonhuman primates,
125 including chimpanzee, gorilla, as well as an Old World rhesus macaque monkey (**Table 1, Methods**).

126 A sequence comparison analysis among these newly generated haplotypes confirms the complex
127 organization of these *TCAF* SDs⁸ and further reveals considerable structural diversity in humans when
128 compared to nonhuman primates (**Figure 2 and S4-S12**). While all the nonhuman primate haplotypes
129 carry a single copy for the *TCAF* SD cassette and single full-length *TCAF1* and *TCAF2* genes (**Figures**
130 **S4, S11-S12**), among the 12 human haplotypes, we identify five distinct haplogroups that carry one to
131 three copies for the SD cassettes, which range from 145–406 kbp in length (**Figure 2, Table 1**). The two
132 largest human haplogroups (Haplogroups 4 and 5) both carry three copies of the SD cassette (**Figures 3,**
133 **S5**) but differ by a 100 kbp inversion (Haplogroup5:337,305-437,635; **Figure S5**). The Haplogroup 5
134 organization is mostly consistent with the current human reference genome (GRCh38), which has a gap in
135 the middle of expanded complex SDs in the reference (**Figure 3**). The newly generated Haplogroup 5
136 sequence, thus, eliminates the gap by adding 103,616 bp (Haplogroup5:216,541-320,157) in the human
137 reference.

138 Most of the sequence-resolved haplotypes (8 out of 15) consist of two copies of the *TCAF* SD
139 cassette (Haplogroups 2 and 3; **Figure 2, S6-S9**). These two haplogroups can be structurally distinguished
140 by a similar 100 kbp inversion as identified between Haplogroups 4 and 5. We, therefore, use this
141 inversion polymorphism to classify two super haplogroups (H4 and H5, **Figure 2**). The shortest human
142 BAC haplotypes (Haplogroup 1; **Figures 2 and S10**) that we assembled consist of a single copy of the

143 SD cassette and are similar in structure to those found in nonhuman primates. To further refine the
144 relationship among these assembled haplotypes, we performed a series of sequence alignment analyses
145 among different pairs of haplotypes (**Figure 2**). Based on the organization of the *TCAF* SDs, we
146 conservatively determined that Haplogroups 2 and 5 are closely related and so are Haplogroups 3 and 4.
147 Among Haplogroup 2, two subgroups, Haplogroup 2-1 (n = 1) and Haplogroup 2-2 (n = 5), were further
148 separated and likely emerged due to two independent events, where each involves a different 130 kbp
149 deletion, with respect to Haplogroup 5 (**Figure 2, Table S3**). Similarly, Haplogroup 3 can also be further
150 classified into two subgroups, Haplogroup 3-1 and Haplogroup 3-2, because we predict that they likely
151 derived from two different large deletions with respect to Haplogroup 4 (**Figure 2, Table S3**). Finally,
152 while Haplogroup1 has a similar structure to the nonhuman primate haplotypes, we showed that this
153 group is more closely related to Haplogroup 3 and likely a result of a deletion event over Haplogroup 3
154 (**Figure S10, Table S3**).

155 The breakpoints of these six large structural differences map within nearly identical (>99.5%)
156 SDs, consistent with the action of non-allelic homologous recombination (NAHR) (**Table S3**). By
157 identifying the longest, high-identity sequences around the putative breakpoints (**Table S3**), we show that
158 four associate with DupA SDs, while two map around DupC SDs. For example, the inferred inversion
159 breakpoints of Haplogroup 4 are immediately flanking two inversely oriented, nearly identical sequences
160 within DupC SDs (~39.6 kbp; sequence identity > 99.8%). In addition, our analysis of sequence
161 composition reveals the presence of nearly identical LINE/SINE elements in sequences flanking the
162 putative breakpoints (**Table S3**). The inferred breakpoints for the deletion event between Haplogroups 3-1
163 and 4 overlap with two perfectly identical (100%) 2.5 kbp *CTAGE* (Cutaneous T Cell Lymphoma-
164 Associated Antigen 1) family sequences (**Figure S8, Table S3**). Remarkably, most of the deletion events
165 are similar in size ranging in length from 129.7 to 134.4 kbp. Thus, our data highlight a rapid
166 evolutionary process that involves large-scale duplication, deletion, and inversion events via NAHR but
167 where deletion events are constrained.

168

169 **Discovery of new *TCAF* genes/isoforms.** The newly generated long-read BAC haplotypes reported allow
170 us to revisit *TCAF* annotation. Functional studies have shown that *TCAF1/TCAF2* are highly expressed in
171 the prostate, brain, esophagus, prostate, and skin tissues¹² (GTEx Portal, <https://gtexportal.org>). To this
172 end, we targeted capture of FLNC *TCAF* transcripts and generated 480,700 FLNC transcripts from six
173 human tissues, including dorsal root ganglion, esophagus, fibroblast, skin, fetal brain, and testis, as well
174 as 50,885 FLNC transcripts from a chimpanzee lymphoblast cell line (**Table S4**). *TCAF* transcript on-
175 target rates range from the lowest in the testis sample (2.2%) to the highest for the chimpanzee
176 lymphoblast sample (15.5%). To identify *TCAF* models, we aligned high-quality FLNC transcripts to the
177 assembled haplogroup sequences and only retained transcripts that have complete open reading frames
178 with >100 amino acids (aa) in length (**Methods**). Of note, we followed the previously described
179 nomenclature⁸ to delineate *TCAF* models and isoforms. Overall, our analysis indicates the expression of
180 *TCAF* genes in a wide variety of tissues and reveals new *TCAF* paralogs as well as a variety of new
181 isoforms (**Figures 3, S5-S13; Table S5**). In particular, we recovered only single isoforms for individual
182 genes in the *TCAF1* subfamily. In contrast, the *TCAF2* subfamily shows much more diversity with 10, 2,
183 and 4 isoforms being identified for *TCAF2A*, *TCAF2C1*, and *TCAF2C2*, respectively. With respect to the
184 current annotation in the human reference genome (GRCh38), all haplotypes confirm that *TCAF2A* (or
185 *TCAF2* in RefSeq Release 200) is incorrectly annotated due to the presence of a gap at this locus (**Figure**
186 **3**). The gene models and isoforms reported here represent the most correct and comprehensive
187 annotations for the *TCAF* families (**Table S5**).

188 Compared to the chimpanzee, every additional copy of an SD cassette identified among human
189 haplotypes is associated with an additional copy of *TCAF1A* and *TCAF2C* paralogs, although they are
190 incomplete copies to *TCAF1B* and *TCAF2A*, respectively (**Figures 3, S4-S12**). Specifically, all
191 *TCAF1A1/TCAF1A2* paralogs match to the last seven exons of *TCAF1B*, while the *TCAF2C1/TCAF2C2*
192 paralogs consist of the first three exons of *TCAF2A* (**Figure 3**). In addition, the predicted breakpoints
193 between Haplogroup 2-2 and Haplogroup 5 coincide with the *TCAF1A1* and *TCAF1A2* paralogs. The
194 sequence alignment analysis shows that the breakpoint at Haplogroup 2-2 spans from the third to the fifth

195 introns, and thus, the *TCAF1A1* copy of Haplogroup 2-2 is a fusion version of the *TCAF1A1* and
196 *TCAF1A2* copies of Haplogroup 5 (**Figure S8**). Because the amino acid sequences encoded from these
197 three paralogs are 100% identical, this suggests the actual breakpoints are within one of the three introns.
198 Finally, the 100 kbp inversion between Haplogroup 4 and Haplogroup 5 raises the question whether the
199 large-scale structural variation event has any effect on *TCAF* coding sequences. Notably, the inferred
200 inversion breakpoints map to the second introns of *TCAF2A* and *TCAF2C2*. There are only two
201 differences in the *TCAF2A* exon alignment between Haplogroups 4 and 5: one is a synonymous change in
202 the second exon, and the other causes a nonsynonymous change (R479Q) in the third exon, which is
203 beyond the inferred breakpoints (**Figure S13**). Thus, our results predict that this inversion shuffles the last
204 six exons of *TCAF2A*, which were ligated with the first two exons of different *TCAF2A/TCAF2C2*
205 paralogs on the two haplotypes without disrupting or altering the coding sequence potential of
206 *TCAF2A/TCAF2C2*.

207 It should be noted that although the structural changes largely increase and decrease dosage of
208 individual *TCAF* family members with relatively few predicted amino-acid differences, there is evidence
209 that structural mutations and conversion events are likely affecting the regulatory landscape. We searched
210 for candidate sites of interlocus gene conversion (IGC) among paralogs (**Figure S14**) and observed an
211 IGC site that overlaps with a strong H3K27Ac signal (left white box, **Figure 1A**) an epigenetic signature
212 for enhancers. Using paralog-specific copy number genotypes (**Methods**), we found a reduction in copy
213 number at this locus (i.e., the acceptor site; chr7:143,615,002–143,624,482, GRCh38; the left white box,
214 **Figures 1A and S15**), which, interestingly, is significantly negatively correlated with an increase in copy
215 number at its paralogous locus (the donor site; chr7:143,833,163–143,842,658; the right white box,
216 **Figures 1A and S15**) (Pearson's correlation = -0.3, p value = 5×10^{-16}) in multiple continental populations
217 (**Figure S16**). This result is consistent with the hypothesis of an IGC event copying the donor sequence
218 over the acceptor locus, resulting in reciprocal copy number changes at these two loci. Interestingly, we
219 find that this paralog-specific copy number at the donor locus is negatively correlated ($R = -0.18$, p
220 value = 3.6×10^{-7}) with the (absolute) latitudinal location of human populations while the relationship for

221 the acceptor site is positive ($R = 0.23$, p value = 2.5×10^{-11} ; **Figure S17**). While the observed correlations
222 between latitudes of individual populations and paralog-specific copy numbers at the IGC sites may be
223 the result of natural selection, such strong correlations are relatively common among similar loci across
224 the genome (p value = 0.22; 1,000 random SD loci of genomic shuffling), suggesting demographic
225 history among populations also likely contributing to this observation.

226

227 ***Evolution of TCAF segmental duplications.*** To reconstruct evolutionary origin of the *TCAF* haplotypes
228 in modern humans, we leveraged these high-quality human and nonhuman primate haplotypes for
229 phylogenetic reconstruction (**Tables 1 and S6**). The mean sequence identities for DupA, DupB, and
230 DupC SDs are high among human haplogroups (99.80% [s.d.: 0.04], 99.71% [s.d.: 0.10], and 99.85%
231 [s.d.: 0.05], respectively) (**Figure S18**), while the sequence divergence between human and nonhuman
232 primate *TCAF* haplotypes are compatible to published whole-genome estimates (**Figure S18**)^{26,27}.

233 Together, these findings are consistent with a recent *TCAF* SD expansion or IGC. We considered each of
234 the duplicated segments independently guarding against confounders such as IGC and non-allelic
235 homologous recombination events. Because several SDs represent hybrids of two SD paralogs due to
236 NAHR fusion events (**Figure 2**), for example, we restricted our analysis to the larger segments of hybrid
237 SDs. To minimize the effects of IGC (**Figure S14**)⁸, we applied GENECONV (v1.81a) to exclude sites
238 corresponding to 80.5%, 61.7%, and 70.6% of DupA, DupB, and DupC sequences, respectively (**Table**
239 **S7**). Phylogenetic reconstruction and dating estimates on this filtered set showed that all human copies
240 form a single clade with a 100% posterior support and coalesce at 1.06 (95% C.I.:0.87-1.26), 1.72 (95%
241 C.I.: 1.26-2.21), and 1.47 (95% C.I.: 0.74-2.54) million years ago (Mya), for DupA, DupB, and DupC
242 SDs, respectively (**Figures 4A and S19-S21**). While a few differences in topology were noted among the
243 inferred phylogenies (**Figures S19-S21**), these estimates, with the exception of DupB, overlapped
244 coalescent estimates obtained from a 12.3 kbp single-copy unique region (**Figure 2, the right panel**),
245 which suggests that the two supergroups: H5 (Haplogroup 5, Haplogroup 2-1, and Haplogroup 2-2) and
246 H4 (Haplogroup 4, Haplogroup 3-1, Haplogroup 3-2, and Haplogroup 1) diverged from each other ~0.73

247 Mya (kya; 95% C.I.: 0.46-1.03) (**Figure S22A**). Importantly, these findings suggest that the duplications
248 began to occur before humans and archaic hominins diverged.

249 To further refine our inferences on the evolution of *TCAF*, we realigned genome sequence data
250 from four archaic hominins along with human and chimpanzee samples to a custom, gapless human
251 reference chromosome 7 using GRCh38 and sequence-resolved BAC haplotypes to generate a joint
252 single-nucleotide variant call set (**Methods**). Due to the limitations of short-read data, we restricted this
253 analysis to three unique regions (52.3 kbp in total) flanking and internal to the *TCAF* locus (**Figures 1**
254 **and 2**). We identified 1,275 single-nucleotide variants (SNVs) (QV>20) and used these to
255 computationally phase and to construct haplotypes from 802 samples (**Figures S2 and S23**). Using these
256 archaic hominin haplotypes with the sequences of the seven BAC haplogroups, we estimate the time to
257 the most recent common ancestor (TMRCA) of all hominins being 0.78 Mya (95% C.I.:0.55-1.03 Mya),
258 which is consistent with that reported above. Interestingly, all archaic haplotypes and the supergroup H4
259 are more closely related than either is to the supergroup H5, and the clade of archaic and supergroup H4
260 haplotypes coalesces ~0.53 Mya (**Figure 4B**). Thus, our results show that as early as 1.72 Mya the *TCAF*
261 locus (DupB SD) began to duplicate and that by 0.73 Mya most of the of the *TCAF* haplotype diversity
262 that we observe in modern-day humans had already emerged. Archaic hominins, in contrast, show a
263 reduction in genetic diversity without the associated copy number duplication diversity found in modern
264 humans (**Figure 4C**).

