# Supplementary information for: Signals of recent positive selection in a worldwide sample of human populations

Joseph K. Pickrell[1,†], Graham Coop[1,12,†], John Novembre[1,2],
Sridhar Kudaravalli[1], Jun Li[10], Devin Absher[4], Balaji S. Srinivasan[6,7,8,9],
Gregory S. Barsh[3], Richard M. Myers[4], Marcus W. Feldman[5],
and Jonathan K. Pritchard[1,11,†]

[1] Department of Human Genetics, The University of Chicago.
[2] Department of Ecology and Evolutionary Biology, UCLA.
[3] Department of Genetics, Stanford University.
[4] HudsonAlpha Institute for Biotechnology, 601 Genome Way, Huntsville, AL.
[5] Department of Biological Sciences, Stanford University.
[6] Stanford Genome Technology Center, Stanford University
[7] Program in Biomedical Informatics, Stanford University
[8] Dept. of Computer Science, Stanford University
[9] Dept. of Statistics, Stanford University
[10] Dept. of Human Genetics, University of Michigan
[11] Howard Hughes Medical Institute.
[12] Current Address: Department of Ecology and Evolution, University of California, Davis
† To whom correspondence should be addressed: pickrell@uchicago.edu,
gmcoop@ucdavis.edu, pritch@uchicago.edu

January 12, 2009

**Power simulations.** To estimate our power to detect selection in the ascertained Illumina data, we performed simulations under the "cosi" model of human demography [3] with a slight modification– for computational efficiency, the recent increases in population size were dropped. This had little effect on the fit of the model to the HapMap data (results not shown). Other relevant parameters are described in the main text. The selection model is one where a single selected site arises and sweeps up in frequency. Implicitly, this assumes a given selected locus in the human genome has experienced only a single selective sweep in the time from over which our statistics are applicable (approximately 10-80ky); we find this reasonable, and simulations of this kind are a standard method for testing the power of a test for selection [2, 4–7]

To approximate the ascertainment of the Illumina panel, we performed a two-stage ascertainment on all of the simulations; this is described in the main text. In Supplementary Figure 1, we show that this procedure provides a good approximation to the frequency spectrum of the Illumina chip. For calculations of power, in all cases the selected site was excluded from analysis, making the calculations conservative. In Supplementary Figure 2, we show power estimates as a function of derived allele frequency in different demographies, and in Supplementary Figure 3, we show the effect of sample size on power.

About 200 of the simulations have a selected allele that arises on the branch common to both the "European" and "Asian" populations; this allows us to get an upper bound on the expected amount of overlap between the populations for XP-EHH (with a selection coefficient of 1% almost all selected alleles that arose during that branch would have gone to fixation and are not likely to be detected by iHS). In these simulations of truly shared selective sweeps, 76% are detected in both populations.

**Overlap of haplotype-based selection signals between populations.** We calculated iHS and XP-EHH in each individual population, and converted scores to empirical p-values as described in the text. For a signal to be considered overlapping between populations, we required that the iHS or XP-EHH score be in the 1% tail of one population and the 5% tail of the other. In Supplementary Figures 4 and 5, we show the fraction of selection signals that overlap between all pairs of populations. The striking difference between African and non-African populations in XP-EHH is an artifact of the way XP-EHH is calculated–for non-African populations, we used the group of Bantu speakers as a reference, and for African populations we used a group of European chromosomes as the reference. We performed the test this way to most closely follow what was most powerful in the simulations, but this means that the non-overlap between African and non-African populations by XP-EHH is not meaningful.

**Calculation of the CLR statistic and overlap with other signals of selection** The CLR test used in the main text as an alternative to XP-EHH for detecting high-frequency sweeps was originally designed to be applied to data for which SNP ascertainment is constant across the genome and, ideally, can be properly modeled. This is not the case for the HGDP data–tag SNPs were identified in the HapMap database, which has complicated and often unknown ascertainment [1].

2

Because of this, we have used this statistic not as a formal test for significance, but rather to rank regions by the deviance of the allele frequency spectra from the genome-wide average. We then used the CLR statistics as test statistics like XP-EHH and iHS and convert them to empirical p-values after controlling for SNP density, as described in the main text.

As XP-EHH and CLR both detect local deviations in the allele frequency spectrum (XP-EHH detects regions of extended homozygosity), we expect them to largely overlap. This is indeed the case for non-African populations–defining overlap as above, we find that the CLR and XP-EHH statistics overlap 63%, 58%, 67%, 83%, 88%, and 71% for Europe, the Middle East, S. Asian, E. Asia, Oceania, and the Americas, respectively. For the Bantu, however, the overlap is only 28%. We interpret this to mean that many signals of fixed sweeps in the Bantu are less robust than those outside of Africa. This is consistent with the observation that there appear to be few reliable signals of recent fixations in African populations[2, 4].

Figure 1: Marginal frequency spectra in the simulations before and after the application of the ascertainment procedure, as well as in the real data. In the first column are the marginal allele frequency spectra in the simulations before application of the ascertainment procedure, in the second column are the allele frequency spectra in the simulations after the application of the ascertainment procedure, and in the third column are the marginal allele frequency spectra in the HapMap populations at the SNPs typed on the Illumina panel.

Figure 2: **Power of iHS (A) and XP-EHH (B) in three demographies based on the HapMap.** Selected alleles were introduced at a random time with a selection coefficient of 1% and simulated forwards in time. In all simulations, the selected site was excluded from analyses. Simulations were binned into six bins according the the final frequency of the selected allele. The type I error was set to 1%, with the threshold determined by neutral simulations.

Figure 3: **Power of iHS (A) and XP-EHH (B) in the YRI demography for different sample sizes (in number of chromosomes)**. Selected alleles were introduced at a random time with a selection coefficient of 1% and simulated forwards in time. Simulations were binned into six bins according the the final frequency of the selected allele. The type I error was set to 1%, with the threshold determined by neutral simulations.

Figure 4: Overlap of iHS signals across populations. In each cell $(i, j)$ the fraction of windows with scores in the 1% tail of iHS in population $i$ and in the 5% tail of iHS in population $j$ is colored according the legend on the right.

Figure 5: Overlap of XP-EHH signals across populations. In each cell $(i, j)$ the fraction of windows with scores in the 1% tail of XP-EHH in population $i$ and in the 5% tail of XP-EHH in population $j$ is colored according the legend on the right.

Figure 6: **Top 10 iHS (A) and XP-EHH (B) signals by population cluster**. This is analogous to Figure 1 in the main text, but uses a sliding window of 200kb that slides in steps of 100kb.

Figure 7: **Top 10 iHS (A) and XP-EHH (B) signals by population cluster**. This is analogous to Figure 1 in the main text, but uses a window size of 100kb.

Figure 8: **Top 10 CLR signals by population cluster**. This is analogous to Figure 1 in the main text.

Figure 9: **F$_{ST}$ surrounding loci involved in natural variation in lipid levels**. This figure was generated as in Figures 3 and 4 in the main text.

Figure 10: **$F_{ST}$ surrounding loci involved in natural variation in susceptibility to Crohn's disease**. This figure was generated as in Figures 3 and 4 in the main text.

13

Figure 11: **F$_{ST}$ surrounding loci involved in natural variation in height**. This figure was generated as in Figures 3 and 4 in the main text.

Figure 12: $\mathbf{F_{ST}}$ **between Bantu and Pygmy populations surrounding loci involved in natural variation in height**. This figure was generated as in Figures 3 and 4 in the main text.

