

**Evolutionary constraints in the  $\beta$ -globin cluster: the signature of purifying selection at the  $\delta$ -globin (*HBD*) locus and its role in developmental gene regulation**

**Supplementary Material**

Ana Moleirinho, Susana Seixas, Alexandra M. Lopes, Celeste Bento, Maria J. Prata , António Amorim

**Corresponding author:** Ana Moleirinho, IPATIMUP, Rua Dr Roberto Frias s/n, 4200-465

Porto, Portugal; Phone: (+351) 22 5570700; Fax: (+351) 22 5570799; E-mail:

[amoleirinho@ipatimup.pt](mailto:amoleirinho@ipatimup.pt)

Keywords:  $\beta$ -globin cluster, Hemoglobin switch, gene diversity, chromatin interactions

Running head: **Signature of purifying selection at  $\delta$ -globin**

Table S1 – Populations sampled by the 1000 Genomes Project

<b>Population</b>	<b>Sample size</b>
CEU - Utah residents (CEPH) with Northern and Western European ancestry (CEU)	85
<b>European</b> FIN-Finnish from Finland	93
GBR - British from England and Scotland	89
IBS - Iberian populations in Spain	14
TSI - Toscani in Italia	98
<b>African</b> ASW - African Ancestry in Southwest US	61
LWK - Luhya in Webuye, Kenya	97
YRI - Yoruba in Ibadan, Nigeria	88
<b>Asian</b> CHB - Han Chinese in Beijing, China	97
CHS - Han Chinese South	100
JPT - Japanese in Toyko, Japan	89
<b>American (admixed)</b> CLM - Colombian in Medellin, Colombia	60
MXL - Mexican Ancestry in Los Angeles, CA	66
PUR - Puerto Rican in Puerto Rico	55

Table S2 – Summary Statistics of Population Variation sequence variation data from 10 Western chimpanzees (*Pan troglodytes verus*), generated by the PanMap Project

	<b>N<sup>a</sup></b>	<b>L<sup>b</sup></b>	<b>S<sup>c</sup></b>	<b>NH<sup>d</sup></b>	<b>Hd<sup>e</sup></b>
<b><i>HBB</i></b>		1843	6	8	0.84
<b><i>HBD</i></b>	20	1877	9	3	0.57
<b><i>HBBP1</i></b>		1925	5	5	0.69