265
266 ***Hominin diversity and natural selection of TCAF haplotypes.*** We used the combined haplotypes
267 constructed from the 53.2 kbp of unique regions to further explore genetic diversity among archaic and
268 present-day hominins. A haplotype-based principal component analysis (PCA) reveals distinct clusters
269 (**Figures S24-S25**) corresponding in part to the two human haplogroups H4 and H5 identified from the
270 BAC sequence analysis (**Figure S25**). Using a supervised clustering method based on haplotype PCA and
271 machine-learning algorithm t-SNE²⁸ for dimensionality reduction and visualization (**Methods**), we show
272 that contemporary hominin haplotypes can be classified into 12 different clusters (**Figures 5A and S26**)

273 with an estimated haplotype heterozygosity of 0.885 (**Methods**), which is much higher than the reported
274 genome-wide estimated mean of 0.534²⁹. Phylogenetic reconstruction using 10 haplotypes drawn at
275 random from each cluster reproducibly reconstructs a maximum likelihood tree closely resembling this
276 cluster-based relationship (**Figure 5B**). Considering the two flanking and unique regions independently
277 also confirms a phylogeny of 12 distinct clusters with different BAC haplogroup representing different
278 clusters (**Figures S26-S27**). Among humans, no single population is restricted to a specific cluster,
279 although some clusters appear enriched. For example, clusters 2 and 7 largely consist of African (60.3%)
280 and South Asian (50.6%) haplotypes, respectively (**Figure S28**). Assuming an equal split of the overall
281 diploid CN estimates between the two haplotypes in an individual, we found significant differences in
282 *TCAF* copy number among haplotype clusters (Bonferroni corrected p values < 0.05; Mann-Whitney *U*
283 tests) (**Figure S29**).

284 In sharp contrast to the extensive haplotype diversity in modern human samples, the eight archaic
285 hominin haplotypes are virtually identical belonging to a single homogenous group (**Figures 5, S24-S29**).
286 Under a panmictic model and assuming an ancestral hominin origin of ~0.73 Mya for *TCAF* haplotypes, it
287 is highly unlikely that all eight archaic haplotypes would form a single cluster given modern-day human
288 diversity (p value = 2.1×10^{-7} , 100×10^6 permutations). To formally test whether natural selection plays a
289 role in these unusual patterns in both modern and archaic humans, we performed a selection test using the
290 Tajima's *D* statistic and coalescent simulations under plausible demographic models to assess
291 significance (**Methods**). The results suggest significant distinct signals of selection in humans and archaic
292 hominins separately (**Figure 5C**). Among modern human populations, Native American, Melanesian, and
293 Siberian populations show evidence for excess of heterozygosity and significantly larger Tajima's *D*s
294 compared to those from coalescent simulations under neutral demographic models (**Figure 5C**); though
295 the distributions for Melanesians and Siberians center on zero. Especially among Native Americans
296 (Tajima's *D* = 1.24, p value = 4.3×10^{-9}), these results are consistent with the action of balancing selection
297 at this locus (**Figure 5C**). In contrast, the Tajima's *D* distribution of the archaic hominin samples are
298 strictly negative and significantly lower than the neutral expectation, and this result holds even when the

299 Denisovan sample was removed (Tajima's $D = -1.05$; p value $< 1 \times 10^{-6}$) (**Figure 5C**). Our results suggest
300 that the observed archaic hominin haplotype variation pattern is unlikely explained by their demographic
301 history, thus, favoring a hypothesis of positive selection.

302 These dramatic differences in diversity and potentially distinct selection forces pose the question
303 of whether there are any functional differences among the different *TCAF* haplogroups. Because our
304 structural analyses suggested dosage and regulatory differences among common haplogroups (**Figures**
305 **S28-S29**), we investigated *TCAF2* expression using the GTEx multi-tissue data (release v8) and possible
306 genetic associations between traits and *TCAF* haplotypes using the UK Biobank data. Based on the BAC
307 references, we identified four tagging SNVs in the unique diploid region embedded among the *TCAF* SDs
308 that could distinguish super haplogroups H4 and H5 (**Figure 6**). Because the two super haplogroups
309 structurally differ by the 100 kbp inversion, we expect a reduction of homologous recombination events
310 between the two groups and therefore an increase in linkage disequilibrium (LD) in hominins. As
311 expected, we observe substantial LD across all three unique diploid regions in all archaic and modern
312 samples and confirm nearly complete LD among the four tagging SNVs ($D' > 0.98$, $R^2 > 0.89$) (**Figure**
313 **S30**). Using these data, we again found that all archaic samples are fixed for the H4 supergroup and
314 inferred the frequencies of H4, H5, and other forms in modern humans being 57.1%, 40.3%, and 2.6%,
315 respectively. While the H4 haplogroup is found at high frequency among non-Africans, especially in
316 Melanesians (72%), some Native Americans (e.g., the Karitiana, >72%), and Northeast Asians (e.g., the
317 Mongolians, >68%), the H5 group segregates in higher frequency among Europeans (e.g., the Basques,
318 >81%), Middle Easterners (e.g., the Druze, >61%), and some Native Americans (e.g., the Surui, >83%)
319 (**Figure 6A**). Consistent with our haplotype-based PCA clustering results, we observed high frequencies
320 of the H4/H5 heterozygous form (orange in **Figure 6B**) across populations, with the highest found in
321 Southeast Asians and Native Americans (>60%) and the lowest in some Melanesian and African
322 populations (<20%).

323 Given the common presence of these two super haplogroups in humans (**Figure 6A-6B**), we
324 collected multi-tissue eQTL data from the GTEx Project (release v8) and found that all four tagging

325 SNVs are *TCAF2* eQTLs (**Figure 6C**). The H4 haplogroup (reference) alleles are significantly associated
326 with increasing *TCAF2* (p values $< 1 \times 10^{-10}$) expression in thyroid and tibial nerve tissues in addition to
327 others (**Figure 6C**). Thus, our inferences of the differential *TCAF2* expression between the H4 and H5
328 haplogroups imply possibly functional differences in the peripheral nervous system and thyroid among
329 hominin groups. Finally, we explored traits that are potentially associated with the *TCAF* haplogroups
330 using the recently released genetic association data from the ancestry-based genetic analysis of the UK
331 Biobank (Pan-UK Biobank as of 2020.08.06; <https://pan.ukbb.broadinstitute.org/>). We examined SNVs
332 from the three unique diploid regions at the *TCAF* locus and focused on nerve- and thyroid-related traits,
333 including diabetes as thyroid dysfunctions often result in diabetes mellitus³⁰. While in general no SNVs
334 reach genome-wide significance for traits that we interrogated, we found several suggestive association
335 signals among specific human populations. For example, we found nominal associations (uncorrected p
336 values < 0.05) for SNVs, including the four eQTLs, for hypothyroidism (phecode-244) in both African
337 and East Asian populations as well as for diabetes medication (Glimepiride) and compressed nerves in
338 samples of European descent (**Figure S31**) consistent with the expression pattern of these genes.

339

340 DISCUSSION

341 Genetic changes in SDs are now well recognized as an important source of functional innovation in
342 human evolution due, in part, to subsequent rounds of duplication, deletion, or inversion^{13,14}. Such
343 regions, however, are often frequently overlooked as part of comparative studies, association studies, or
344 scans of selection because of their structural complexity and high-degree of homology. In this study, we
345 first delineate the complex structural diversity of the locus in humans and compare it to nonhuman
346 primates at the sequence level in order to reconstruct its evolutionary history¹². Of note, complete
347 assembly of 15 haplotypes by tedious BAC-based sequencing was sufficient to capture more than half of
348 the genetic diversity of *TCAF* SDs in humans. Based on these data, we develop a model of rapid
349 diversification with respect to both haplotype and *TCAF* SD copy number diversity (**Figures 1-3**). It is
350 worth noting that despite the recent advances in long-read sequence assembly, the most complex

351 haplotypes still remain collapsed and unassembled in some of the recently released human genomes³¹⁻³³
352 (**Figure S4, Table S2**). We believe that this is due partly to the read-length limitation of the data as well
353 as the near-perfect sequence identity among *TCAF* SDs that lead to ambiguous read partitions during
354 haplotype assembly. This highlights the importance of alternate approaches, such as large-insert BAC
355 clones, to construct highly accurate haplotype-resolved assemblies. More sophisticated assembly
356 approaches and long-read technology for generating longer and even more accurate sequence reads will
357 likely make such investigations much more routine in the future³⁴.

358 The structural complexity and copy number variation of this region underscores its mutational
359 dynamics. The high degree of sequence identity among *TCAF* paralogs (>99.7%), their clustered and
360 tandem orientation, and their large size (10–60 kbp) promote recurrent structural rearrangements making
361 the entire locus genetically unstable as a result of NAHR¹³. Human Haplogroups 4 and 5, in particular,
362 carry more *TCAF* SD paralogs than others, and, thus, both are predicted to be subject to a higher
363 likelihood of structural mutation. Our analysis predicts that H4 would be more predisposed to structural
364 rearrangements than H5 due to a greater number of directly orientated duplications (DupA and DupB)
365 (**Figure 2**). This prediction is consistent with the observation of more subgroups within supergroup H4 in
366 the haplotype-based PCA (**Figures S24-S27**). Surprisingly, this predisposition to recurrent structural
367 diversity is specific to modern-day humans as our sampling of eight archaic hominin haplotypes suggests
368 a static H4 haplogroup without a single duplication event. Given our estimates of origin of the
369 duplications (~1.7 Mya) and haplogroup coalescence (~0.73 Mya), the lack of diversity in the
370 Neanderthal and Denisovan lineages is highly unlikely.

371 One of the challenges studying this region is the wide-spread IGC we identified among the *TCAF*
372 SDs. IGC essentially erases signatures left by previous mutation events in the acceptor site, thus biasing
373 towards younger TMRCA estimates^{35,36}. While we mitigated such a bias by focusing our analysis on non-
374 IGC sequences, this reduced the available dataset limiting our power to infer its evolution at a finer level
375 of resolution. Interestingly, the analysis uncovered a significant correlation between paralog-specific copy
376 number and the latitudinal location of human populations precisely over a potential enhancer region,

377 suggesting that IGC may have contributed to regulatory differences of *TCAF* genes among human
378 populations. Such co-option of IGC has been noted before among duplicate genes. For example, a human-
379 specific, pseudogene-mediated IGC of *SIGLEC11* replaced the 5' upstream region and coding exons of
380 *SIGLEC11* by its adjacent pseudogene *SIGLEC16P* in the human lineage. This event has been
381 hypothesized as adaptive during human brain development³⁷. While we cannot rule out the possibility of
382 confounding effects of demography resulting in the observed correlation, given the importance of *TCAF*
383 in thermal regulation and sensing, it is tempting to hypothesize that IGC contributes to local adaptation by
384 changing dosages through rapid replacement of regulatory sequences in addition to SD.

385 Our more complete reference genomes coupled with targeted cDNA sequencing also provides the
386 most comprehensive *TCAF* annotation to date for humans and nonhuman primates (**Figures 3 and S4-**
387 **S12**). We show that each additional *TCAF* SD cassette (DupA, DupB, and DupC) generates additional
388 copies of *TCAF1A* and *TCAF2C* through incomplete gene duplication. Gene duplication has been known
389 to be an important source of genetic novelty and adaptive evolution^{1,38}, especially for the great ape
390 lineages given the surge of duplications predicted to have occurred in the hominid ancestor³⁹. Throughout
391 the evolution of mammals, *TCAF* genes have been highly conserved, and human is the only lineage that is
392 known to carry additional copies. The coding sequences among the individual paralogs are highly
393 conserved as a result of their recent origin arguing for the relative importance of *TCAF* dosage differences
394 in regulating the TRPM8 ion channel as suggested by a previous studies¹². We also observe, however,
395 considerable isoform diversity as a result of the incomplete nature of the SD particularly among *TCAF2*
396 paralogs suggesting a potentially more complex relationship with the ancestral genes. This pattern of
397 incomplete duplication resembles the evolution of the neural human-specific duplication *SRGAP2* genes
398⁵, where the younger, incomplete duplicated paralog *SRGAP2C* inhibits the function of the ancestral copy
399 *SRGAP2A*, and in mice experiments this in turn prolongs the maturation of synapses and subsequently
400 increases the density of the length of synapses⁴. Taken together, our results bolster the hypothesis of
401 recent incomplete gene duplications and adaptive human evolution.

402 Remarkably, our data support a model of two distinct forces of natural selection operating on the
403 same locus over the last half million years of hominin evolution. We propose that diversifying or
404 balancing selection is acting in at least some human populations, particularly out-of-African populations
405 such as Native Americans, to maintain and expand haplotype and structural diversity. While the exact
406 mechanism and phenotype that balancing selection are acting upon are not clear, the differential eQTL
407 signal related to *TCAF2* expression between supergroup H4 and H5 alleles in multiple tissues would be
408 consistent with antagonistic pleiotropy, where two haplotypes have opposing effects on two different
409 traits depending on either environmental or life-history conditions^{40,41}. In contrast, Neanderthal and
410 Denisovan show a paucity of genetic variation, and while the sample size is still limited, this observation
411 is unlikely to change with the sequencing of additional archaic genomes. Moreover, both the significantly
412 negative Tajima's *D* statistic and the single-haplotype without duplication are highly unlikely observed
413 under neutrality after accounting for differences explained by demographic history. Our inference results,
414 thus, strongly argue that positive selection has reduced genetic diversity at the *TCAF* locus in these
415 archaic hominin lineages.

416 The opposing effects on the *TCAF2* expression found between H4 and H5 haplogroups in the
417 thyroid is particularly intriguing. Thyroid activity is one of the most important determinants of overall
418 energy expenditure and basal metabolic rate in human body⁴². Not only does the thyroid hormone
419 maintain core body temperature when exposed to cold, but it also regulates metabolism through action in
420 the brain, brown fat, skeletal muscle, etc.⁴³. Thyroid dysfunctions, both hyperthyroidism and
421 hypothyroidism, often lead to metabolic disorders, such as diabetes mellitus³⁰, and also cause neurological
422 disorders during fetal and brain development due to low levels of maternal thyroid hormones from
423 maternal hypothyroidism or lack of iodine in diet^{44,45}. Despite no genome-wide significant genetic
424 associations, we did identify suggestive trends between specifically those eQTL variants and
425 hypothyroidism, diabetes medication, as well as pinched nerves. Interestingly, the TRPM8 ion channel, in
426 addition to its role in detecting environmental cold, is also a known regulator of insulin homeostasis
427 through neuronal control of liver insulin clearance and, thus, required for precise thermoregulatory

428 responses to cold and fasting^{46,47}. In addition, a rare variant within the *TCAF* locus (chr7:143735931) was
429 recently reported to be associated with type 2 diabetes in sub-Saharan African populations⁴⁸. Given this
430 variant's restricted geographic distribution and low frequency in the GTEx data (allele frequency: ~0.4%;
431 GTEx portal as of 2020.07.21), it explains neither the observed patterns of diversity nor the selection and
432 the eQTL signals observed here. Therefore, if the *TCAF* haplotype and structural diversity affect an
433 individual's ability to regulate TRPM8 ion-channel gating and, thus, body temperature, it is possible that
434 the antagonistic forces of selection promote adaptations to a cold environment or an enhanced ability to
435 dynamically adapt to changing conditions, such as diets.