Figure 13: The top 1% of iHS signals in the Bantu, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

0.001

0.05

BiakaPygmy Bantu Europe Mideast S.Asia E.Asia Oceania America

| Position | Genes |
|---|---|
| chr9 31.2:31.6 | |
| chr7 66.6:67 | COX5B,ACTR1B,ZAP70,KIAA1641 |
| chr2 97.6:97.8 | MAN2A1 |
| chr5 109:109.2 | HIST1H4G,HIST1H3F,HIST1H2BH,HIST1H3G,HIST1H2BI,HIST1H4H,BTN3A2,BTN2A2,BTN3A1,BTN2A3,BTN3A3,BT |
| chr6 26.4:26.6 | RBMS3 |
| chr3 29.4:29.8 | OR10Q1,OR10W1,OR5B17,OR5B3,OR5B2,OR5B12,OR5B21 |
| chr11 57.8:58 | CPNE4 |
| chr3 133.2:133.4 | TBX18 |
| chr6 85.4:85.6 | DCAMKL1,SOHLH2 |
| chr13 35.4:35.6 | AFF3 |
| chr2 99.6:99.8 | SLC28A3 |
| chr9 86:86.2 | |
| chr9 75.2:75.4 | PRTG,NEDD4 |
| chr15 53.8:54 | |
| chr3 75:75.2 | SUHW4,TCF12 |
| chr15 54.8:55 | |
| chr9 13.6:13.8 | TBC1D15,TPH2 |
| chr12 70.6:70.8 | TCF12 |
| chr15 55.2:55.4 | RTN4,FLJ42562 |
| chr2 55:55.2 | |
| chr4 160.8:161.2 | |
| chr5 130:130.2 | SGEF |
| chr3 155.2:155.4 | |
| chr1 5:5.2 | KIAA0350 |
| chr16 11:11.2 | ITGA6,PDK1 |
| chr2 173:173.2 | CREB5 |
| chr7 28.6:28.8 | |
| chr19 36.6:36.8 | NSUN6,ARL5B |
| chr10 19:19.2 | NEK10,SLC4A7 |
| chr3 27.2:27.4 | SCRN1,FKBP14,PLEKHA8,WIPF3 |
| chr7 29.8:30 | C7orf10 |
| chr7 40.6:40.8 | CSNK1G3 |
| chr5 123:123.2 | SETBP1 |
| chr18 40.8:41 | PHACTR3 |
| chr20 57.6:57.8 | THADA,PLEKHH2 |
| chr2 43.6:43.8 | CTNNA3 |
| chr10 67.4:67.6 | FLJ37357 |
| chr2 84.6:84.8 | |
| chr7 84.2:84.4 | SOX6 |
| chr11 16.4:16.6 | GDPD4,PAK1,MYO7A |
| chr11 76.6:76.8 | RIOK3,C18orf8,NPC1,ANKRD29,C18orf45 |
| chr18 19.2:19.4 | C3orf48,EFHB,RAB5A,PCAF,SGOL1 |
| chr3 20:20.2 | PFKP,PITRM1 |
| chr10 3:3.2 | |
| chr1 189.8:190 | |
| chr13 57.4:57.6 | TRIO |
| chr5 14.2:14.4 | |
| chr7 9.4:9.6 | |
| chr5 26.6:26.8 | EML4,COX7A2L |
| chr2 42.2:42.4 | DNER,PID1 |
| chr2 229.8:230 | SPOCK1 |
| chr5 136.2:136.4 | EPAS1,ATP6V1E2,RHOQ |
| chr2 46.4:46.6 | LOC388965 |
| chr2 84.2:84.4 | INPP4A,MGC26733,CNGA3 |
| chr2 98.2:98.4 | |
| chr6 23.8:24 | |
| chr6 115.8:116 | OPRM1 |
| chr6 154.2:154.4 | |
| chr7 125.4:125.6 | DKFZP586P0123,PPME1,P4HA3,PGM2L1,KCNE3 |
| chr11 73.6:73.8 | B2M,TRIM69,C15orf43 |
| chr15 42.8:43 | NRXN1 |
| chr2 50.2:50.4 | |
| chr15 51.2:51.4 | SOD2,WTAP,ACAT2,TCP1,MRPL18,PNLDC1,MAS1 |
| chr6 159.8:160.2 | |
| chr6 19:19.2 | GNA12,CARD11 |
| chr7 2.8:3 | TTC23,LRRC28,DMN |
| chr15 97.4:97.6 | MAP2 |
| chr2 209.8:210 | PCDHGA1,PCDHGA2,PCDHGA3,PCDHGB1,PCDHGA4,PCDHGB2,PCDHGA5,PCDHGB3,PCDHGA6,PCDHGA7,PCD |
| chr5 140.8:141 | PDCD7,CLPX,CILP,PARP16,PUNC |
| chr15 63.2:63.4 | LRRN1 |
| chr3 76.8:77 | FER |
| chr3 3.6:3.8 | EPB41L2 |
| chr5 108.4:108.6 | |
| chr6 131.2:131.4 | |
| chr5 18.2:18.4 | |
| chr13 90.4:90.6 | |
| chr3 63:63.2 | TGFBI,SMAD5,TRPC7 |
| chr20 37.8:38 | FUT8 |
| chr5 135.4:135.6 | C14orf145,TSHR |
| chr14 64.8:65 | TSCOT,ZFP37 |
| chr14 80.4:80.6 | |
| chr9 114.6:114.8 | POMC,DNMT3A |
| chrX 145.6:145.8 | SYN3,TIMP3 |
| chr2 25.2:25.4 | C8orf34 |
| chr22 31.6:31.8 | NPHP4,KCNAB2 |
| chr8 69.6:69.8 | |
| chr1 5.8:6 | DSCAM |
| chr11 38.4:38.6 | EPB41L4B,C9orf4,C9orf5,CTNNAL1 |
| chr11 104.6:104.8 | SCLT1,C4orf33 |
| chr1 79:79.2 | |
| chr21 40.8:41 | CSS3 |
| chr9 110.8:111 | COL21A1 |
| chr4 130.2:130.4 | |
| chr5 160.2:160.4 | GMDS |
| chr16 78.6:78.8 | TAAR2,TAAR5,TAAR1,VNN1,VNN3,VNN2,C6orf192,RPS12 |
| chr5 129.2:129.4 | |
| chr6 56:56.2 | |
| chr18 60.2:60.4 | |
| chr6 2.2:2.4 | |
| chr6 133:133.2 | |
| chr10 119.4:119.6 | |

Figure 14: The top 1% of iHS signals in the Biaka Pygmies, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

17

Figure 15: The top 1% of iHS signals in the Europeans, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 16: The top 1% of iHS signals in the Middle East, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 17: The top 1% of iHS signals in South Asia, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 18: The top 1% of iHS signals in East Asia, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 19: The top 1% of iHS signals in the Americas, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 20: The top 1% of iHS signals in Oceania, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