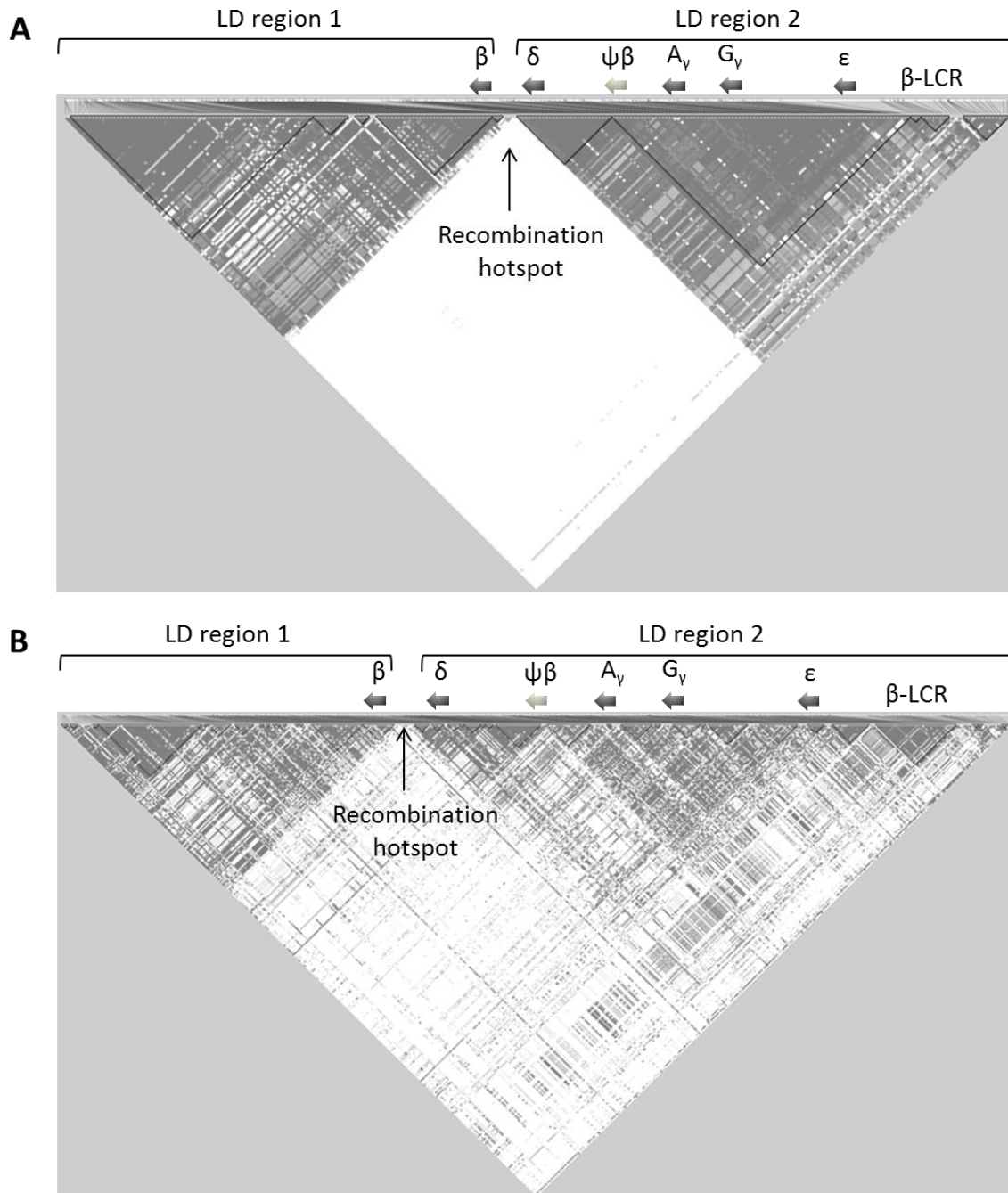
<sup>a</sup> Number of chromosomes

<sup>b</sup> Total number of sites surveyed

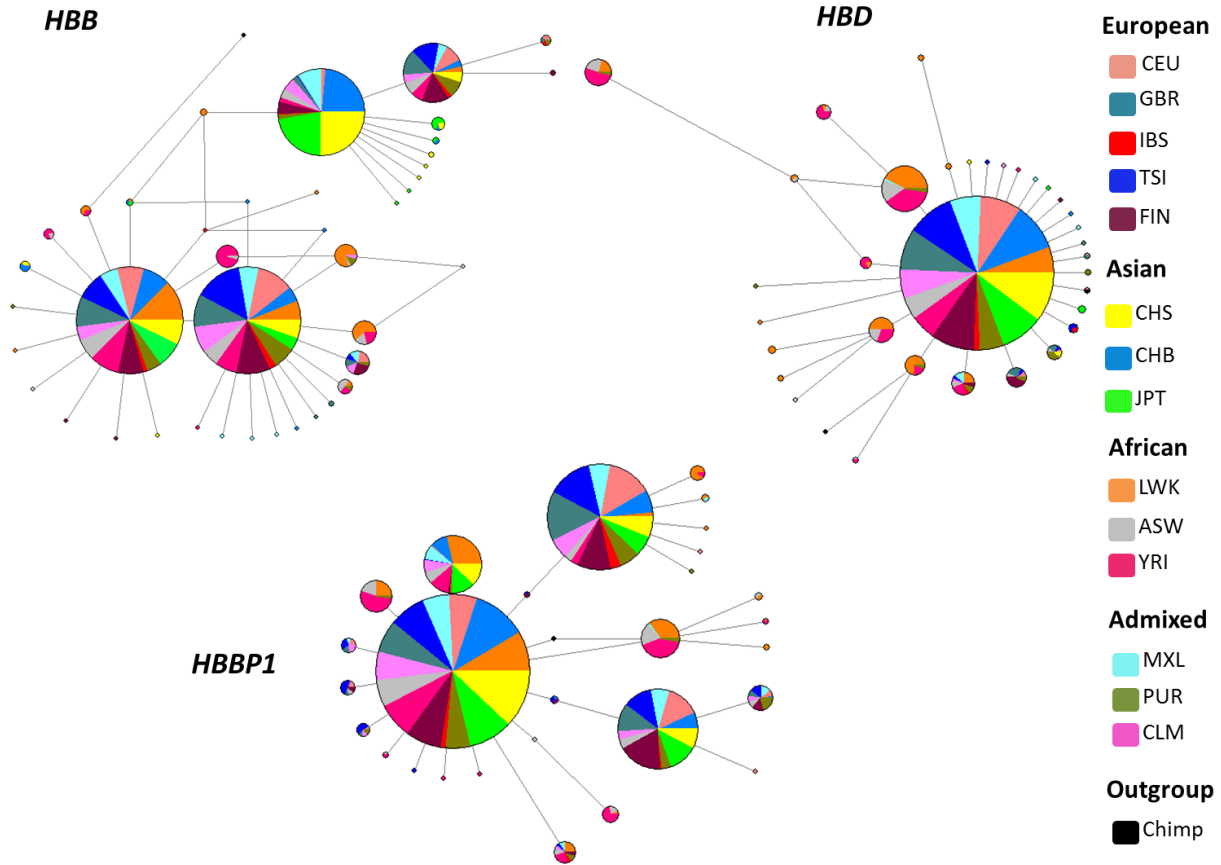