436 While it is still controversial, cold weather is thought to have forced Neanderthals to move from
437 coastal areas to further inland, where they would have had to adapt to iodine-poor diets, which may in
438 turn have affected brain development in this species^{44,49}. Iodine deficiency is one of the common causes
439 of hypothyroidism (reduced thyroid hormone levels), which is characterized by lower core body
440 temperature and decreased resting energy expenditure³⁰. Despite the limited understanding of the
441 molecular bases underlying the interactions between *TCAF* and thyroid hormones, one possibility may be
442 that selection on H4 increased *TCAF2* expression and subsequently suppressed the activity of TRPM8 in
443 the thyroid of archaic hominins. This may have improved physiologically adaptation to cold weather and
444 low iodine-intake diets.

445

446 **METHODS**

447 Copy number genotyping in high-coverage Illumina short-read data

448 To explore the temporal and spatial copy number variation of *TCAF* SDs, we leveraged a large collection
449 of publicly available high-coverage Illumina genomes from 1,102 modern human samples, four archaic
450 hominins, and 71 nonhuman great apes¹⁷⁻²³. In short, sequencing reads were divided into multiples of 36-
451 mer and then mapped to a repeat-masked human reference genome (GRCh38) using mrsFAST (v3.4.1)⁵⁰.
452 We allowed up to two mismatches per alignment to increase our mapping sensitivity and corrected for
453 possible biases in read depth due to different levels of GC content in sequencing reads. Finally, the copy

454 number estimate was computed by summarizing over all mappable bases at the locus of interest for each
455 sample.

456

457 Library processing and assembly for *TCAF* BAC haplotypes

458 We obtained the human BAC libraries used in this study from the Virginia Mason Research Center and
459 the nonhuman primate ones from the Children’s Hospital Oakland Research Institute. Probes for regions
460 of interest were designed, radioactively labeled and hybridized to the Performa filters, washed, exposed to
461 Phosphor screens, and scanned on a Typhoon scanner. Positives are called and corresponding clones
462 selected from the BAC library. We prepared barcoded libraries from clone DNA using Illumina-
463 compatible Nextera DNA sample prep kits (Epicentre, catalog number GA09115) as described
464 previously⁵¹ and carried out paired-end sequencing (125 bp reads) on an Illumina HiSeq 2500. Reads
465 were then mapped to the reference genome, GRCh38, to identify singly unique nucleotide k-mers
466 (SUNKs)¹⁶. Non-overlapping BACs were pooled and sheared as described previously⁵². Libraries
467 were processed using the PacBio SMRTbell Template Prep kit following the protocol “Procedure and
468 Checklist—20 kb Template Preparation Using BluePippin Size-Selection System” with the addition of
469 barcoded adaptors during ligation. Up to ten barcoded libraries were then pooled at equimolar
470 amounts and size-selected as a pool on the Sage PippinHT with a start value of 10,000–12,000 and an
471 end value of 50,000. The resulting library was then sequenced on one Sequel SMRT cell 1M by
472 diffusion using Sequel v3.0 chemistry. We performed *de novo* assembly of pooled BAC inserts using
473 Canu (v1.5). Reads were masked for vector sequence (pCC1BAC) and assembled with Canu, then
474 subjected to consensus sequence calling with Arrow
475 (<https://github.com/PacificBiosciences/GenomicConsensus>). We reviewed PacBio assemblies for
476 misassembly by visualizing the read depth of PacBio reads in Parasight (v7.6,
477 <http://baileylab.brown.edu/parasight/download.html>), using coverage summaries generated during the
478 resequencing protocol. We performed BLAST (v2.10.1) pairwise alignment between the BAC sequences
479 for the construction of haplotypes. Complete contigs for 15 *TCAF* haplotypes must satisfy the following

480 two conditions: (1) overlapped DNA segment between BAC sequences must be at least 99.99% identical
481 with no tolerated mismatch and (2) contig must anchor to unique location on both ends when mapping
482 back to the human genome reference GRCh38.

483

484 Identification of haplogroups and interlocus gene conversion (IGC) tracts

485 All sequence alignments in this study were performed using MAFFT (v7.453). To determine the
486 structural differences among the haplogroups, we built single-base level sequence identity profiles using
487 pairwise sequence alignment, followed by smoothing over windows of 500 bp with a step size of 100 bp.
488 Because we are interested in recent, large IGC events, empirical IGC tracts are determined as more than
489 two consecutive windows that have 100% sequence identity. To find additional support, we use a model-
490 based program, GENECONV (v1.81a), to identify pairs of uninterrupted sequences with 100% sequence
491 identity that are longer than expected given the overall pattern of variable sites in an alignment. To be
492 conservative, we considered IGC tracts that have both simulation- and Karlin-Altschul- p values < 0.05
493 reported by GENECONV.

494

495 Full-length non-chimeric (FLNC) transcripts for *TCAF* genes using long-read cDNA sequencing

496 Total RNA was harvested from six human tissues (dorsal root ganglion: Clontech Laboratories, Inc.,
497 Takara [catalog #636150]; esophagus: Clontech Laboratories, Inc., Takara [catalog #636178]; fibroblast
498 (Coriell); skin: Biochain Institute Inc. [catalog #R1234218-P], fetal brain: NCBI BioSample:
499 SAMN09459150; testis: Clontech Laboratories, Inc., Takara [catalog #636533]) and one chimpanzee
500 lymphoblast cell line. We purified polyA RNA using oligo-dT magnetic beads (Dynal: Thermo Fisher
501 Scientific Inc.). Double-stranded cDNA with 96 random barcodes was prepared, amplified and subjected
502 to hybridization capture followed the protocols detailed in (Dougherty et al. 2018). Hybridization probes
503 were designed to target the duplicated exons of *TCAF2*, and sequences are given in Supplemental Table
504 S2b of Dougherty et al., 2018.

505 Following post-capture PCR, the amplified dsDNA was purified on magnetic beads (AMPure PB;
506 PacBio) and then subjected to library preparation for long-read sequencing (SMRTbell Template Prep Kit
507 1.0; PacBio with barcoded SMRTbell adapters). SMRT sequencing was performed on the Sequel
508 platform with Sequel v3.0 chemistry (PacBio) with LR SMRT Cells with 2-hour pre-extension and 20-
509 hour movies. Reads corresponding to each sample were extracted by their SMRTbell barcodes and
510 circular consensus sequences were generated from the raw subreads using SMRT Link with minimum
511 number of pass set to 1. The program lima in the Iso-Seq3 pipeline
512 (<https://github.com/PacificBiosciences/IsoSeq3>) was applied to remove the 5' and 3' dual barcodes and
513 also obtain the unclustered FLNC reads. Parameters used were: lima --isoseq --dump-clips. We did not
514 cluster the FLNC reads further because highly identical paralogous transcripts could undesirably cluster
515 together in this step. Due to the variability of the yields of FLNC transcripts across samples and loci, we
516 used data from various combinations of samples in subsequent analyses when applicable.

517 To determine the coding potential for each FLNC transcript, we use the program ANGEL
518 (<https://github.com/PacificBiosciences/ANGEL>) to call open reading frames with a minimum length
519 cutoff of 200 amino acids and proper stop codons. All transcripts with open reading frames were mapped
520 to the seven *TCAF* haplogroups (**Tables S4 and S5**) using minimap2 (v2.17-r941) and individual best
521 placements to specific haplogroups were determined using sequence identity. Given the recent history of
522 the *TCAF* duplications, we required a >99% sequence identity for an alignment to be considered a match
523 between an FLNC read and the haplogroup contig and most have at least an overlap of 200 bases. To
524 improve the sensitivity of FLNC read assignments, we allowed multi-mapping placements for an FLNC
525 read as long as the difference in sequence identity between alignments is < 0.01%. Finally, to determine
526 gene models and isoforms in each haplogroup, we focused our analysis on FLNC transcripts with more
527 than 10 reads (**Table S5**).

528

529 Single-nucleotide variant (SNV) calling for unique diploid sequences at *TCAF* locus

530 To perform population genetic analyses, we generated a joint call set for 840 published high-coverage
531 short-read Illumina genomes, including 828 Human Genome Diversity Project (HGDP) samples, eight
532 chimpanzee samples, and the four high-coverage archaic genomes. Because the current human reference
533 GRCh38 has a gap at the *TCAF* locus, we constructed two custom references for chromosome 7: one with
534 Haplogroup 4 and the other with Haplogroup 5, in addition to a *TCAF*-SDs-hardmasked
535 (chr7:143521769-143874696, GRCh38) chromosome 7. Paired-end data were aligned to the two custom
536 references using BWA-MEM (v0.7.12), while ancient DNA data were mapped using BWA (v0.7.17) with
537 parameters with parameters: `-n 0.01 -o 2 -l 16500`¹⁹. SNV genotypes were jointly generated using
538 haplotype caller FreeBayes (v1.0.2-6-g3ce827d) with the following parameters: “`--min-coverage 10 --`
539 `use-best-n-alleles 4`”. To ensure genotype quality, we excluded variants that were found within 10 bp of
540 indels and have quality score (QUAL) < 20. We identified a total of 1,275 and 1,295 SNVs from the call
541 sets of Haplogroups 4 and 5, respectively, among the 840 samples. Note that we excluded 38 and 30
542 samples from the call sets of Haplogroups 4 and 5, respectively, due to >10% missing data to ensure the
543 quality of genotype calls and downstream analyses. Because the results and interpretations of all
544 downstream analyses are highly compatible, we concluded no obvious reference biases between the two
545 sets. Unless otherwise stated, we reported results from the Haplogroup 4 call set because it has slightly
546 better mapping quality.

547

548 Phylogenetic analyses

549 To infer the phylogenetic relationships for *TCAF* SDs and the unique diploid sequences in primates, we
550 performed both maximum likelihood (IQ-TREE, v1.6.11)⁵³ and Bayesian phylogenetic-based (BEAST
551 v2.5.0)⁵⁴ analyses. First, we determined parameters and models using ModelGenerator (v0.85)⁵⁵ and
552 compared with the recommendations from the model-selection feature of IQ-TREE (using the parameters:
553 “`-bb 1000 -redo -nt 1 -m MFPMERGE`”). Based on the best-fit model recommended, to run BEAST, in
554 general we set 1) GAMMA Category Count = 5, shape parameter = 0.1, Proportion Invariant = 0.38, and

555 using the GTR substitution model (all rates = 1.0) for the Site Model, and 2) Relaxed Clock Log Normal
556 over branches on the tree for Clock Model. For tree priors, we used Calibrated Yule Model, and while
557 kept most of the parameters of the priors as default, for birthRate and clockRate we used Gamma(0.001,
558 1000) and set calibrations based on human–chimpanzee and human–rhesus macaque divergence
559 following the distributions of a log-normal(M = 6500000, S = 0.09) and log-normal(M = 10000000,
560 S = 0.09), respectively. For the analysis involving the inferred archaic hominin haplotypes, we also set
561 dates for terminal samples to account for the differences between contemporary and fossil samples. For
562 each run, we drew a date (in years) from a uniform [36000, 100000] for the Neanderthals and a uniform
563 [50000, 80000]. For each locus, we performed five independent runs to infer the phylogeny using a chain
564 length of 50,000,000 samples and recorded every 1,000 samples. We used the accompany program Tracer
565 (v.1.7.1) to determine the quality of each run and, in general, we used the first 10% as burn-in and only
566 kept runs with ESS > 500. All phylogenetic trees were plotted using Figtree (v1.4.3).

567

568 Population genetic analyses

569 All population genetic analyses in this study were based on SNVs from the three unique diploid regions at
570 the *TCAF* locus (chr7:143,501,000-143,521,000, chr7:143,729,525-143,741,525, chr7:143,875,000-
571 143,895,000 [GRCh38]) using the HGDP, archaic hominin, and chimpanzee samples. Haplotypes were
572 inferred by applying BEAGLE (v5.1-25Nov19.28d) with default settings to modern human, archaic
573 hominin, and chimpanzee samples separately. We explored the diversity of the *TCAF* haplotype using a
574 haplotype-based PCA to classify and assign haplotypes into individual clusters. Cluster assignments were
575 based on the supervised k-means algorithm using the top N informative PCs, where $\sum_{i=1}^N Var(PC_i) >$
576 0.9. To determine the best number of cluster k, for each possible k between 1 and 20, we computed the
577 within groups sum of squares (WSS) using a distance matrix among the top N PCs and determined the
578 best k when $WSS_{k+1} - WSS_k < \delta$, where delta was chosen as 500 according to our analysis. Data
579 visualization of the haplotype clustering was performed using the machine learning technique t-
580 distributed stochastic neighbor embedding (t-SNE)²⁸ implemented in the R package Rtsne

581 (<https://github.com/jkrijthe/Rtsne>) using the following parameters: “perplexity=50, max_iter=5000,
582 early_exag_coeff=12, exaggeration=4, stop_lying_iter=1000, check_duplicates = FALSE, dims=3”. We
583 optimized the t-SNE analysis by running five replicates of this computation and identifying the smallest
584 Kullback-Leibler distance as the candidate for final visualization ²⁸.