0.001

0.05

| | BiakaPygmy | Bantu | Europe | Mideast | S.Asia | E.Asia | Oceania | America | |
|---|---|---|---|---|---|---|---|---|---|
chr13 24.2:24.4 | | | | | | | | | ATP12A,CENPJ,RNF17
chr20 34.8:35 | | | | | | | | | NDRG3,DSN1,C20orf117,C20orf118,SAMHD1
chr17 0.6:0.8 | | | | | | | | | VPS53,FAM57A,C17orf25,RNMTL1,NXN,TIMM22,GEMIN4
chr6 164.6:164.8 | | | | | | | | | CEACAM6,LYPD4,RPS19,ARHGEF1,RABAC1,ATP1A3,GRIK5,DMRTC2,CEACAM3,CD79A
chr19 47:47.2 | | | | | | | | | PKHD1,IL17A,IL17F,MCM3
chr6 52:52.2 | | | | | | | | | NDUFB4,GTF2E1,HGD,RABL3
chr3 121.8:122 | | | | | | | | |
chr14 84.6:84.8 | | | | | | | | | SPTY2D1,TMEM86A,IGSF22,PTPN5,UEVLD
chr11 18.6:18.8 | | | | | | | | | MAP2K3,KCNJ12
chr17 21.2:21.6 | | | | | | | | | HRASLS5,LGALS12,RARRES3,HRASLS2,HRASLS3,DKFZP564J0863,RTN3
chr11 63:63.2 | | | | | | | | | DIAPH3
chr13 59.6:59.8 | | | | | | | | |
chr1 207:207.2 | | | | | | | | |
chr2 59:59.2 | | | | | | | | | HRASLS
chr3 194.2:194.4 | | | | | | | | | PLD3,HIPK4,SERTAD1,SERTAD3,BLVRB,SHKBP1,PRX,SPTBN4,LTBP4
chr19 45.6:45.8 | | | | | | | | | BTBD11
chr12 106.2:106.4 | | | | | | | | | DYM
chr18 44.8:45 | | | | | | | | | CYB561,ACE,KCNH6,CCDC44,WDR68
chr17 58.8:59 | | | | | | | | |
chr10 92:92.2 | | | | | | | | |
chr13 58.6:58.8 | | | | | | | | | LRRC21,LRRC22,PCDH21,RGR,KIAA1128
chr10 86:86.2 | | | | | | | | | RASGEF1A,RET,GALNACT-2
chr10 42.8:43 | | | | | | | | | ZIC4,ZIC1
chr3 148.4:148.6 | | | | | | | | | HMG2L1,TOM1,HMOX1,MCM5
chr22 34:34.2 | | | | | | | | | CAMK1D
chr10 12.4:12.6 | | | | | | | | | NGFB
chr1 115.6:115.8 | | | | | | | | | SNF1LK,FLJ41733,C21orf125,HSF2BP
chr21 43.6:43.8 | | | | | | | | | IKZF4,WIBG,DGKA,SILV,CDK2,RAB5B,SUOX,RPS26,ERBB3,RPL41,ZC3H10,FAM62A,MYL6B,MYL6,SMARCC2,PA2
chr12 54.6:54.8 | | | | | | | | | TULP4
chr6 72.4:72.6 | | | | | | | | | GPR45,TGFBRAP1,C2orf49,FHL2
chr6 158.6:158.8 | | | | | | | | | DLX6,DLX5,ACN9
chr2 105.2:105.4 | | | | | | | | | DVL3,AP2M1,ABCF3,ALG3,MGC2408,ECE2,CAMK2N2,PSMD2,EIF4G1,C3orf40,CLCN2,POLR2H,THPO,CHRD
chr7 96.4:96.6 | | | | | | | | | TTC7B
chr3 185.4:185.6 | | | | | | | | | DPP6
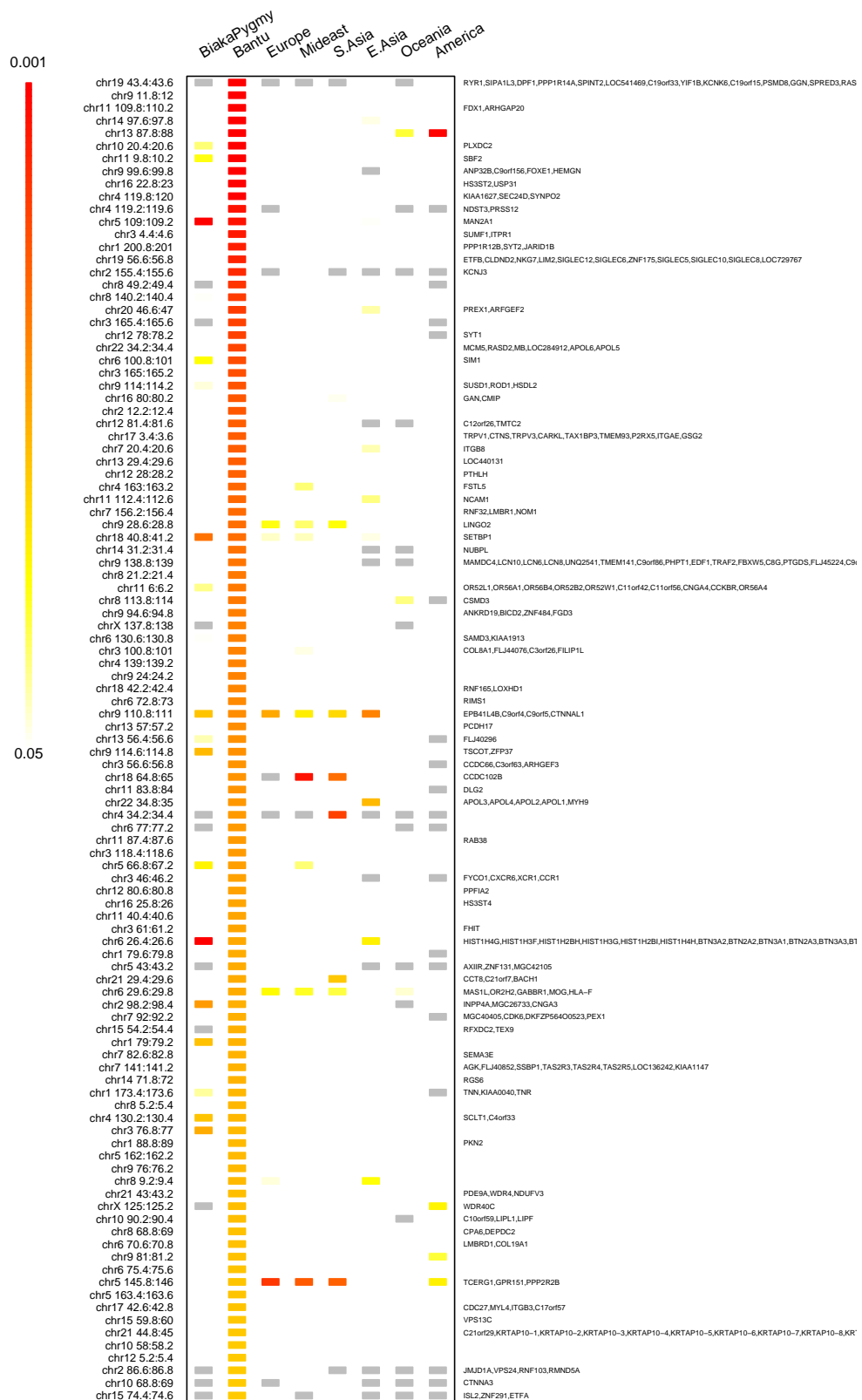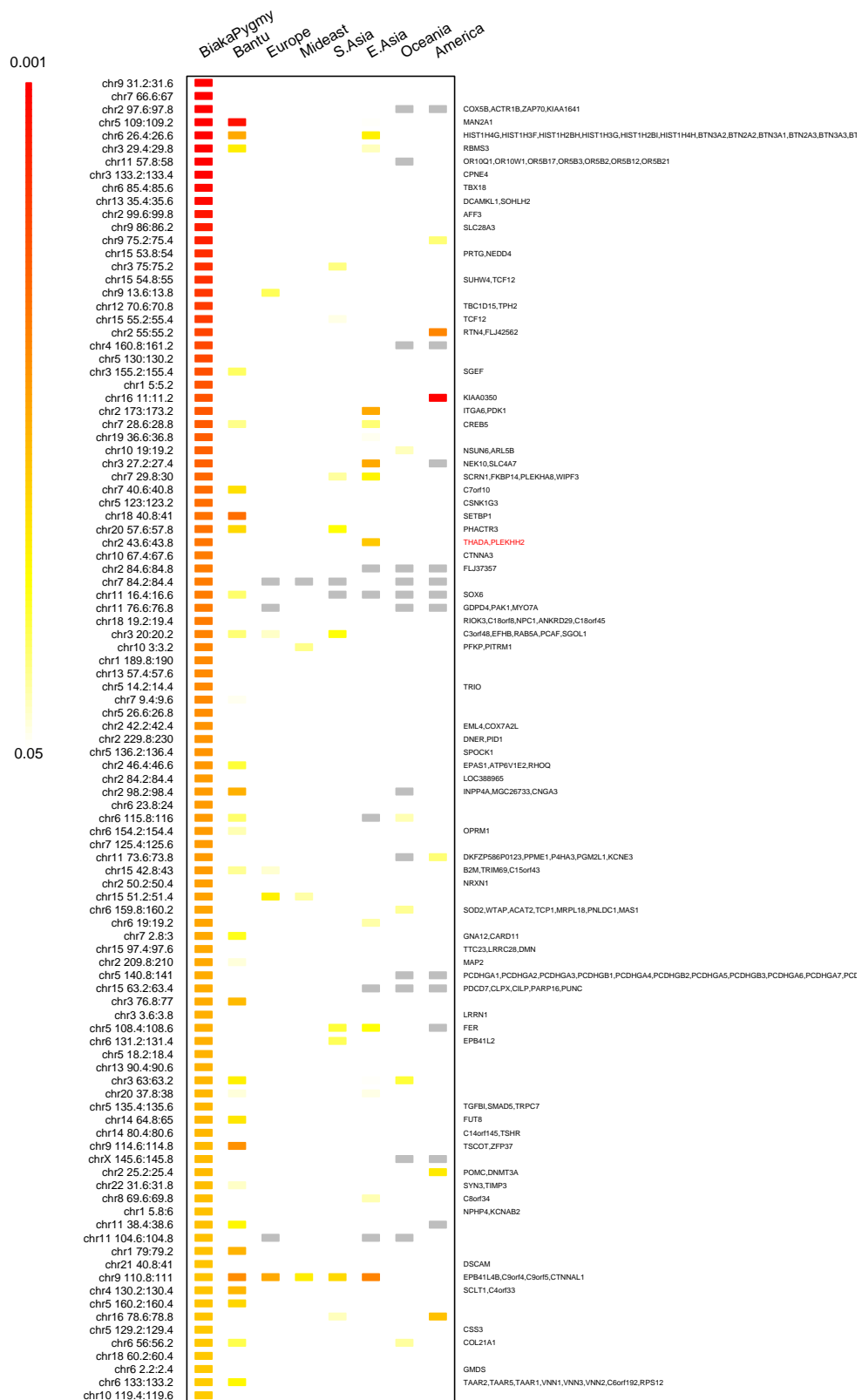chr14 90:90.2 | | | | | | | | | KLHDC4,SLC7A5,CA5A,BANP
chr7 153.4:153.6 | | | | | | | | | CASK,GPR34,GPR82
chr16 86.4:86.6 | | | | | | | | | DMD