<sup>c</sup> Polymorphic sites

<sup>d</sup> Number of haplotypes

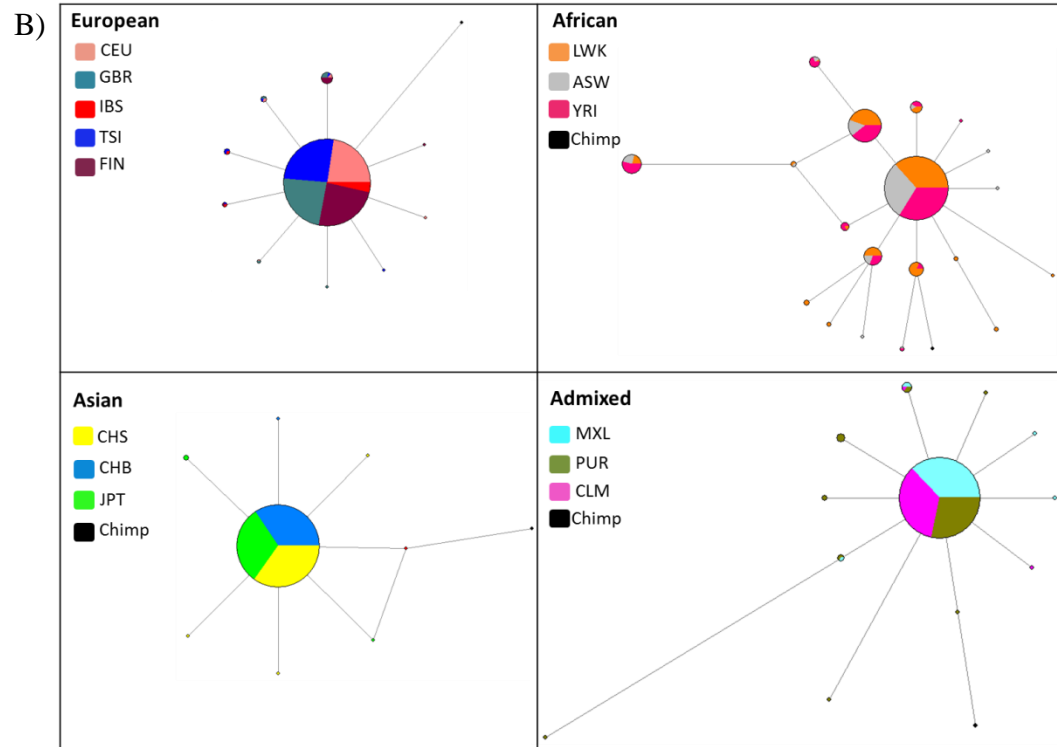
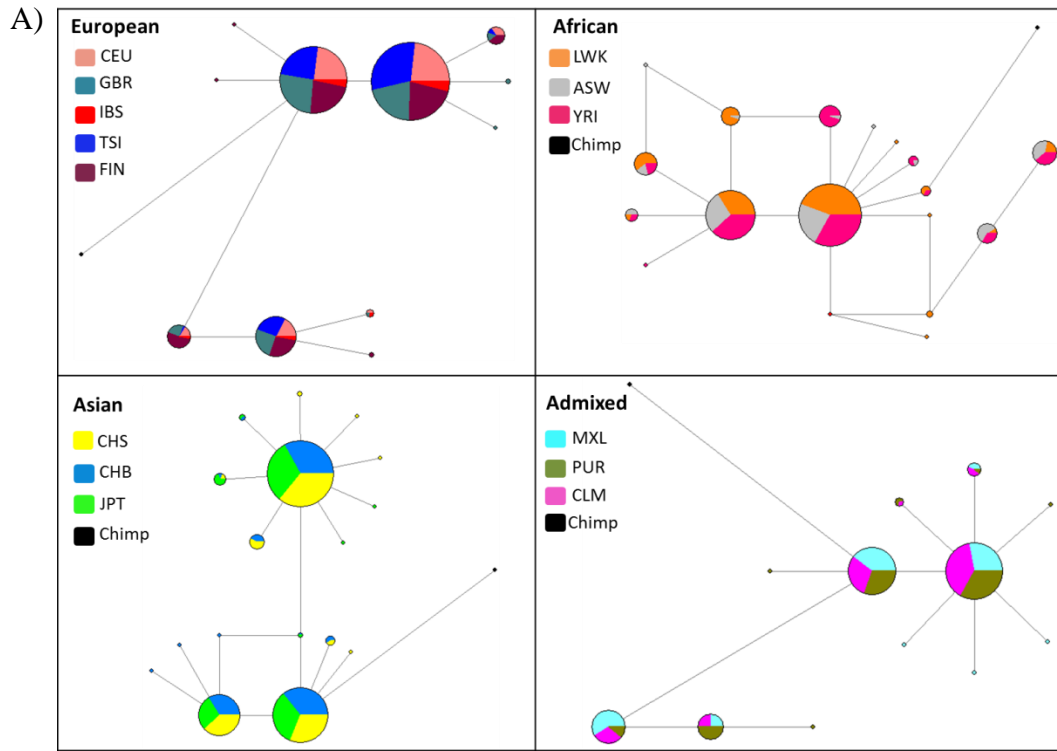
<sup>e</sup> Haplotype diversity

