

Supplementary Table S5. Polymorphism and divergence in hominids at homologous positions in the 1'002 bp sequence including the *NAT* pseudogene (*NATP*).

The screened sequence spans from 18'228'116 to 18'229'117 on chromosome 8 in the human reference sequence GRCh37/hg19. The 77 substitutions corresponding to inter-species divergence (between two species at least) are highlighted by stars to the left part of the table. Polymorphisms shared between species or subspecies are highlighted in blue font. Deletions are indicated by del; (N) indicates that the position is undefined in some samples. Variants also recorded in sequenced ancient genomes of hominids (Denisova, Neanderthal or 45'000 years old *Homo sapiens* from Ust'-Ishim) are boxed (see Supplementary Table S8).

			A	A	A	A	A	A	A	A	A	A/T	A/T
26	18228303		C	C	C	C	C	T	C	C	C	C	C
* 27	18228304		C	C	C	G	G	C	C	C	C	C	C
28	18228306		C	G	G	G	G	G	G	G	G	C/T	C
* 29	18228315		G	G	G	G	G	G	G	G	T	T	T
* 30	18228320		T	G	G	G	G	G	G	T	T	T	T
* 31	18228327		G	G	G	G	G	G	G	G	C	C	C
32	18228329	rs547876229	G/A	G	G	G	G	G	G	G	G	G	G
* 33	18228332		G	G	G	G	G	G	G	G	G	A	A
* 34	18228335		C	C	C	C	C	C	C	C	C	G	G
* 35	18228336		G	A	A	A	A	A	A	A	A	A	A
* 36	18228367		C	T	T	T	T	T	T	C	C	C	C
37	18228368		T	T/C	T	T	T	T/C	T	T	T	T	T
38	18228373		G	G	G	G	G	G	G	G	G	G/T(N)	G/T
* 39	18228382		A	G	G	G	G	G	G	G	G	G	G
* 40	18228385		G	G	G	G	G	G	G	G	G	A	A
41	18228389	rs561082251	G/A	G	G	G	G	G	G	G	G	G	G
* 42	18228392		A	T	T	T	T	T	T	T	T	T	T
* 43	18228395		G	G	G	G	G	G	G	A	A	G	G
44	18228404		C	C	C/T	C	C	C	C	C	C	C	C
45	18228405	rs530061116	C/T	C	C	C	C	G	G	G	G	C	C
46	18228409		G/A	G	G	G	G	G	G	G	G	G	G
* 47	18228418		A	A	A	A	A	A	A	A	A	A	G
* 48	18228419		T	C	C	C	C	C	C	T	T	T	T
* 49	18228420		A	A	A	A	A	A	A	A	A	G	G
* 50	18228424	rs184770517	C/T	T	T	T	T	T	T	C	C	T	T
51	18228425	rs570049254	G/A	G	G	G	G	G	G	G	G	G	G
* 52	18228437		C	C	C	C	C	C	C	A	A	C	C
53	18228442	rs538826711	T/A	T	T	T	T	T	T	T	T	T	T
* 54	18228458	rs372738250	G/A	A	A	A	A	A	A	G/A	G	G	G
* 55	18228463		C	C	C	C	C	C	C	T	T	T	T
* 56	18228470		T	T	T	T	T	T	T	A	A	T	T
* 57	18228480		C	C	C	C	C	C	C	C	C	T	T
58	18228483	rs188523053	A/T	A	A	A	A	A	A	A	A	A	A
59	18228501		C	C/T	C/T	C/T	C/T	C/T	C	C	C	C	C
60	18228502	rs534373318	C/T	C	C	C	C	C	C	C	C	C	C
* 61	18228519		C	T	T	T	T	T	T	T	T	T	T
62	18228543		A	A	A/G	A	A	A	A	A	A	A	A
63	18228545	rs554671092	C/T	C	C	C	C	C	C	C	C	C	C
* 64	18228552		A	A	A	A	A	A	A	A	A	G	G
* 65	18228560		C	C/A	C	C	C	C	C	C	C	T	T
66 ⁸	18228571	rs180998104	C/T	C(N)	C(N)	C(N)	C(N)	C(N)	C	C	C	C(N)	C

147	18229103		A	A	A	A/T	A	A	A	A	A
148	18229104	rs74444655	T/C	T	T	T	T	T	T	T(N)	T
149	18229105		A/C	A	A	A	A	A	A	A(N)	A
*	150	18229106		T	T	T	T	T	T	C	C
	151	18229111		A	A	A	A	A	A	A/G	A

¹ Human polymorphisms are those recorded by the consensus gene nomenclature of human *NAT* alleles (<http://nat.mbg.duth.gr/>), complemented with haplotype data from 1000 Genomes Phase 1 data (The Genomes Project Consortium 2012), (Patin et al. 2006a) and (Mortensen et al. 2011), as well as 1000 Genomes SNPs and indels reported for GRCh37.p13 (hg19) in Ensembl (Yates et al. 2016) (http://grch37.ensembl.org/Homo_sapiens/Info/Index, accessed October 2016). Human positions not recorded as polymorphic but associated with a SNP identifier are reported with a highest population MAF < 0.01 in Ensembl (http://www.ensembl.org/Homo_sapiens/Info/Index). Note that the list is not exhaustive: a polymorphic position (A/G) at 18'228'182 is reported for Denisova in the ancient genome browser of the Max Planck Institute for Evolutionary Anthropology (see Supplementary Table S8); it is not recorded in the table since all humans and other hominoids are fixed on A, and an additional polymorphic position (T/C) at 18'228'748, which defines the *Pan NATP*13* haplotype, is not recorded in the table either as it was only observed in the hybrid *P. t. verus/troglodytes* individual (all humans and other hominids are fixed on T).