chrX 41.4:41.8 | | | | | | | | | VPS41,POU6F2
chrX 33.2:33.4 | | | | | | | | | C20orf23
chr7 38.8:39 | | | | | | | | | ZFPM2
chr20 16.4:16.6 | | | | | | | | | RAMP3
chr8 106.4:106.6 | | | | | | | | | AKAP1,COIL,SCPEP1
chr7 45.2:45.4 | | | | | | | | | TTC17
chr17 52.4:52.6 | | | | | | | | | KIAA1303,RNF213,FLJ35220,NPTX1
chr11 43.4:43.6 | | | | | | | | | CSMD1
chr17 76:76.2 | | | | | | | | | PSG5,PSG4,PSG9,TEX101,CD177
chr8 4.8:5 | | | | | | | | | FLJ37543
chr19 48.4:48.6 | | | | | | | | | MAGI3,PHTF1,RSBN1,PTPN22,C1orf178,AP4B1,DCLRE1B
chr5 60.8:61 | | | | | | | | | KIAA0895,ANLN,AOAH
chr1 113.8:114.2 | | | | | | | | | RHBDL2,C1orf108,NDUFS5,MACF1
chr7 36.4:36.6 | | | | | | | | | NRXN1
chr1 39.2:39.4 | | | | | | | | | PTHLH
chr2 49.8:50 | | | | | | | | | PRKAR2B,HBP1,COG5
chr12 28:28.2 | | | | | | | | | RALGPS2,FAM20B,ABL2,TOR3A
chr7 106.4:106.8 | | | | | | | | | ADRA1A,DPYSL2
chr1 177.2:177.4 | | | | | | | | | NRXN1
chr8 26.6:26.8 | | | | | | | | | CYP7A1,SDCBP,NSMAF
chr2 50.8:51 | | | | | | | | | LARS2,SACM1L,SLC6A20,LZTFL1,LIMD1
chr8 59.6:59.8 | | | | | | | | | TFAP4,Magmas,CORO7,VASN,DNAJA3,GLIS2
chr3 45.6:45.8 | | | | | | | | | C9orf9,TSC1,GFI1B,GTF3C5,CEL,RALGDS,GBGT1
chr16 4.2:4.4 | | | | | | | | | TRPV1,CTNS,TRPV3,CARKL,TAX1BP3,TMEM93,P2RX5,ITGAE,GSG2
chr9 134.8:135 | | | | | | | | | MLSTD2
chr17 3.4:3.6 | | | | | | | | | HCN2,POLRMT,FGF22,FLJ45684,RNF126,PRSSL1,C19orf21,PTBP1,PRG2,AZU1,PRTN3,ELA2,CFD,THRAP5,C19or
chr11 13.6:13.8 | | | | | | | | |
chr19 0.6:0.8 | | | | | | | | | IKIP,APAF1,ANKS1B
chr1 82.4:82.8 | | | | | | | | | PTPRS,ZNRF4
chr12 97.6:97.8 | | | | | | | | | OPCML
chr19 5.2:5.4 | | | | | | | | | ITGA8,C10orf97
chr11 132.4:132.6 | | | | | | | | |
chr10 15.8:16.2 | | | | | | | | | SETBP1
chr2 62.2:62.4 | | | | | | | | | HIST1H4G,HIST1H3F,HIST1H2BH,HIST1H3G,HIST1H2BI,HIST1H4H,BTN3A2,BTN2A2,BTN3A1,BTN2A3,BTN3A3,BT
chr18 40.8:41 | | | | | | | | | TRAPPC6A,EXOC3L2,MARK4,CKM,ERCC2,PPP1R13L,CD3EAP,ERCC1,BLOC1S3,KLC3
chr6 26.4:26.6 | | | | | | | | | PHCA,CAPN5,OMP,GDPD4,B3GNT6,MYO7A
chr19 50.4:50.6 | | | | | | | | |
chr11 76.4:76.6 | | | | | | | | | C14orf70,DLK1
chr9 25:25.2 | | | | | | | | | FAM126A,KLHL7,NUPL2
chr14 100.2:100.4 | | | | | | | | | LEF1,HADH
chr7 23:23.2 | | | | | | | | | ARHGEF9
chr4 109.2:109.4 | | | | | | | | |
chrX 62.8:63 | | | | | | | | | VCX2