585 To test for signals of natural selection, we computed the Tajima’s *D* statistic ⁵⁶ over 2000 bp
586 windows over the 52.3 kbp unique diploid sequences for individual populations as well as the archaic
587 samples. Because the Tajima’s *D* statistic is sensitive to population history, we performed a large
588 coalescent simulation using MaCS ⁵⁷ and 1,000 demographic models as described previously ^{25,58} to
589 account for past histories between modern and archaic hominins and among multiple African and non-
590 African populations. Note that our simulations carefully match the local mutation rate and recombination
591 rate variation at these regions to avoid possible biases to our selection inferences. To estimate the
592 TMRCA for a locus of interest, we used the Thomson’s estimator ⁵⁹. LD measurements among SNVs
593 were computed using Lewontin’s D' ⁶⁰ and R^2 implemented in PLINK (v1.09, www.cog-
594 genomics.org/plink/1.9/) ⁶¹. Note only SNVs with minor allele frequencies > 10% were included for this
595 analysis because lower frequency variants are less informative about linkage.

596

597 ACKNOWLEDGEMENTS

598 The authors thank T. Brown for assistance in editing this manuscript. We also thank S.C. Murali and D.S.
599 Gordon for help submitting data to the NCBI database. **Funding:** This work was supported, in part, by
600 the US National Institutes of Health (NIH) grant R01HG002385 to E.E.E. P.H. is supported by the NIH
601 Pathway to Independence Award (NHGRI, K99HG011041). S.C. was supported by a National Health and
602 Medical Research Council (NHMRC) C. J. Martin Biomedical Fellowship (1073726). E.E.E. is an
603 investigator of the Howard Hughes Medical Institute. **Author contributions:** P.H. and E.E.E. designed
604 and planned experiments. V.D., C.B., S.C., A.P.L., K.M.M., M.S., A.E.W., and J.G.U. prepared libraries
605 and generated and analyzed sequencing data. P.H., V.D., M.R.V., and T.H. performed variant calling and
606 bioinformatics analyses. P.H., V.D., M.R.V., and P.C.D. analyzed long-read sequencing data and

607 assembled contigs. P.H. and Y.M. performed population genetic and phylogenetic inferences. A.P.L.,
608 K.M.M., P.C.D., and J.G.U. generated Iso-Seq transcript data. P.H. and E.E.E. wrote the manuscript.

609 **Declaration of interests:** J.G.U. is an employee of Pacific Biosciences, Inc.

610

611 **DATA AVAILABILITY:** All data used in this study, including assembled BAC contigs and Iso-Seq

612 capture transcript data, are deposited in NCBI under the BioProject: PRJNA657884, BioSample:

613 SAMN09459150; these data are available to anyone for purposes of reproducing or extending the

614 analysis.

615

616

617 REFERENCES

- 618 1. Ohno, S. *Evolution by gene duplication*, xv, 160 p. (Springer-Verlag, New York, 1970).
- 619 2. Kimura, M. & Ohta, T. On some principles governing molecular evolution. *Proc Natl*
620 *Acad Sci U S A* **71**, 2848-52 (1974).
- 621 3. Boyd, J.L. *et al.* Human-chimpanzee differences in a FZD8 enhancer alter cell-cycle
622 dynamics in the developing neocortex. *Curr Biol* **25**, 772-779 (2015).
- 623 4. Charrier, C. *et al.* Inhibition of SRGAP2 function by its human-specific paralogs induces
624 neoteny during spine maturation. *Cell* **149**, 923-35 (2012).
- 625 5. Dennis, M.Y. *et al.* Evolution of human-specific neural SRGAP2 genes by incomplete
626 segmental duplication. *Cell* **149**, 912-22 (2012).
- 627 6. Florio, M. *et al.* Human-specific gene ARHGAP11B promotes basal progenitor
628 amplification and neocortex expansion. *Science* **347**, 1465-70 (2015).
- 629 7. Fiddes, I.T. *et al.* Human-Specific NOTCH2NL Genes Affect Notch Signaling and
630 Cortical Neurogenesis. *Cell* **173**, 1356-1369 e22 (2018).
- 631 8. Dennis, M.Y. & Eichler, E.E. Human adaptation and evolution by segmental duplication.
632 *Curr Opin Genet Dev* **41**, 44-52 (2016).
- 633 9. Nuttle, X. *et al.* Emergence of a Homo sapiens-specific gene family and chromosome
634 16p11.2 CNV susceptibility. *Nature* **536**, 205-9 (2016).
- 635 10. Colburn, R.W. *et al.* Attenuated cold sensitivity in TRPM8 null mice. *Neuron* **54**, 379-86
636 (2007).
- 637 11. Key, F.M. *et al.* Human local adaptation of the TRPM8 cold receptor along a latitudinal
638 cline. *PLoS Genet* **14**, e1007298 (2018).
- 639 12. Gkika, D. *et al.* TRP channel-associated factors are a novel protein family that regulates
640 TRPM8 trafficking and activity. *J Cell Biol* **208**, 89-107 (2015).
- 641 13. Eichler, E.E. Recent duplication, domain accretion and the dynamic mutation of the
642 human genome. *Trends Genet* **17**, 661-9 (2001).
- 643 14. Chaisson, M.J., Wilson, R.K. & Eichler, E.E. Genetic variation and the de novo assembly
644 of human genomes. *Nat Rev Genet* **16**, 627-40 (2015).

- 645 15. Chaisson, M.J.P. *et al.* Multi-platform discovery of haplotype-resolved structural
646 variation in human genomes. *Nat Commun* **10**, 1784 (2019).
- 647 16. Sudmant, P.H. *et al.* Diversity of human copy number variation and multicopy genes.
648 *Science* **330**, 641-6 (2010).
- 649 17. Meyer, M. *et al.* A high-coverage genome sequence from an archaic Denisovan
650 individual. *Science* **338**, 222-6 (2012).
- 651 18. Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature* **499**,
652 471-5 (2013).
- 653 19. Prufer, K. *et al.* A high-coverage Neandertal genome from Vindija Cave in Croatia.
654 *Science* **358**, 655-658 (2017).
- 655 20. Prufer, K. *et al.* The complete genome sequence of a Neanderthal from the Altai
656 Mountains. *Nature* **505**, 43-9 (2014).
- 657 21. Mafessoni, F. *et al.* A high-coverage Neandertal genome from Chagyrskaya Cave. *Proc*
658 *Natl Acad Sci U S A* **117**, 15132-15136 (2020).
- 659 22. Bergstrom, A. *et al.* Insights into human genetic variation and population history from
660 929 diverse genomes. *Science* **367**(2020).
- 661 23. Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142 diverse
662 populations. *Nature* **538**, 201-206 (2016).
- 663 24. Redon, R. *et al.* Global variation in copy number in the human genome. *Nature* **444**, 444-
664 54 (2006).
- 665 25. Hsieh, P. *et al.* Adaptive archaic introgression of copy number variants and the discovery
666 of previously unknown human genes. *Science* **366**(2019).
- 667 26. Magness, C.L. *et al.* Analysis of the *Macaca mulatta* transcriptome and the sequence
668 divergence between *Macaca* and human. *Genome Biol* **6**, R60 (2005).
- 669 27. Kronenberg, Z.N. *et al.* High-resolution comparative analysis of great ape genomes.
670 *Science* **360**(2018).
- 671 28. Maaten, L.v.d. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning*
672 *research* **9**, 2579-2605 (2008).
- 673 29. Stephens, J.C. *et al.* Haplotype variation and linkage disequilibrium in 313 human genes.
674 *Science* **293**, 489-93 (2001).
- 675 30. Biondi, B., Kahaly, G.J. & Robertson, R.P. Thyroid Dysfunction and Diabetes Mellitus:
676 Two Closely Associated Disorders. *Endocr Rev* **40**, 789-824 (2019).
- 677 31. Nurk, S. *et al.* HiCanu: accurate assembly of segmental duplications, satellites, and allelic
678 variants from high-fidelity long reads. *bioRxiv*, 2020.03.14.992248 (2020).
- 679 32. Garg, S. *et al.* Accurate chromosome-scale haplotype-resolved assembly of human
680 genomes. *bioRxiv*, 810341 (2020).
- 681 33. Porubsky, D. *et al.* A fully phased accurate assembly of an individual human genome.
682 *bioRxiv*, 855049 (2019).
- 683 34. Vollger, M.R. *et al.* Improved assembly and variant detection of a haploid human
684 genome using single-molecule, high-fidelity long reads. *Ann Hum Genet* **84**, 125-140
685 (2020).
- 686 35. Chen, J.M., Cooper, D.N., Chuzhanova, N., Ferec, C. & Patrinos, G.P. Gene conversion:
687 mechanisms, evolution and human disease. *Nat Rev Genet* **8**, 762-75 (2007).
- 688 36. Osada, N. & Innan, H. Duplication and gene conversion in the *Drosophila melanogaster*
689 genome. *PLoS Genet* **4**, e1000305 (2008).
- 690 37. Hayakawa, T. *et al.* A human-specific gene in microglia. *Science* **309**, 1693 (2005).

- 691 38. Lynch, M. & Katju, V. The altered evolutionary trajectories of gene duplicates. *Trends*
692 *Genet* **20**, 544-9 (2004).
- 693 39. Marques-Bonet, T., Girirajan, S. & Eichler, E.E. The origins and impact of primate
694 segmental duplications. *Trends Genet* **25**, 443-54 (2009).
- 695 40. Faria, R., Johannesson, K., Butlin, R.K. & Westram, A.M. Evolving Inversions. *Trends*
696 *Ecol Evol* **34**, 239-248 (2019).
- 697 41. Merot, C., Llaurens, V., Normandeau, E., Bernatchez, L. & Wellenreuther, M. Balancing
698 selection via life-history trade-offs maintains an inversion polymorphism in a seaweed
699 fly. *Nat Commun* **11**, 670 (2020).
- 700 42. Kim, B. Thyroid hormone as a determinant of energy expenditure and the basal metabolic
701 rate. *Thyroid* **18**, 141-4 (2008).
- 702 43. Mullur, R., Liu, Y.Y. & Brent, G.A. Thyroid hormone regulation of metabolism. *Physiol*
703 *Rev* **94**, 355-82 (2014).
- 704 44. Stenzel, D. & Huttner, W.B. Role of maternal thyroid hormones in the developing
705 neocortex and during human evolution. *Front Neuroanat* **7**, 19 (2013).
- 706 45. Forhead, A.J. & Fowden, A.L. Thyroid hormones in fetal growth and parturition
707 maturation. *J Endocrinol* **221**, R87-R103 (2014).
- 708 46. McCoy, D.D. *et al.* Enhanced insulin clearance in mice lacking TRPM8 channels. *Am J*
709 *Physiol Endocrinol Metab* **305**, E78-88 (2013).
- 710 47. Reimundez, A. *et al.* Deletion of the Cold Thermoreceptor TRPM8 Increases Heat Loss
711 and Food Intake Leading to Reduced Body Temperature and Obesity in Mice. *J Neurosci*
712 **38**, 3643-3656 (2018).
- 713 48. Chen, J. *et al.* Genome-wide association study of type 2 diabetes in Africa. *Diabetologia*
714 **62**, 1204-1211 (2019).
- 715 49. Dobson, J.E. The iodine factor in health and evolution. *Geographical Review* **88**, 3-28
716 (1998).
- 717 50. Hach, F. *et al.* mrsFAST: a cache-oblivious algorithm for short-read mapping. *Nat*
718 *Methods* **7**, 576-7 (2010).
- 719 51. Steinberg, K.M. *et al.* Structural diversity and African origin of the 17q21.31 inversion
720 polymorphism. *Nat Genet* **44**, 872-80 (2012).
- 721 52. Vollger, M.R. *et al.* Long-read sequence and assembly of segmental duplications. *Nat*
722 *Methods* **16**, 88-94 (2019).
- 723 53. Kalyanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A. & Jermini, L.S.
724 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* **14**,
725 587-589 (2017).
- 726 54. Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis.
727 *PLoS Comput Biol* **10**, e1003537 (2014).
- 728 55. Keane, T.M., Creevey, C.J., Pentony, M.M., Naughton, T.J. & McLnerney, J.O.
729 Assessment of methods for amino acid matrix selection and their use on empirical data
730 shows that ad hoc assumptions for choice of matrix are not justified. *BMC Evol Biol* **6**, 29
731 (2006).
- 732 56. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
733 polymorphism. *Genetics* **123**, 585-95 (1989).
- 734 57. Chen, G.K., Marjoram, P. & Wall, J.D. Fast and flexible simulation of DNA sequence
735 data. *Genome Res* **19**, 136-42 (2009).

- 736 58. Gutenkunst, R.N., Hernandez, R.D., Williamson, S.H. & Bustamante, C.D. Inferring the
737 joint demographic history of multiple populations from multidimensional SNP frequency
738 data. *PLoS Genet* **5**, e1000695 (2009).
- 739 59. Thomson, R., Pritchard, J.K., Shen, P., Oefner, P.J. & Feldman, M.W. Recent common
740 ancestry of human Y chromosomes: evidence from DNA sequence data. *Proc Natl Acad*
741 *Sci U S A* **97**, 7360-5 (2000).
- 742 60. Lewontin, R.C. The Interaction of Selection and Linkage. I. General Considerations;
743 Heterotic Models. *Genetics* **49**, 49-67 (1964).
- 744 61. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
745 datasets. *Gigascience* **4**, 7 (2015).
- 746

Table 1. Summary of 15 assembled *TCAF* haplotypes constructed using large-insert BAC libraries and long-read sequencing. BAC clones were selected and sequenced using the PacBio long-read sequencing technology and assembled into individual haplotypes (**Methods**). Copy number of *TCAF* segmental duplication (SD) cassettes and the classification for individual haplotypes were determined by Miropeats and sequence alignment analysis (**Figures 2 and S5-S13**).