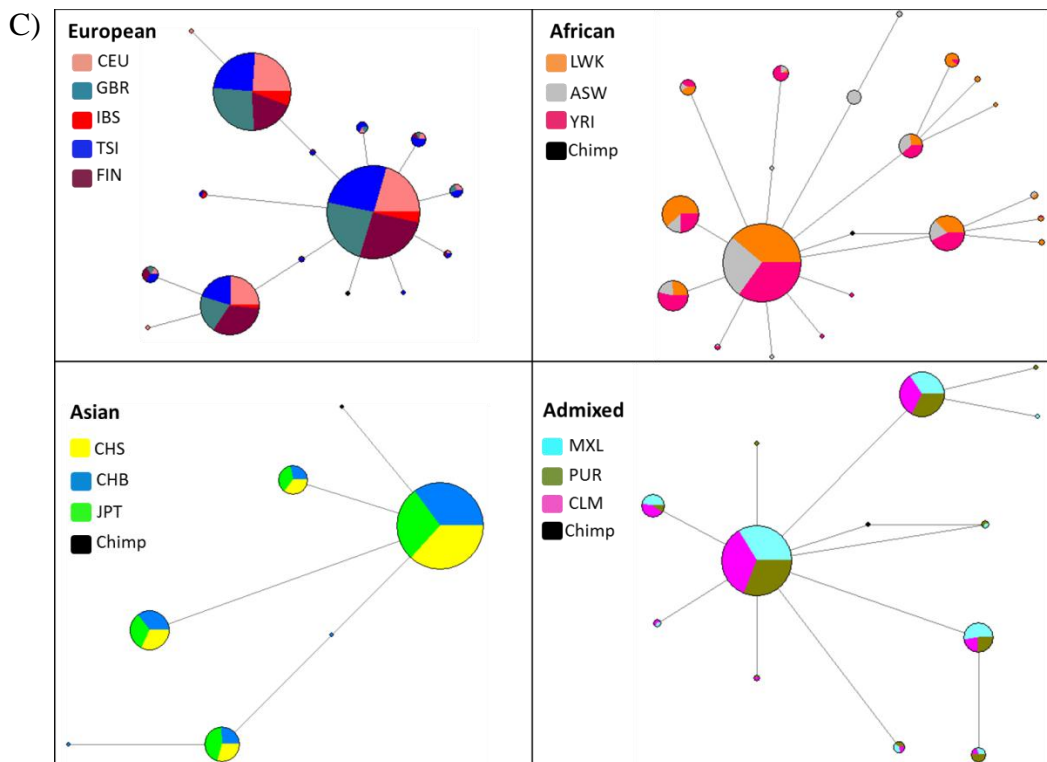


**Figure S1 – LD plot of the  $\beta$ -globin cluster for the data from 1000 Genomes phase 1 release v3 for A) CEU and B) YRI.** The image was constructed using Haploview 4.1 software. The triangles represent LD blocks. Two distinct regions with strong LD are identified: one containing *HBB* (LD region 1) and the other extending from *HBD* to the LCR (LD region 2). In this analysis only variants with a frequency  $\geq 0.5\%$  were considered.