chr20 37.8:38 | | | | | | | | |
chrX 8:8.2 | | | | | | | | | LOC152485
chr4 180.4:180.6 | | | | | | | | | IKZF2
chr4 147:147.2 | | | | | | | | | MN1,PITPNB
chr2 213.6:213.8 | | | | | | | | | NFKB2,ELOVL3,PITX3,GBF1,PSD,FBXL15,CUEDC2,C10orf95,TMEM180,ACTR1A
chr22 26.4:26.6 | | | | | | | | | SDHC,C1orf192,FCGR2A,HSPA6,FCGR3A,FCGR2C
chr10 104:104.2 | | | | | | | | | FBXO2,FBXO44,FBXO6,MAD2L2,C1orf187,MTHFR,CLCN6,NPPA,NPPB,AGTRAP
chr1 159.6:159.8 | | | | | | | | | PLEKHG1,MTHFD1L
chr1 11.6:11.8 | | | | | | | | | LIPA,IFIT2,IFIT3,IFIT1L,IFIT5,SLC16A12,IFIT1,CH25H
chr6 151:151.2 | | | | | | | | | PHACTR2,LTV1,PLAGL1
chr10 91:91.2 | | | | | | | | | GCK,YKT6,CAMK2B,NPC1L1,DDX56,TMED4,OGDH,NUDCD3
chr6 144.2:144.4 | | | | | | | | | C17orf80,CDC42EP4,SDK2
chr14 100.8:101 | | | | | | | | | MME,FLJ46120
chr7 44.2:44.6 | | | | | | | | | KCNA4
chr17 68.8:69 | | | | | | | | | GALNTL4
chr3 156.2:156.4 | | | | | | | | | HEATR2,C7orf20,CENTA1,CYP2W1,MGC11257,FLJ44124,UNC84A,COX19
chr11 29.8:30 | | | | | | | | |
chr11 11.4:11.6 | | | | | | | | | CDH19
chr7 0.8:1 | | | | | | | | | FOS,JDP2
chr11 115.8:116 | | | | | | | | | ITSN1,ATP5O
chr11 121.6:121.8 | | | | | | | | |
chr18 62.4:62.6 | | | | | | | | | C6orf118
chr14 74.8:75 | | | | | | | | | ACOT6,DNAL1,PNMA1,C14orf43,ZADH1,ZNF410
chr21 34:34.2 | | | | | | | | |
chr10 9.6:9.8 | | | | | | | | | LMO3
chr6 165.4:165.6 | | | | | | | | | SNX13
chr14 73.2:73.4 | | | | | | | | | SNAPC4,GPSM1,CARD9,SDCCAG3,PMPCA,INPP5E,C9orf163,NOTCH1,LOC728489,KIAA0310
chr5 91.8:92 | | | | | | | | | LOC220686,LOC375133,UBE2L3,LOC150223,CCDC116,SDF2L1,PPIL2,YPEL1,MAPK1
chr12 16.6:16.8 | | | | | | | | | C11orf11,C11orf9,C11orf10,FEN1,FADS1,FADS2,FADS3,RAB3IL1
chr17 6:17.8 | | | | | | | | | ANP32B,C9orf156,FOXE1,HEMGN
chr9 138.4:138.6 | | | | | | | | | CDC27,MYL4,ITGB3,C17orf57
chr22 20.2:20.4 | | | | | | | | | IFNK,MOBKL2B,C9orf72
chr11 61.2:61.4 | | | | | | | | | ABTB2,CAT
chr9 99.6:99.8 | | | | | | | | | CTPS,SLFNL1,SCMH1
chr17 42.6:42.8 | | | | | | | | | SLC35F2,RAB39,CUL5
chr9 27.4:27.6 | | | | | | | | | C7orf44,HECW1,STK17A
chr11 34.2:34.4 | | | | | | | | | ABCG8,LRPPRC,PPM1B
chr1 41.2:41.4 | | | | | | | | | KIAA1920
chr11 107.2:107.4 | | | | | | | | |
chr7 43.4:43.6 | | | | | | | | |
chr2 44:44.2 | | | | | | | | |
chr15 82.6:82.8 | | | | | | | | |
chr8 135.2:135.4 | | | | | | | | |