² Western, Niger-Cameroon, Eastern and Central chimpanzees are *Pan troglodytes verus*, *P. t. ellioti*, *P. t. schweinfurthii* and *P. t. troglodytes*, respectively. Polymorphism recording is based on the individuals of the present study (Supplementary Figure S1) and the chimpanzees of Prado-Martinez et al. (2013) cross-checked with the *P. t. verus* assembly reference sequence (panTro4, February 2011).

³ Based on the individual of this study (Bonobo), the bonobos of Prado-Martinez et al. (2013) and the *Pan paniscus* draft assembly reference sequence (panPan1, May 2012).

⁴ Based on the individuals of this study, the gorillas of Prado-Martinez et al. (2013) and the *Gorilla gorilla gorilla* draft assembly reference sequence (gorGor4, December 2014).

⁵ Based on the individuals of this study, the orangutans of Prado-Martinez et al. (2013) and the *Pongo pygmaeus abelii* draft assembly reference sequence (ponAbe2, July 2007).

⁶ All individuals from the *Pan (troglodytes and paniscus)* *Gorilla* and *Pongo* species and subspecies sequenced by Prado-Martinez et al. (2013) have undefined positions in the stretch between 18'228'271 and 18'228'282 (included). For this reason, positions in the stretch 18'228'271-18'228'282 were not considered for median-joining network construction.

⁷ We speculate that a T insertion at position 18'228'278 probably occurred on the human lineage, as it was absent in all Sanger sequenced samples of chimpanzees, bonobos, gorillas and orangutans in this study (not represented in the median-joining network, see footnote 6).

⁸ All individuals from the *Pan troglodytes* and *Pongo abelii* species and subspecies sequenced by Prado-Martinez et al. (2013) have undefined positions in the stretch between 18'228'567 and 18'228'581 (included). For this reason, positions in the stretch 18'228'567-18'228'581 were not considered for median-joining network construction.

⁹ From positions 18'228'670 onwards, all individuals from the *Pan (troglodytes and paniscus)* and *Pongo abelii* species and subspecies sequenced by Prado-Martinez et al. (2013) have undefined positions: up to 18'228'701 (included) for *P. troglodytes*, up to 18'228'681 (included) for *P. paniscus* and up to 18'228'682 (included) for *Pongo abelii*. For *Gorilla*, positions 18'228'668 and 18'228'669 are undefined, as well as all positions between 18'228'672 and 18'228'689 (included). For *P. pygmaeus*, all positions between 18'228'672 and 18'228'682 (included) are undefined. Polymorphisms in this region were detected for the individuals sequenced in this study. However,

these are all indel polymorphisms of repetitive motives (CA, CAA, etc.) and we chose not to consider the stretch running from 18'228'670 to 18'228'701 for median-joining network construction.

¹⁰ At position 18'228'678, an insertion of a CAACAAA motif was observed in all Sanger sequenced gorillas, and an insertion of a CA/AA polymorphic motif in all Sanger sequenced orangutans. In the latter, a T deletion was also found in all Sanger sequenced individuals at position 12'228'886 (not represented in the median-joining network, see footnote 9).

¹¹ Position 18'228'838 is polymorphic (A/G) in *Gorilla gorilla*, but not in *Gorilla beringei*, for which only nucleotide G (or N) is reported. At this position, the gorilla reference sequence (gorGor4) indicates nucleotide A, similarly to the human, chimpanzee, bonobo and orangutan reference sequences (hg19, panTro4, panPan1, ponAbe2). The same ambiguous pattern is found at positions 18'228'841, 18'228'852 and 18'229'005, although for the latter two, the gorilla reference sequence differs from that of the other species (at position 18'228'852, gorGor4 is T, whereas it is C in all other reference sequences, and at position 18'229'005, gorGor4 is C, whereas it is A in all other reference sequences). Thus none of these four positions, 18'228'838, 18'228'841, 18'228'852 and 18'229'005, were considered for median-joining network construction.

¹² All *Pongo* individuals sequenced by Prado-Martinez et al. (2013) have undefined positions between 18'228'876 and 18'228'887 (included). For this reason, positions in the stretch 18'228'876-18'228'887 were not considered for median-joining network construction.

¹³ All *Pongo* individuals sequenced by Prado-Martinez et al. (2013) have undefined positions between 18'228'993 and 18'229'003 (included). For this reason, positions in the stretch 18'228'993-18'229'003 were not considered for median-joining network construction.

References

- Mortensen HM, Froment A, Lema G, Bodo JM, Ibrahim M, Nyambo TB, Omar SA, and Tishkoff SA. 2011. Characterization of genetic variation and natural selection at the arylamine N-acetyltransferase genes in global human populations. *Pharmacogenomics* 12(11):1545-1558.
- Patin E, Barreiro LB, Sabeti PC, Austerlitz F, Luca F, Sajantila A, Behar DM, Semino O, Sakuntabhai A, Guiso N et al. . 2006. Deciphering the ancient and complex evolutionary history of human arylamine N-acetyltransferase genes. *American journal of human genetics* 78(3):423-436.
- Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, Veeramah KR, Woerner AE, O'Connor TD, Santpere G et al. . 2013. Great ape genetic diversity and population history. *Nature* 499(7459):471-475.
- The Genomes Project C. 2012. An integrated map of genetic variation from 1,092 human genomes. *Nature* 491:56.
- Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, Cummins C, Clapham P, Fitzgerald S, Gil L et al. . 2016. Ensembl 2016. *Nucleic Acids Res* 44(D1):D710-D716.