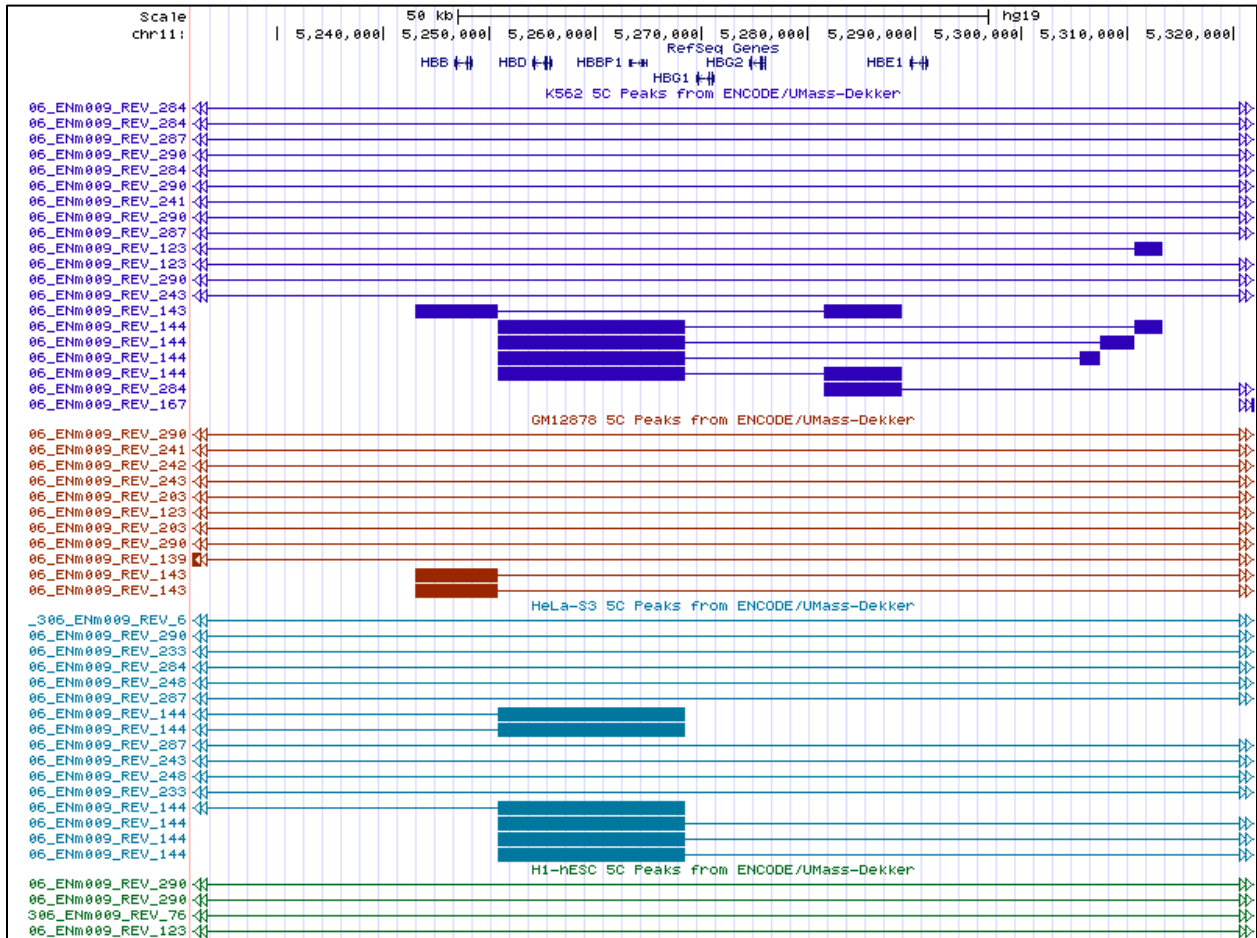
**Figure S2 – Haplotype genealogies for *HBB*, *HBD* and *HBBP1*.** Networks were built using haplotypes from 1000 Genomes Project (14 populations combined). Haplotypes are shown as circles with an area proportional to their frequency; lines connect different haplotypes and the number of mutations is proportional to their length. Full description of population’s acronyms is available in table S1.