Figure 21: The top 1% of XP-EHH signals in the Bantu, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 22: The top 1% of XP-EHH signals in the Biaka Pygmies, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.
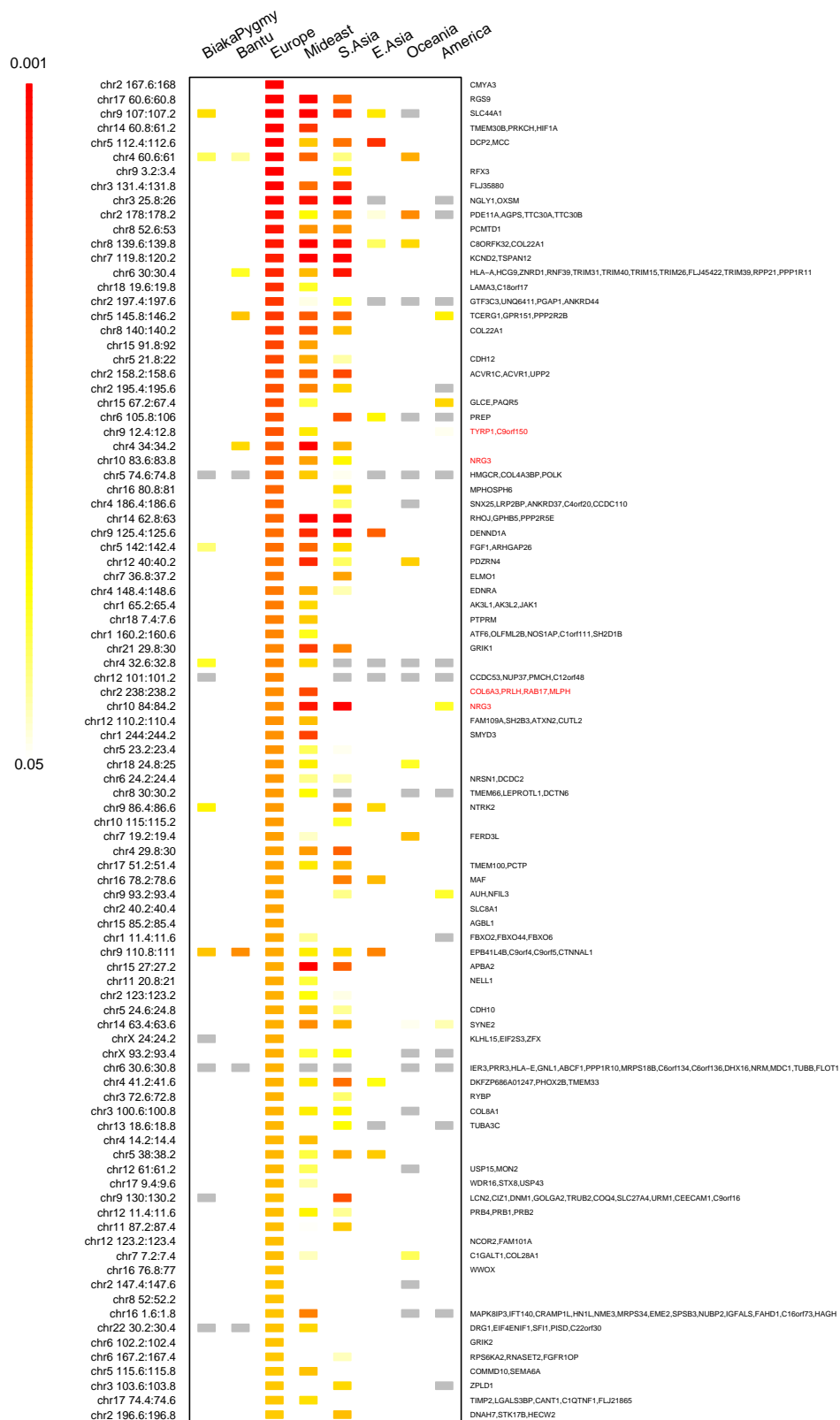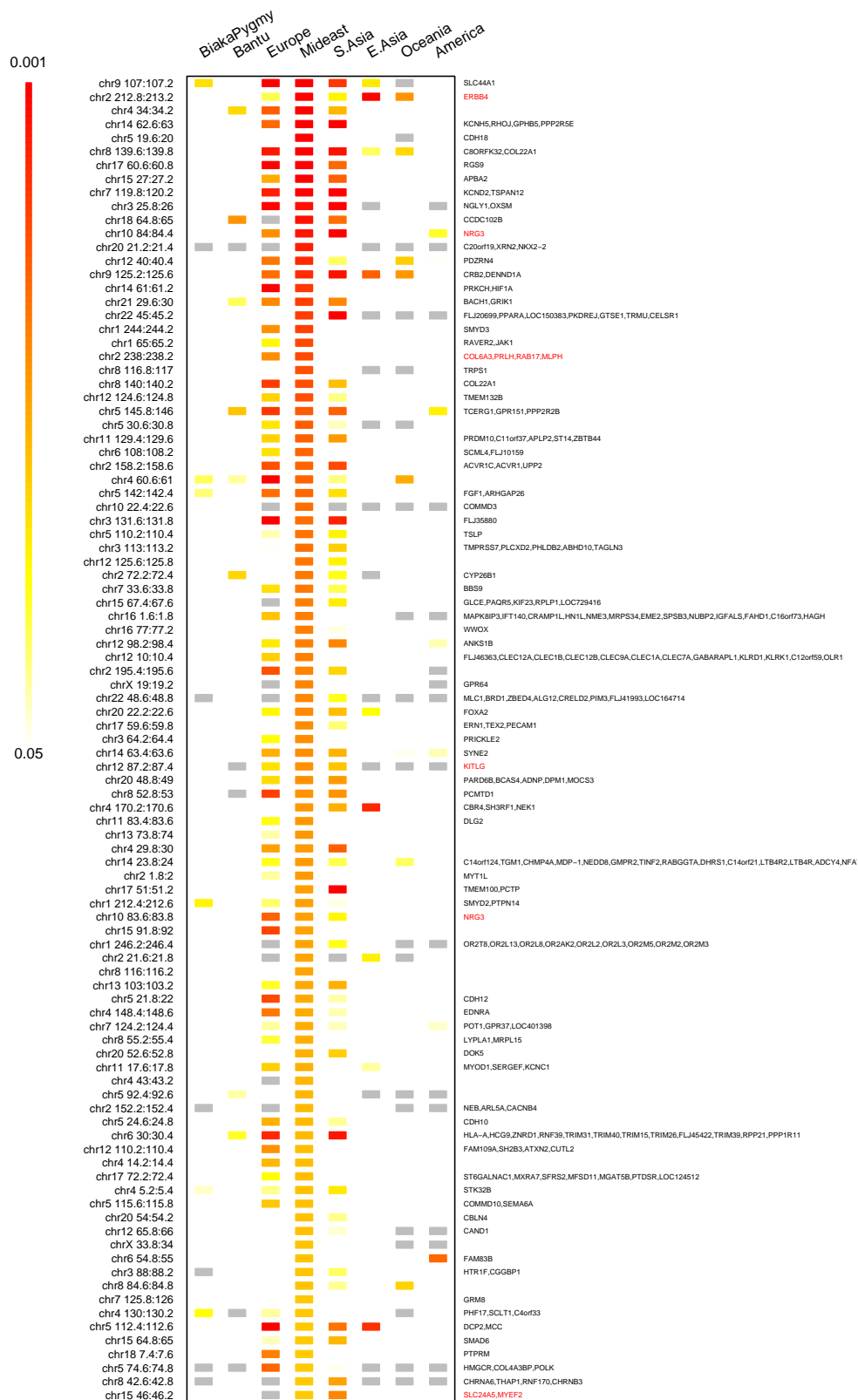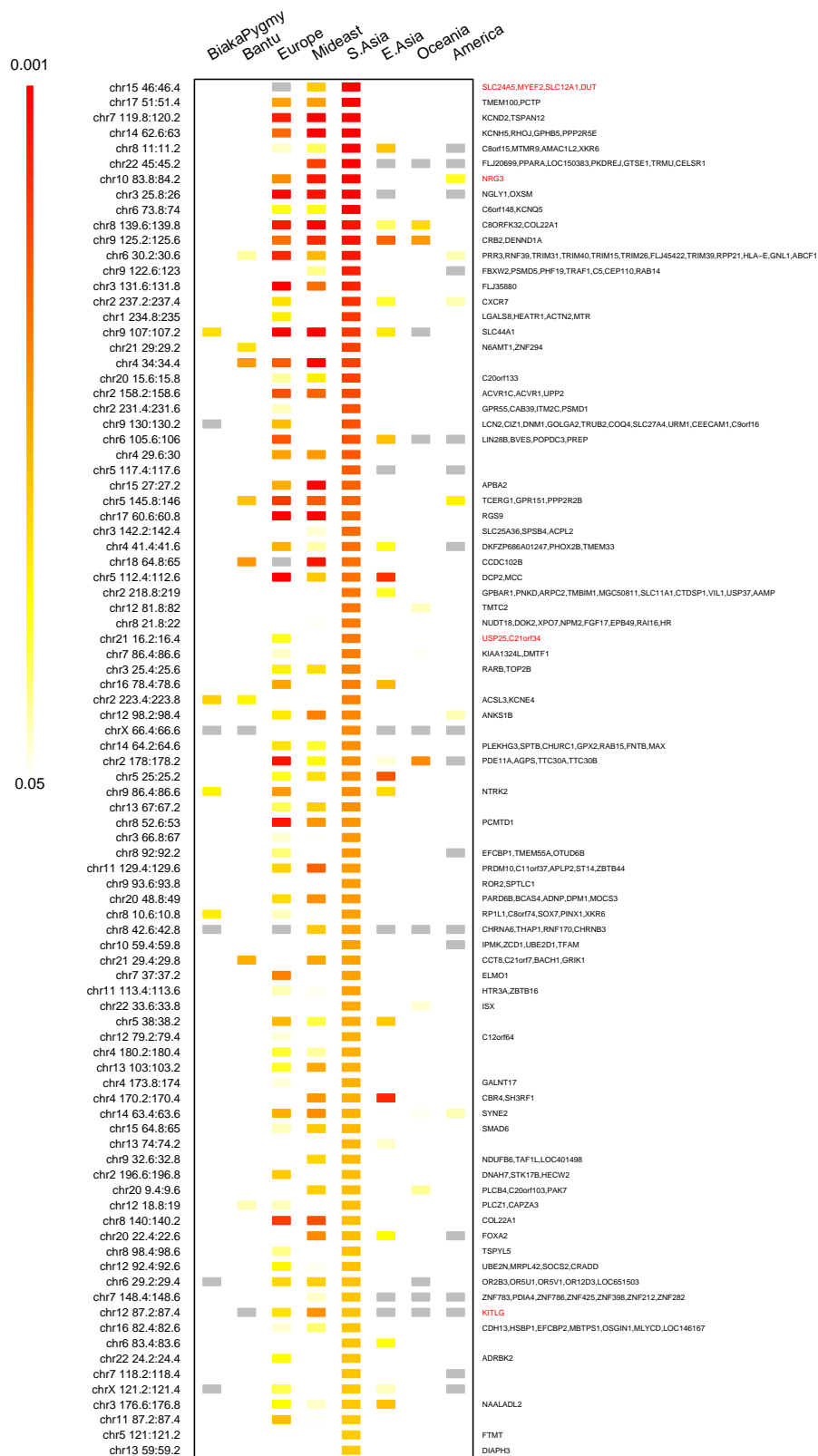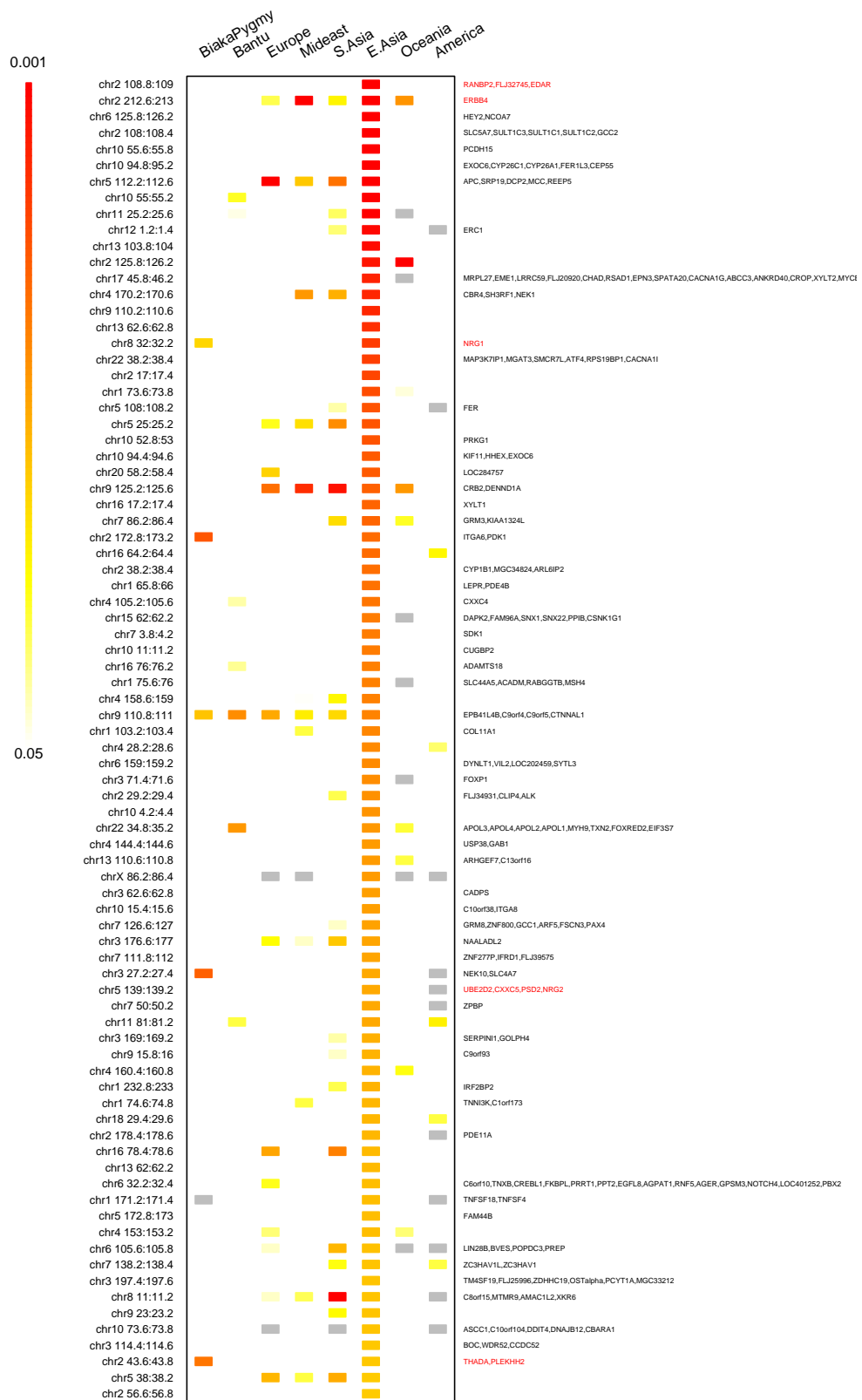