Haplotype ID	BAC library (species or population)	Length (bp)	Length of <i>TCAF</i> SD cassettes (bp)	Copy number of <i>TCAF</i> SD cassettes	%GC	Haplogroup
CHM1	CHM1	368,013	277,806	2	39.83	Haplogroup 2-1
VMRC53_hapA	NA12878 (European)	433,048	277,806	2	39.56	Haplogroup 2-2
VMRC53_hapB	NA12878 (European)	425,306	273,483	2	39.63	Haplogroup 3-2
VMRC61_hapA	HG00732 (Puerto Rican)	337,690	277,856	2	39.91	Haplogroup 2-2
VMRC61_hapB	HG00732 (Puerto Rican)	435,583	405,366	3	39.71	Haplogroup 4
VMRC62_hapA	HG00733 (Puerto Rican)	395,405	277,854	2	39.60	Haplogroup 2-2
VMRC64_hapA	NA19240 (Yoruba)	323,367	260,853	2	39.94	Haplogroup 2-2
VMRC64_hapB	NA19240 (Yoruba)	348,654	277,808	2	40.07	Haplogroup 2-2
VMRC66_hapA	NA19434 (Luhya)	496,357	406,131	3	40.00	Haplogroup 5
VMRC69_hapA	HG00514 (Han Chinese)	387,079	277,712	2	39.29	Haplogroup 3-1
VMRC73_hapA	GM10539 (Melanesian)	247,628	145,427	1	39.93	Haplogroup 1
VMRC73_hapB	GM10539 (Melanesian)	222,558	145,424	1	39.85	Haplogroup 1
CH251_contig	CH251 (Pan troglodytes)	273,442	127,988	1	39.90	Ancestral-like
CH277_contig	CH277 (Gorilla gorilla)	241,956	140,234	1	39.98	Ancestral-like
CH250_contig	CH250 (Rhesus macaque)	225,909	140,184	1	40.49	Ancestral-like

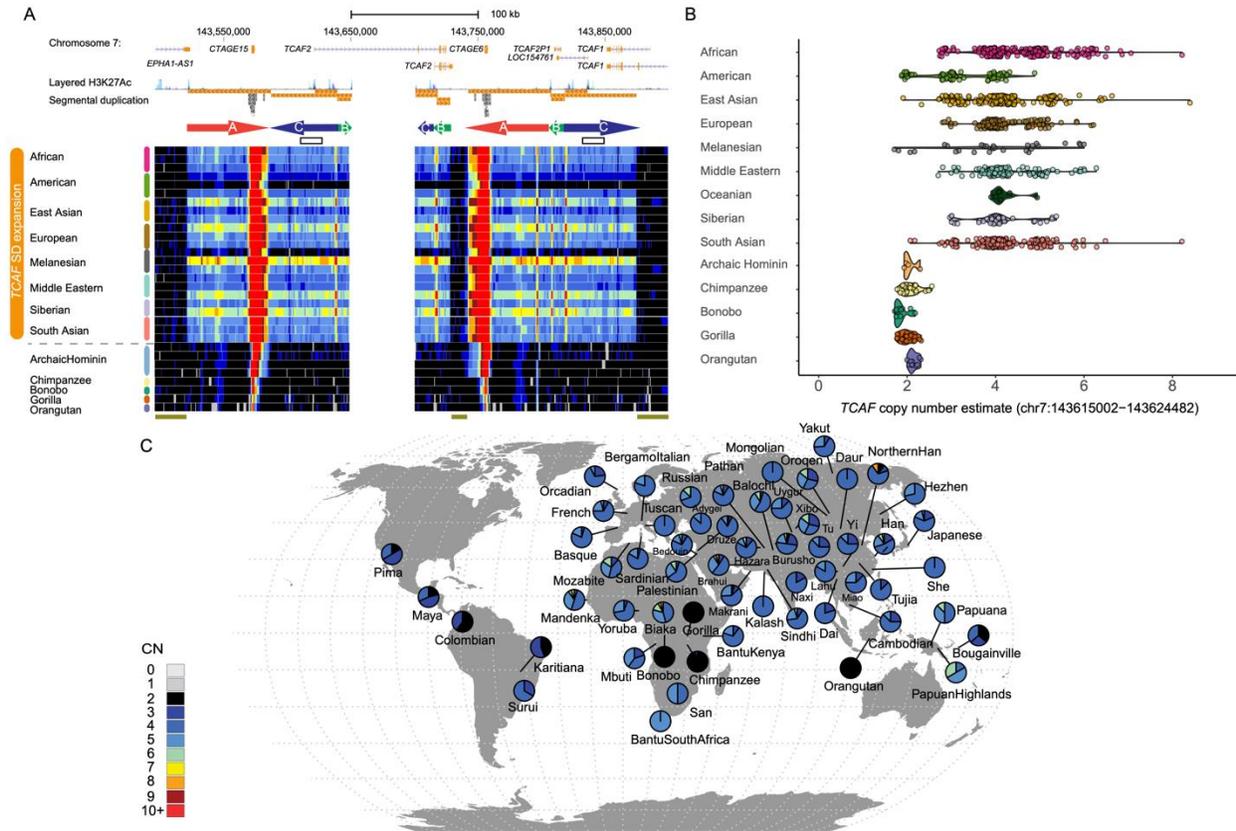


Figure 1. Copy number variation of *TCAF* SDs in a collection of diverse human and nonhuman samples. Copy numbers were estimated using read-depth based genotyping method. (A) Copy number (CN) heat maps, where each row represents the CN of a sample over the *TCAF* locus. The colored arrows (A, B, and C) represent the three major SD blocks in this region. The white area in the middle represents the gap present in the human reference genome (GRCh38). The two white boxes show a putative interlocus gene conversion event that correlates with latitudinal locations of populations (**Figure S16-S18**). Dark green bars at the bottom indicate the unique diploid sequences used for downstream phylogenetic and population genetic analyses (chr7:143,501,000-143,521,000, chr7:143,729,525-143,741,525, chr7:143,875,000-143,895,000 [GRCh38]). (B) Distributions of the overall *TCAF* CN genotypes among samples from nonhuman great apes, archaic hominins, and modern humans using a representative region (chr7:143,615,002–143,624,482). (C) Geographic distribution of the overall *TCAF* CN genotypes in the 54 Human Genome Diversity Project (HGDP) populations and the nonhuman great ape samples. Geo-coordinates for the populations are listed in **Table S8**. Pie charts shows the CN frequency distribution for a given population (colors correspond to those in the CN heat maps).

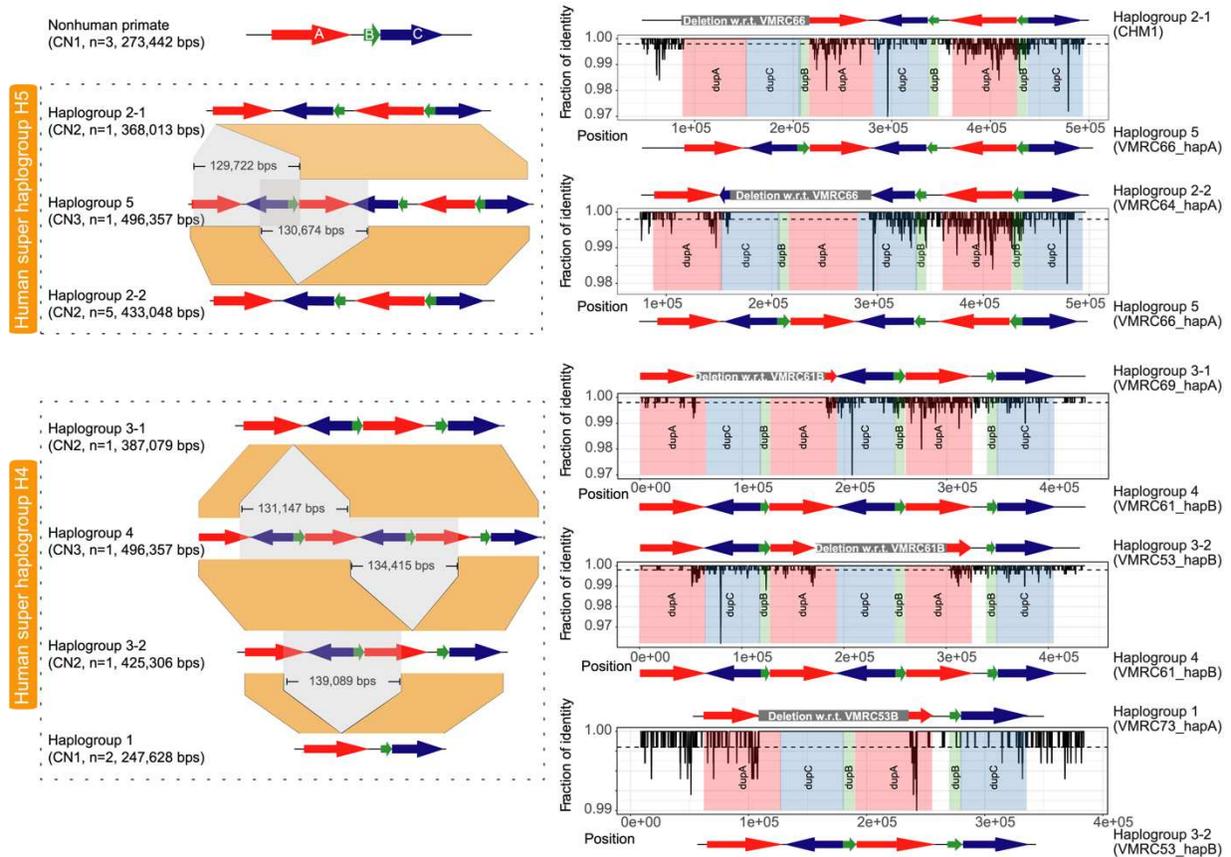


Figure 2. Complex structural haplotype diversity at the 7q35 *TCAF* locus in humans. Left panel: Schematic representations of the major- (super haplogroups H4 and H5) and sub-haplogroups identified from 15 targeted BAC long-read assemblies. Colored arrows are *TCAF* SDs, and the gray and orange areas indicate putative deletion events and orthologous sequences between haplogroups. For example, despite the similarity in structure, Haplogroup 2-1 and Haplogroup 2-2 are likely resulted from different deletion events from Haplogroup 5 based on comparative sequence data. Right panel: The comparative sequence analysis between individual haplogroups shows pairwise sequence identities (black lines) over 500 bp windows (sliding by 100 bp). With each plot, colored rectangles correspond to *TCAF* SDs, while the white region embedded between a DupA (red) and a DupB (green) blocks represents a 12.3 kbp single-copy unique sequence. Gray rectangles above each sequence identity plots indicate the locations of putative structure variants (**Table S3**).

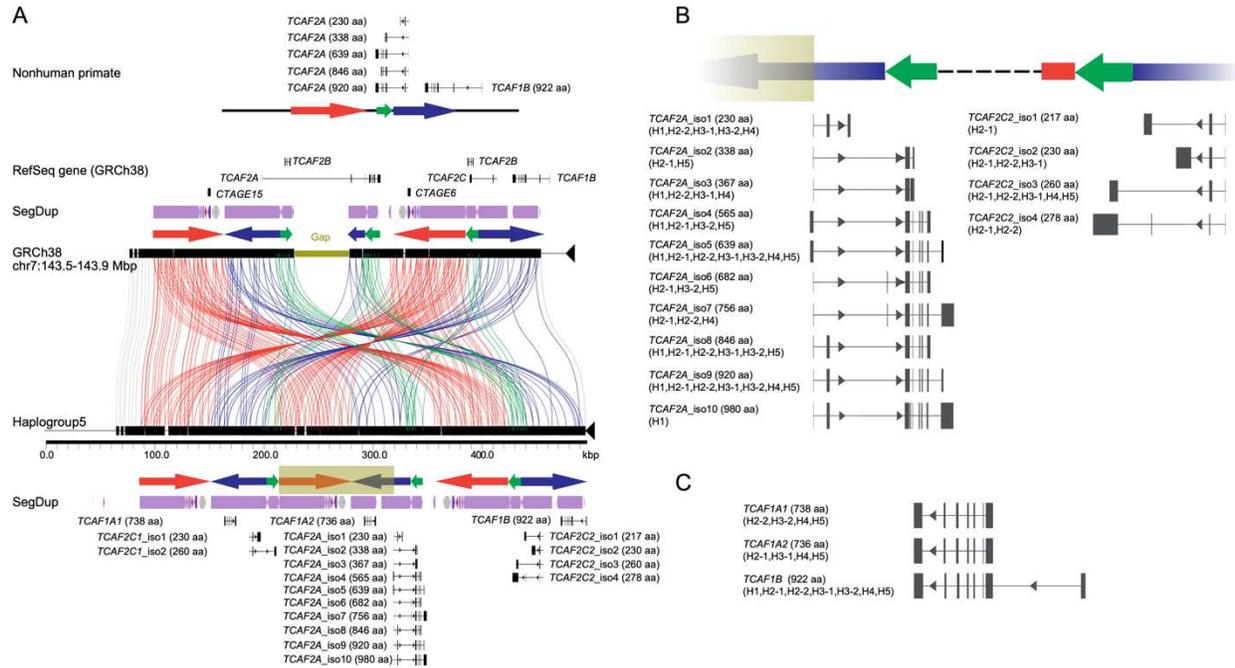


Figure 3. *TCAF* gene and isoform diversity among human haplotypes. (A) Miroppeats analysis reveals a great consistency between Haplogroup 5 and the gapped 7q35 *TCAF* locus in the human reference genome (GRCh38). Colored arrows are annotated *TCAF* SDs and lines connecting the sequences show regions of homology. The dark green rectangle in the sequence-resolved, gap-free Haplogroup 5 contig (bottom) represents the region missing in the GRCh38 sequence (dark green bar, top). Additional annotations include a schematic of *TCAF* sequence structure and gene models/isoforms for nonhuman primate, RefSeq gene track (GRCh38), segmental duplication (SegDup) tracks, BAC clones used in the Haplogroup 5 assembly, and predicted *TCAF* models and isoforms using full-length non-chimeric transcripts from six human tissues (**Methods**). Note that the *TCAF2A* gene model in GRCh38 is incorrect due to the presence of the gap in the middle of the *TCAF* SD region. Also note that the numbers above overlap between two BAC clones and indicate the percent sequence identity (#identical bases/total bases). For illustration, all predicted gene models and isoforms are aligned to Haplogroup 5 sequence for (B) *TCAF2A* and (C) *TCAF2C*, along with annotations for amino acid sequence lengths and haplogroups in which they are observed. For simplicity, we skip sequences between *TCAF2A* and *TCAF2C2* (dashed line). Detailed haplotype-specific and/or tissue-specific gene models and isoforms can be found in **Figures S6-S13** and **Table S5**.