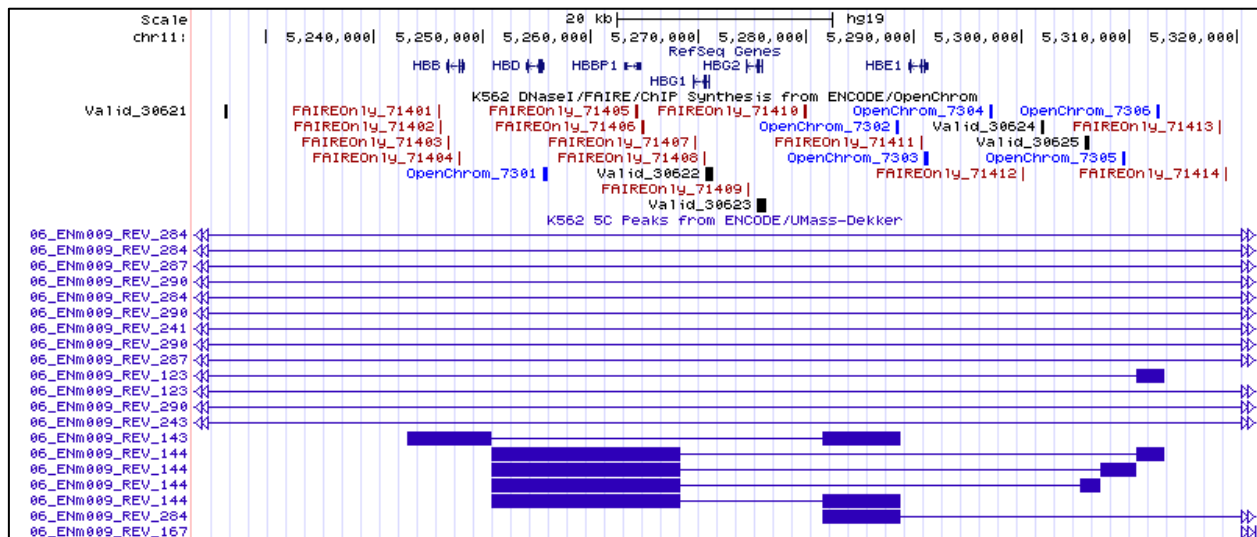


**Figure S3 – Haplotype genealogies for A) *HBB*, B) *HBD* and C) *HBBP1*.** Networks were built using haplotypes from 1000 Genomes Project, representing 1092 individuals from 14 populations: three African (ASW, LWK and YRI), five European (CEU, FIN, GBR, IBS and TSI), three Asian (CHB, CHS and JPT) and 3 American-admixed populations (CLM, MXL and PUR). Haplotypes are shown as circles with an area proportional to their frequency; lines connect different haplotypes and the number of mutations is proportional to their length. Full description of population's acronyms is available in table S1.

A)



B)





**Figure S4 – Chromatin features in the  $\beta$ -globin gene cluster as displayed in the UCSC**

**Genome Browser.** Human RefSeq  $\beta$ -globin genes are labeled and the genomic coordinates of the region displayed are shown (hg19). **A) Chromatin interactions determined by 5C (Chromatin Conformation Capture Carbon Copy).** Each color represents a different cell line (dark blue – K562; red – GM12878; light blue – HeLa-S3; green- H1-hESC). The letters to the right represent the primers used to detect each of the interactions. The regions involved in significant interactions in cis (i.e., from the same ENCODE pilot regions) are represented by blocks and connected by a horizontal line. Interactions displayed were filtered using a z-score interval from 500-1000. **B) Open chromatin regions and/or transcription factor binding sites identified in K562 cells by one or more complementary methodologies: DNaseI hypersensitivity (HS) (Duke DNaseI HS), Formaldehyde-Assisted Isolation of Regulatory Elements (FAIRE) (UNC FAIRE), and chromatin immunoprecipitation (ChIP) for select regulatory factors (UTA TFBS).** Each color represents the assay(s) by which it was detected and its level of validation. Regions that overlap between methodologies identify regulatory elements that are cross-validated indicating high confidence regions (black). In addition, multiple lines of evidence suggest that regions detected by a single assay (e.g., DNase-only or FAIRE-only) are also biologically relevant: high significance regions that indicate open chromatin (blue); low significance regions (red).