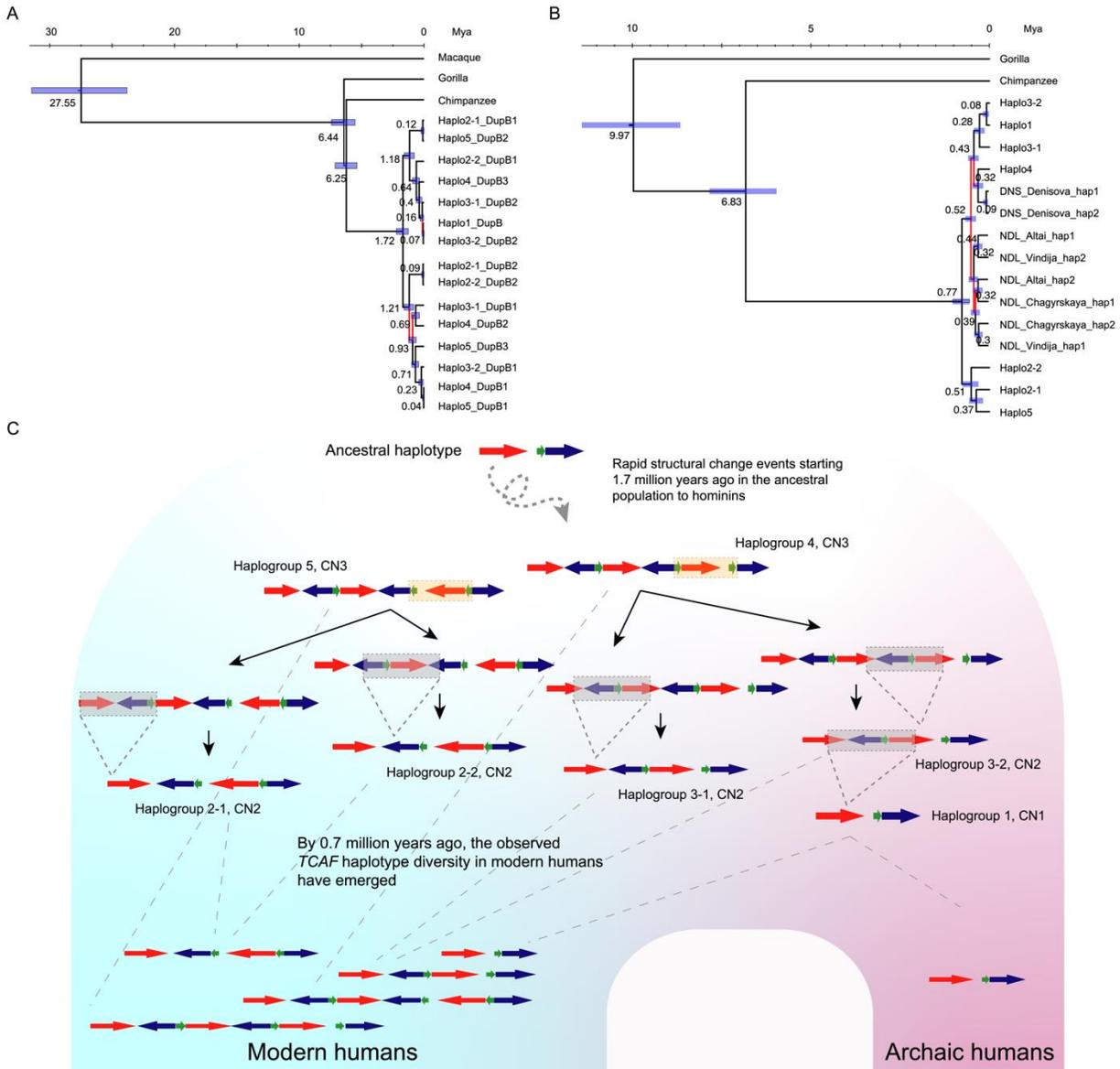


Figure 4. Evolutionary reconstruction of *TCAF* structural diversity. (A) Phylogeny of the haplogroups was inferred using *TCAF* DupB sequences and BEAST (v.2.6.2) with five independent runs of 10 million iterations of Markov Chain Monte Carlo (**Methods**). Numbers and horizontal bars at internal nodes indicate point estimates and 95% highest posterior density intervals for the divergences (in million years ago, Mya), respectively. Branches with posterior probabilities <90% are colored in red. See **Figure S14-S16** for results of other SD sequences. (B) Inferred phylogeny of the modern human haplogroups and archaic hominin haplotypes using the 12 kbp unique sequences embedded within *TCAF* SDs (**Figure 1**). Haplotypes of archaic samples were generated using high-confident single-nucleotide variants (SNVs) called within the unique diploid region. Phylogenetic inference was performed similarly as described above. (C) Schematic model for the evolution of *TCAF* haplotypes in humans based on phylogenetic inferences (**Figures 4A-4B and S14-S17**). Colored arrows are *TCAF* SDs; orange and gray areas indicate relative inversion and deletion events between haplogroups, respectively. Short dashed lines indicate putative breakpoints of structural changes between haplogroups, while the long-dashed lines illustrate lineage sorting.

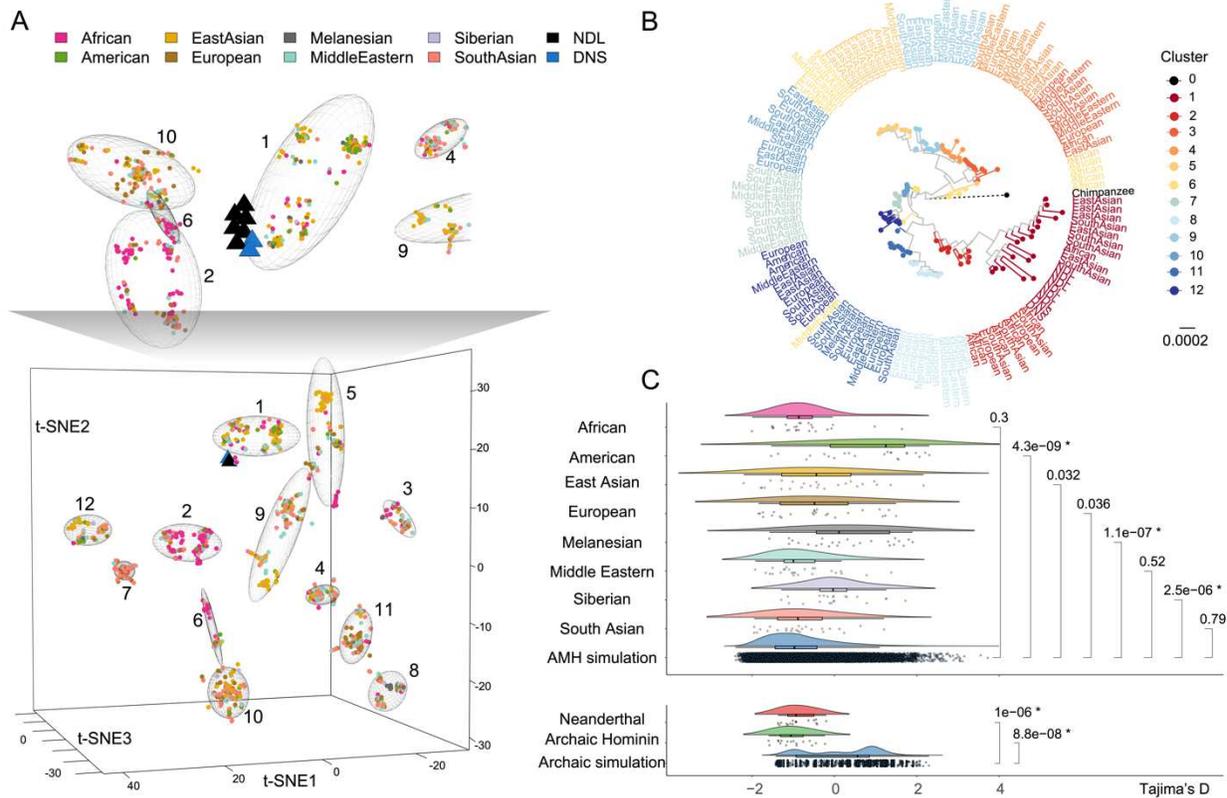


Figure 5. Archaic hominin versus human haplotype diversity. Haplotypes were inferred using 1,275 SNVs in the three unique diploid sequences around the *TCAF* SD region (**Figures 1 and 2**). (A) Haplotype-based principle component analysis was performed, followed by haplotype clustering and cluster visualization using the dimension-reduction technique, t-SNE (**Methods**). On the t-SNE plot, each dot/triangle is a haplotype and colored according to population/species origin. NDL and DNS refer to Neanderthal (black triangle) and Denisovan (blue triangle) haplotypes. Numbers and ellipses in the 3D t-SNE plots indicate individual clusters (see also **Figures S24-S29**). The zoom-in above the 3D t-SNE shows that all archaic haplotypes are in tight proximity to each other. (B) The maximum likelihood phylogeny was constructed using 10 randomly selected haplotypes from the 12 inferred clusters, in addition to eight archaic and one chimpanzee haplotypes. Note that the branch length of chimpanzee (dashed line) is truncated by 90% of its actual length for the purpose of illustration. (C) Distributions of Tajima's D statistic computed using 2 kbp windows sliding over the three unique diploid regions. Significance of natural selection signals was determined using coalescent simulations based on 1,000 different demographic models (**Methods**). The asterisks indicate significant tests after the Bonferroni correction.

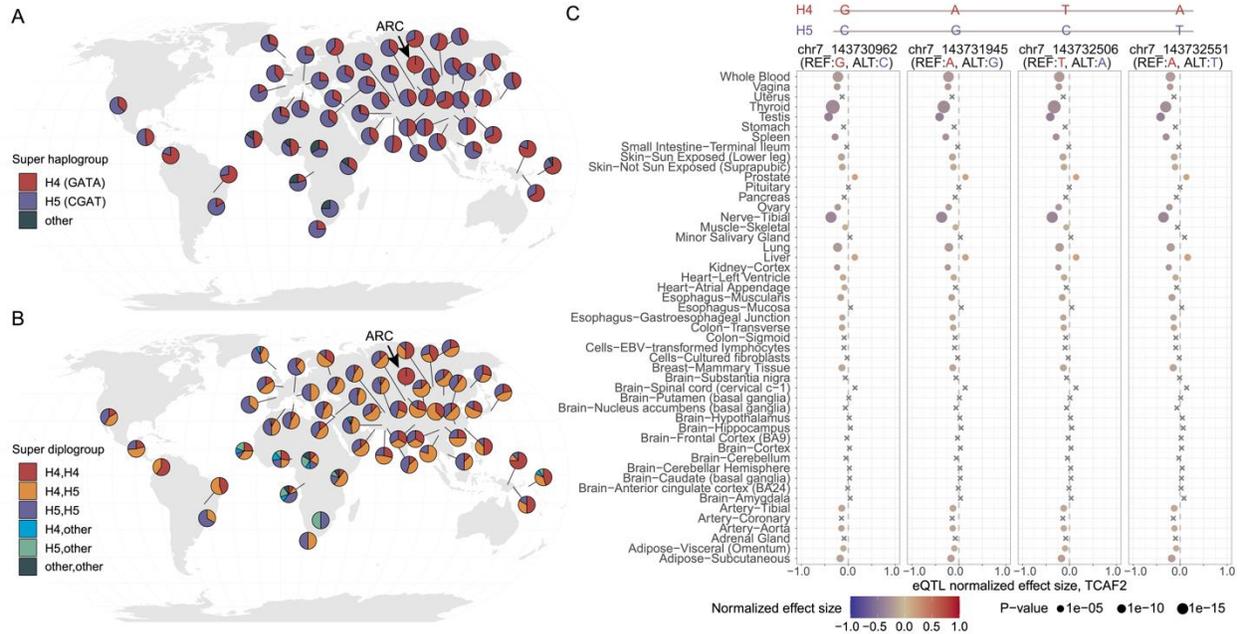


Figure 6. H4 and H5 super haplogroup distribution and eQTL analysis. Four tagging SNVs were first identified by perfectly separating super haplogroups H4 and H5 among the seven BAC haplogroups and confirmed based on patterns of linkage disequilibrium (LD) in samples from the HGDP panel. (A and B) Distributions of super haplogroups H4 (red, 57.1%) and H5 (blue, 40.3%) and their diploid type in the HGDP populations across the world. Other haplotypes were found in ~2.6% of the samples. Note that archaic hominin haplotypes from the Neanderthal ($n = 3$) and Denisovan ($n = 1$) samples all carry H4 haplotypes, and as a representation the geographic location of these samples were placed in the Altai Mountains in Siberia (arrows). (C) Multi-tissue eQTL plots show consistent patterns of associations between the four tagging SNVs and expression levels of *TCAF2* across 50 tissues. Effect sizes were calculated as the effect of the alternative allele (blue) relative to the reference allele (red) as defined in GTEx (release v8) and are scaled using color. The (unadjusted) p values of eQTL association are represented by dot sizes. Note that crosses (x) indicate insignificant associations.

Figures

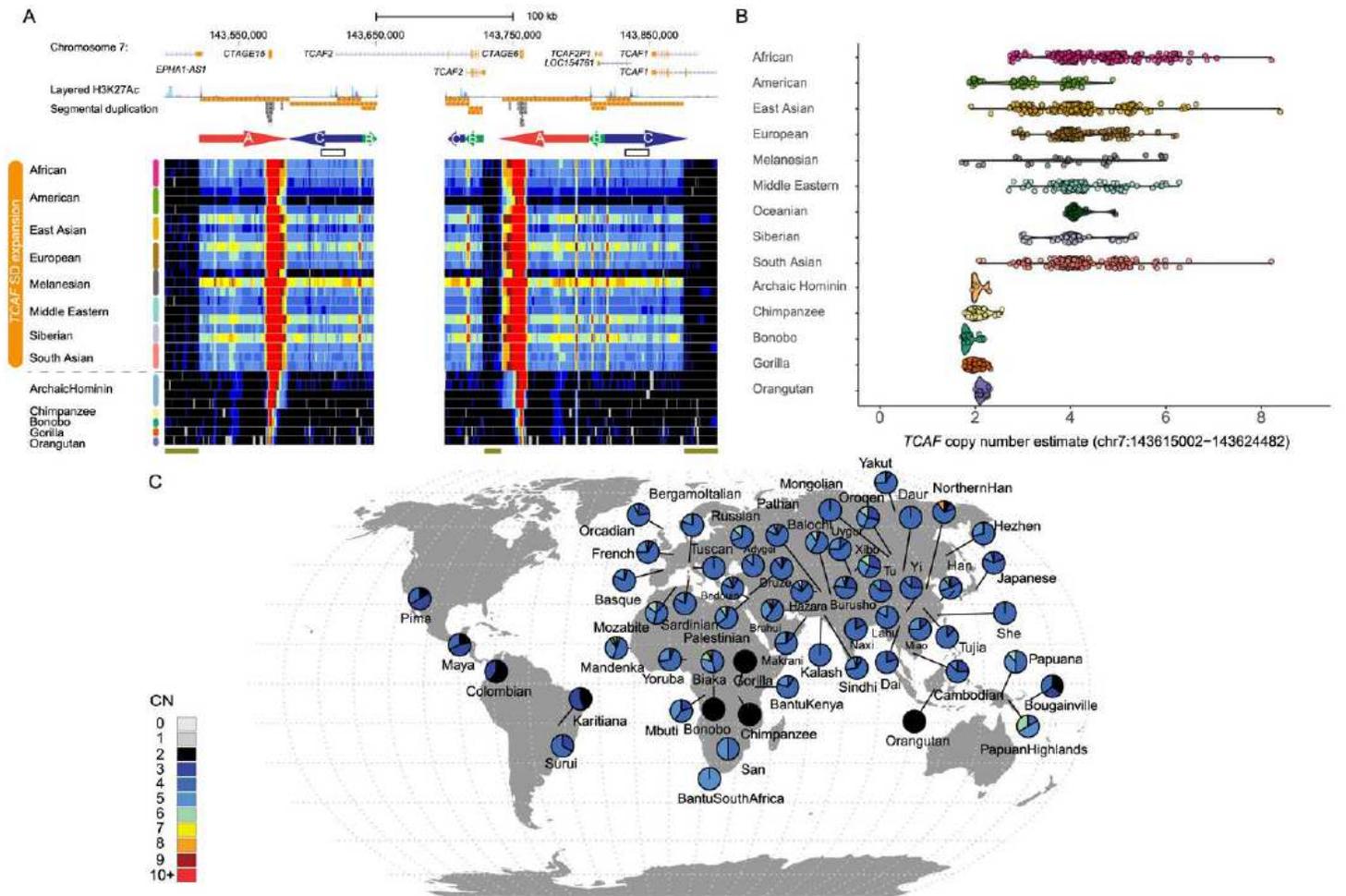


Figure 1

Copy number variation of TCAF SDs in a collection of diverse human and nonhuman samples. Copy numbers were estimated using read-depth based genotyping method. (A) Copy number (CN) heat maps, where each row represents the CN of a sample over the TCAF locus. The colored arrows (A, B, and C) represent the three major SD blocks in this region. The white area in the middle represents the gap present in the human reference genome (GRCh38). The two white boxes show a putative interlocus gene conversion event that correlates with latitudinal locations of populations (Figure S16-S18). Dark green bars at the bottom indicate the unique diploid sequences used for downstream phylogenetic and population genetic analyses (chr7:143,501,000-143,521,000, chr7:143,729,525-143,741,525, chr7:143,875,000-143,895,000 [GRCh38]). (B) Distributions of the overall TCAF CN genotypes among samples from nonhuman great apes, archaic hominins, and modern humans using a representative region (chr7:143,615,002-143,624,482). (C) Geographic distribution of the overall TCAF CN genotypes in the 54 Human Genome Diversity Project (HGDP) populations and the nonhuman great ape samples. Geo-coordinates for the populations are listed in Table S8. Pie charts show the CN frequency distribution for a given population (colors correspond to those in the CN heat maps). Note: The designations employed

and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

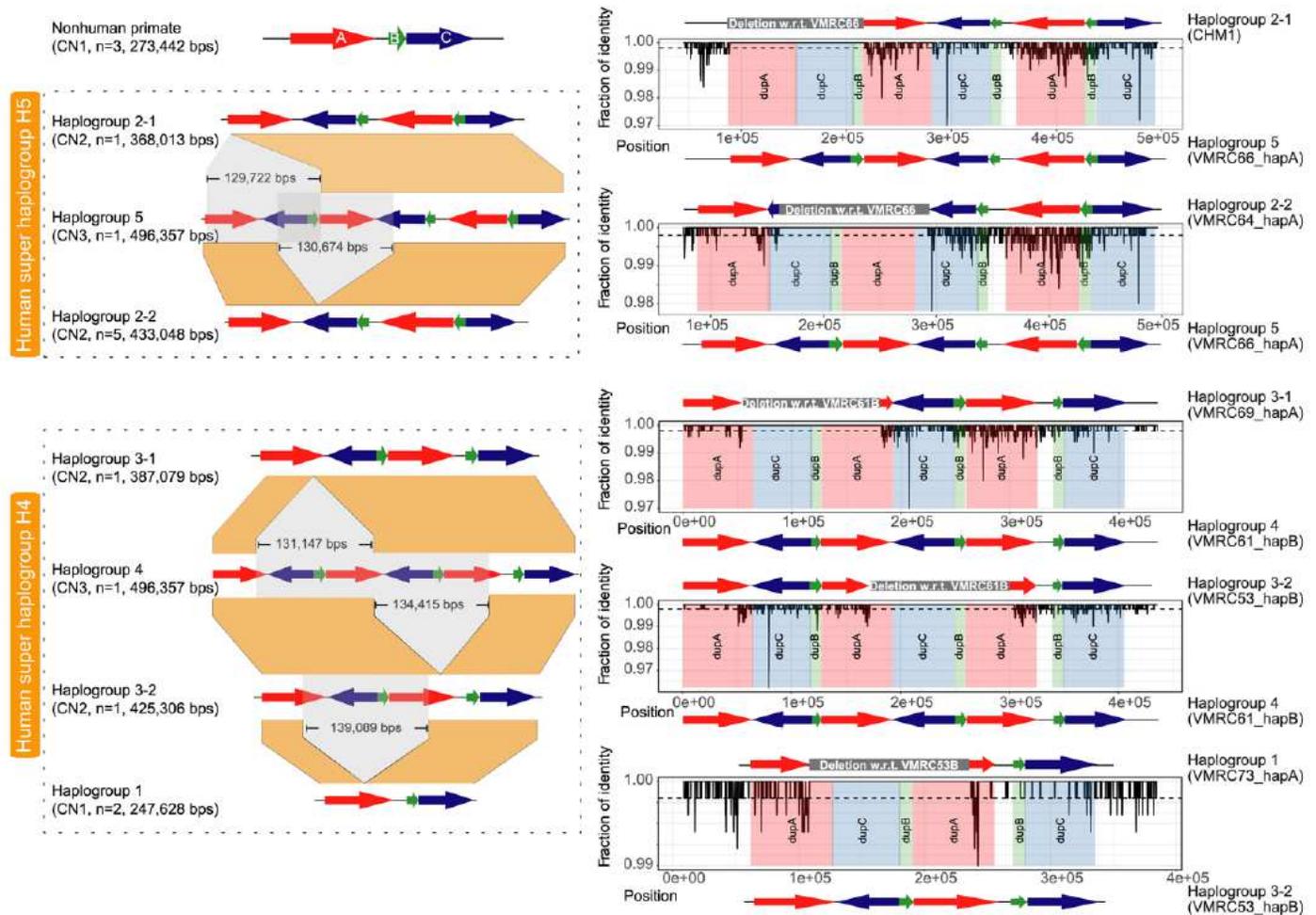


Figure 2

Complex structural haplotype diversity at the 7q35 TCAF locus in humans. Left panel: Schematic representations of the major- (super haplogroups H4 and H5) and sub-haplogroups identified from 15 targeted BAC long-read assemblies. Colored arrows are TCAF SDs, and the gray and orange areas indicate putative deletion events and orthologous sequences between haplogroups. For example, despite the similarity in structure, Haplogroup 2-1 and Haplogroup 2-2 are likely resulted from different deletion events from Haplogroup 5 based on comparative sequence data. Right panel: The comparative sequence analysis between individual haplogroups shows pairwise sequence identities (black lines) over 500 bp windows (sliding by 100 bp). With each plot, colored rectangles correspond to TCAF SDs, while the white region embedded between a DupA (red) and a DupB (green) blocks represents a 12.3 kbp single-copy unique sequence. Gray rectangles above each sequence identity plots indicate the locations of putative structure variants (Table S3).

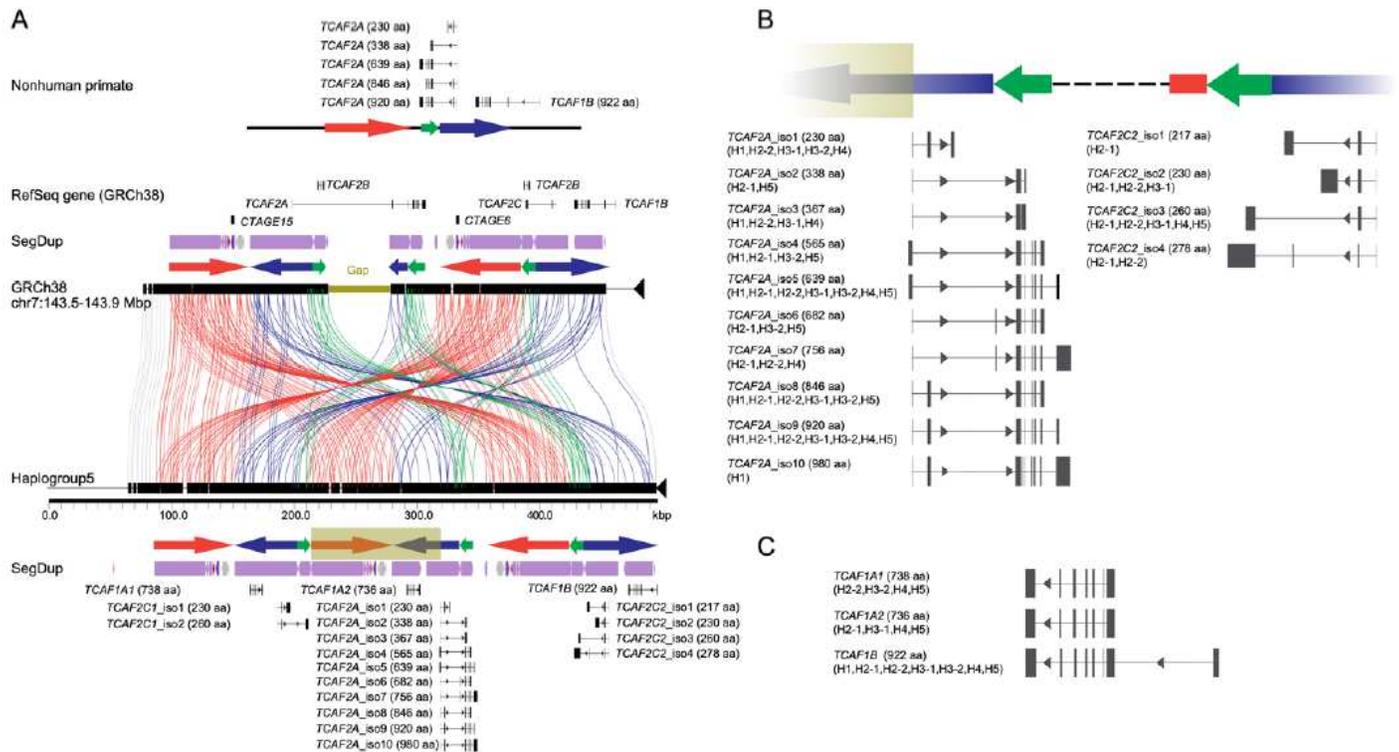


Figure 3

TCAF gene and isoform diversity among human haplotypes. (A) Miropeats analysis reveals a great consistency between Haplogroup 5 and the gapped 7q35 TCAF locus in the human reference genome (GRCh38). Colored arrows are annotated TCAF SDs and lines connecting the sequences show regions of homology. The dark green rectangle in the sequence-resolved, gap-free Haplogroup 5 contig (bottom) represents the region missing in the GRCh38 sequence (dark green bar, top). Additional annotations include a schematic of TCAF sequence structure and gene models/isoforms for nonhuman primate, RefSeq gene track (GRCh38), segmental duplication (SegDup) tracks, BAC clones used in the Haplogroup 5 assembly, and predicted TCAF models and isoforms using full-length non-chimeric transcripts from six human tissues (Methods). Note that the TCAF2A gene model in GRCh38 is incorrect due to the presence of the gap in the middle of the TCAF SD region. Also note that the numbers above overlap between two BAC clones and indicate the percent sequence identity (#identical bases/total bases). For illustration, all predicted gene models and isoforms are aligned to Haplogroup 5 sequence for (B) TCAF2A and (C) TCAF2C, along with annotations for amino acid sequence lengths and haplogroups in which they are observed. For simplicity, we skip sequences between TCAF2A and TCAF2C2 (dashed line). Detailed haplotype-specific and/or tissue-specific gene models and isoforms can be found in Figures S6-S13 and Table S5.

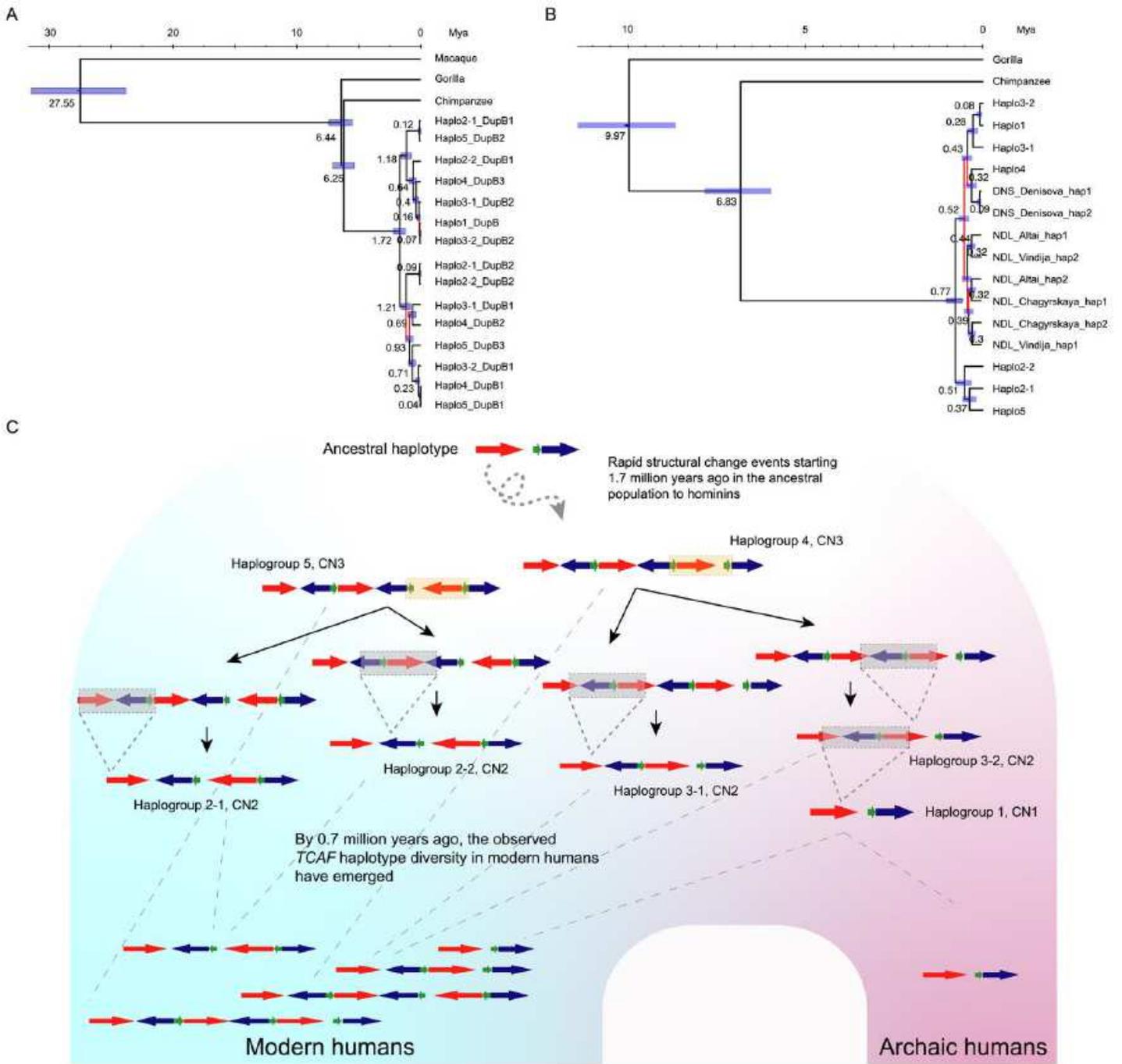


Figure 4

Evolutionary reconstruction of TCAF structural diversity. (A) Phylogeny of the haplogroups was inferred using TCAF DupB sequences and BEAST (v.2.6.2) with five independent runs of 10 million iterations of Markov Chain Monte Carlo (Methods). Numbers and horizontal bars at internal nodes indicate point estimates and 95% highest posterior density intervals for the divergences (in million years ago, Mya), respectively. Branches with posterior probabilities <90% are colored in red. See Figure S14-S16 for results of other SD sequences. (B) Inferred phylogeny of the modern human haplogroups and archaic hominin haplotypes using the 12 kbp unique sequences embedded within TCAF SDs (Figure 1). Haplotypes of archaic samples were generated using high-confident single-nucleotide variants (SNVs) called within the

unique diploid region. Phylogenetic inference was performed similarly as described above. (C) Schematic model for the evolution of TCAF haplotypes in humans based on phylogenetic inferences (Figures 4A-4B and S14-S17). Colored arrows are TCAF SDs; orange and gray areas indicate relative inversion and deletion events between haplogroups, respectively. Short dashed lines indicate putative breakpoints of structural changes between haplogroups, while the long-dashed lines illustrate lineage sorting.

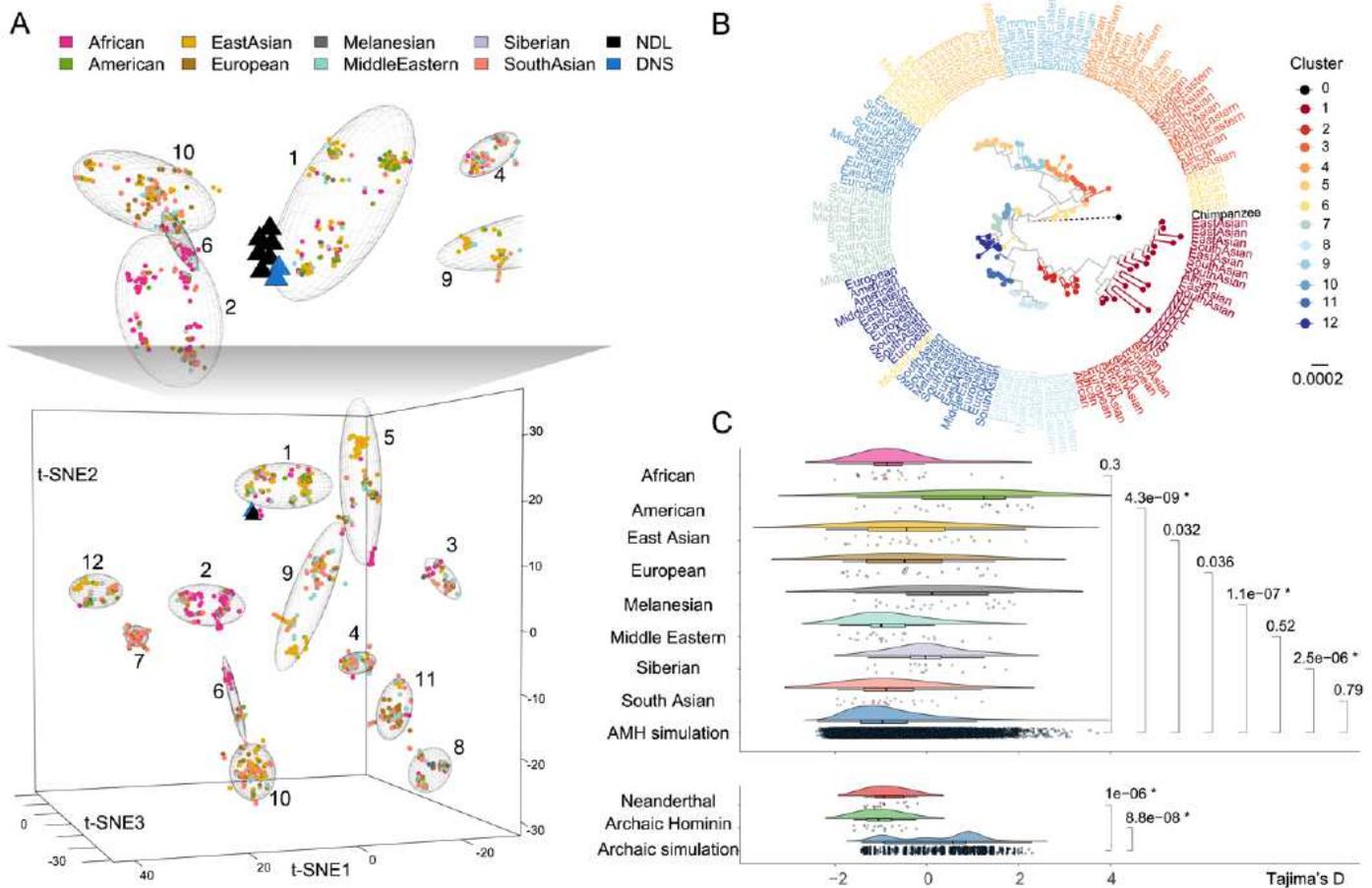


Figure 5

Archaic hominin versus human haplotype diversity. Haplotypes were inferred using 1,275 SNVs in the three unique diploid sequences around the TCAF SD region (Figures 1 and 2). (A) Haplotype-based principle component analysis was performed, followed by haplotype clustering and cluster visualization using the dimension-reduction technique, t-SNE (Methods). On the t-SNE plot, each dot/triangle is a haplotype and colored according to population/species origin. NDL and DNS refer to Neanderthal (black triangle) and Denisovan (blue triangle) haplotypes. Numbers and ellipses in the 3D t-SNE plots indicate individual clusters (see also Figures S24-S29). The zoom-in above the 3D t-SNE shows that all archaic haplotypes are in tight proximity to each other. (B) The maximum likelihood phylogeny was constructed using 10 randomly selected haplotypes from the 12 inferred clusters, in addition to eight archaic and one chimpanzee haplotypes. Note that the branch length of chimpanzee (dashed line) is truncated by 90% of its actual length for the purpose of illustration. (C) Distributions of Tajima's D statistic computed using 2 kbp windows sliding over the three unique diploid regions. Significance of natural selection signals was

determined using coalescent simulations based on 1,000 different demographic models (Methods). The asterisks indicate significant tests after the Bonferroni correction.

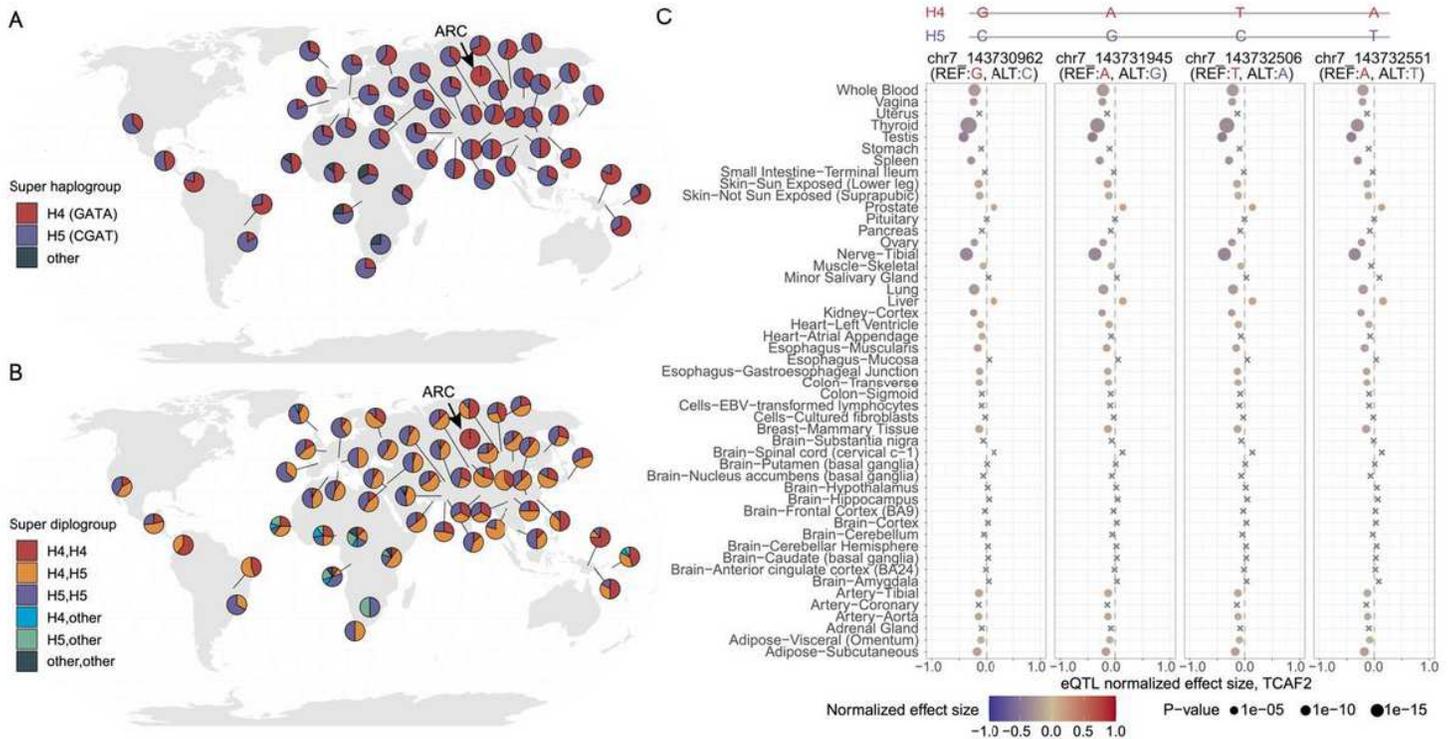


Figure 6

H4 and H5 super haplogroup distribution and eQTL analysis. Four tagging SNVs were first identified by perfectly separating super haplogroups H4 and H5 among the seven BAC haplogroups and confirmed based on patterns of linkage disequilibrium (LD) in samples from the HGDP panel. (A and B) Distributions of super haplogroups H4 (red, 57.1%) and H5 (blue, 40.3%) and their diploid type in the HGDP populations across the world. Other haplotypes were found in ~2.6% of the samples. Note that archaic hominin haplotypes from the Neanderthal ($n = 3$) and Denisovan ($n = 1$) samples all carry H4 haplotypes, and as a representation the geographic location of these samples were placed in the Altai Mountains in Siberia (arrows). (C) Multi-tissue eQTL plots show consistent patterns of associations between the four tagging SNVs and expression levels of TCAF2 across 50 tissues. Effect sizes were calculated as the effect of the alternative allele (blue) relative to the reference allele (red) as defined in GTEx (release v8) and are scaled using color. The (unadjusted) p values of eQTL association are represented by dot sizes. Note that crosses (x) indicate insignificant associations. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Supplementary Files

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

- [TCAFSupplementaryFiguresandTables.v8.pdf](#)
- [TableS1VSTpairwiseSuperPop.sorted.v8.xlsx](#)
- [TableS7GENECONVIGCtracks.v8.xlsx](#)
- [TableS8GeographicCoordinatesHGDPpopsandCNsummary.v8.xlsx](#)