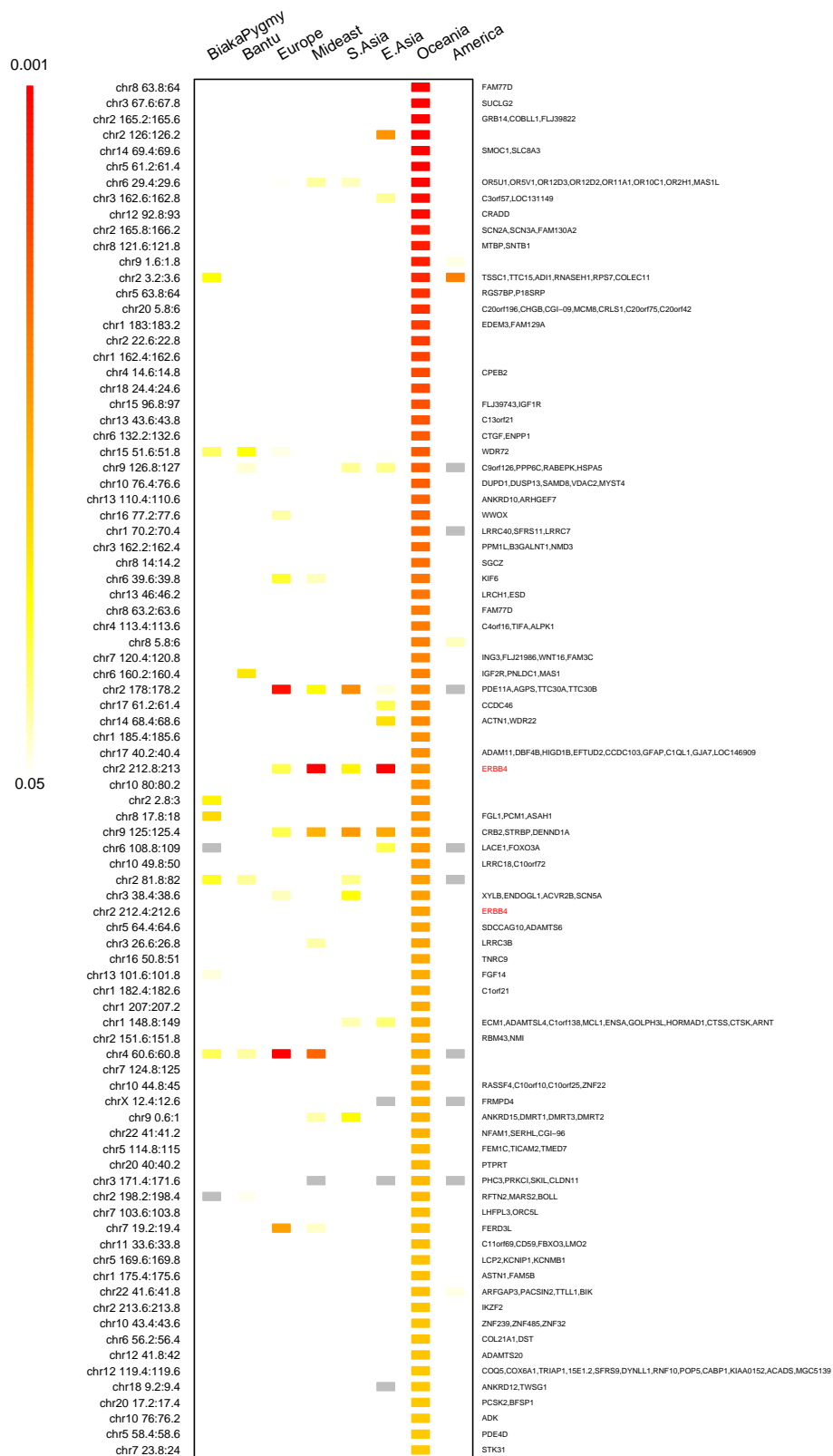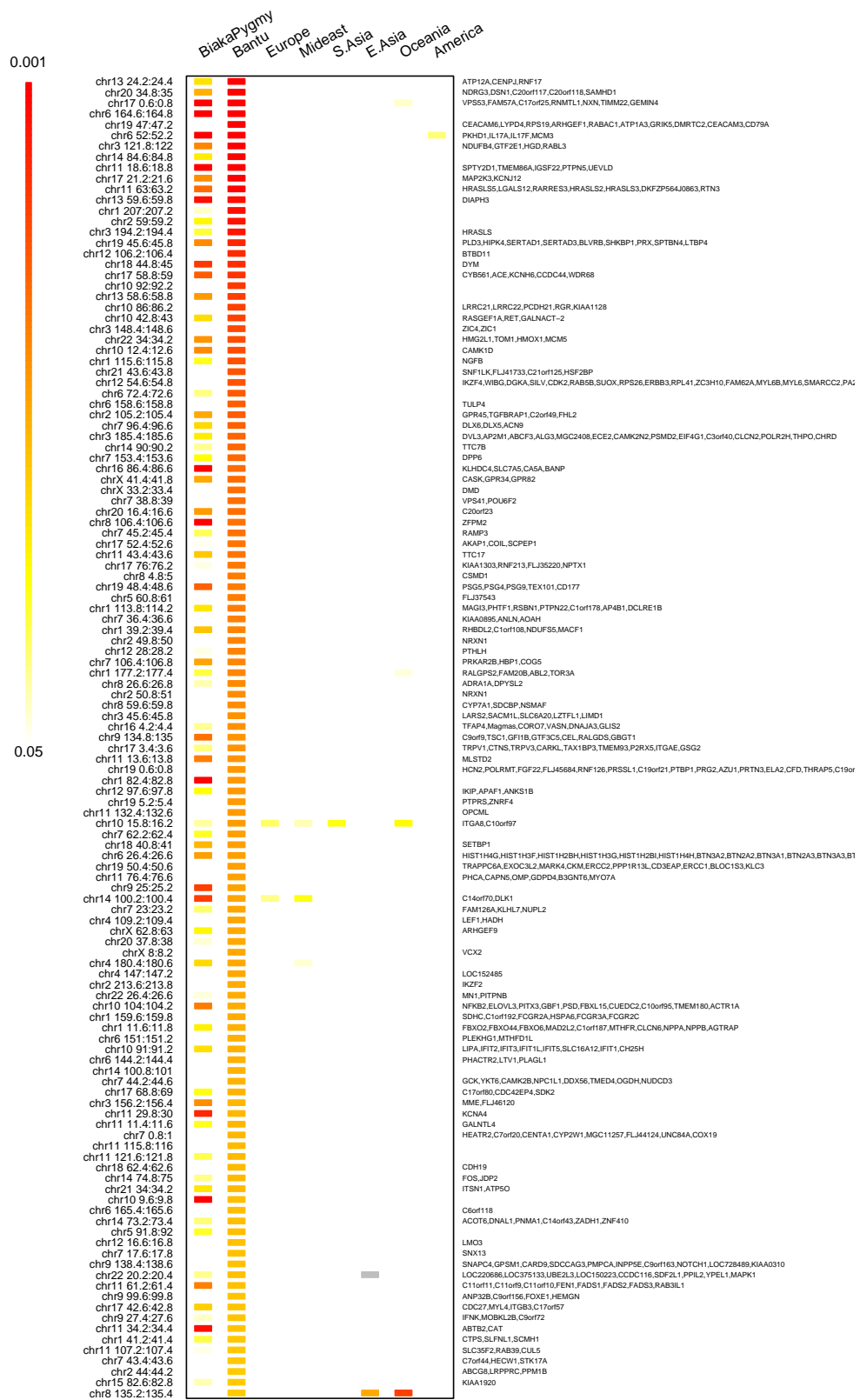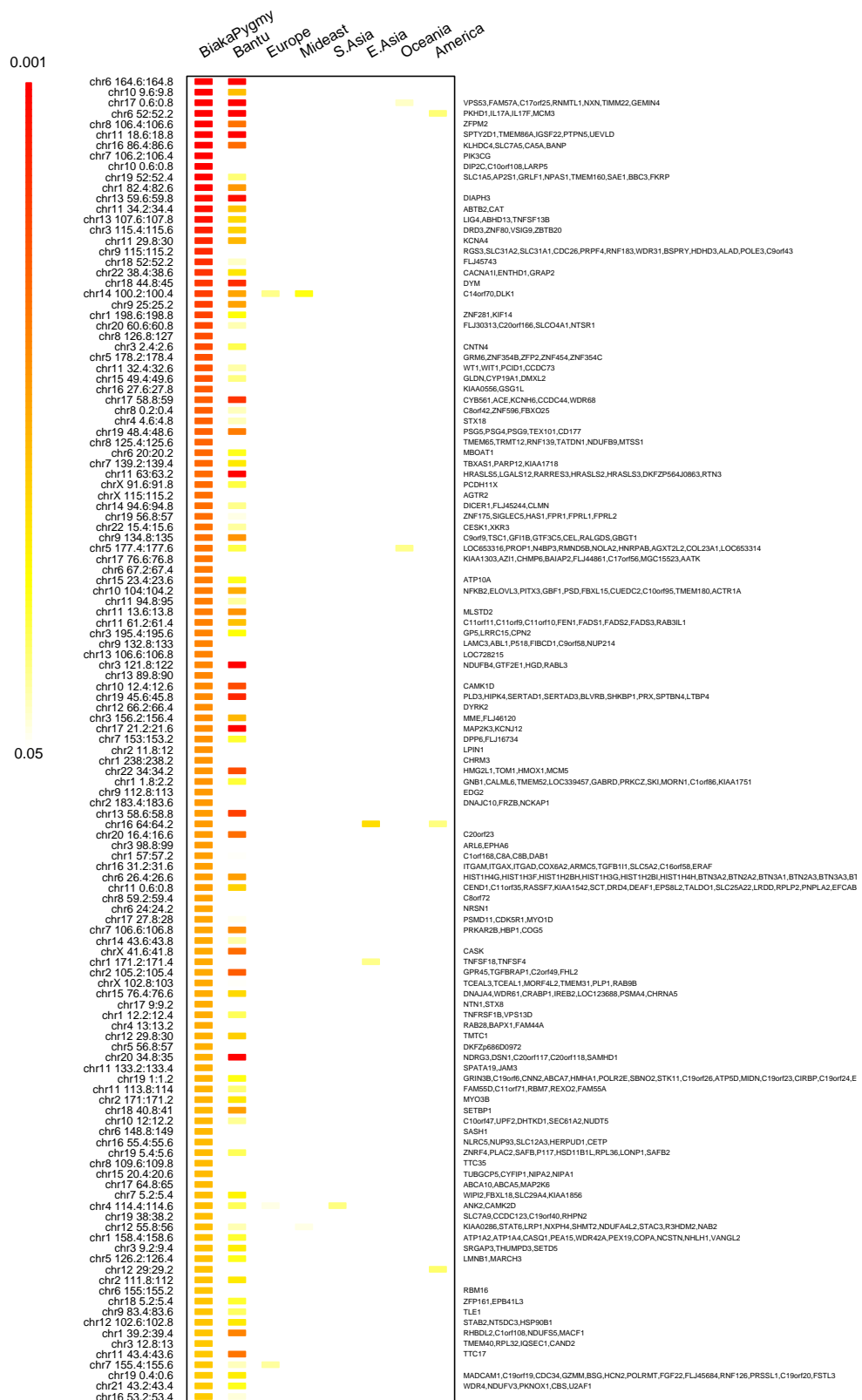
0.001

BiakaPygmy  Bantu  Europe  Mideast  S.Asia  E.Asia  Oceania  America

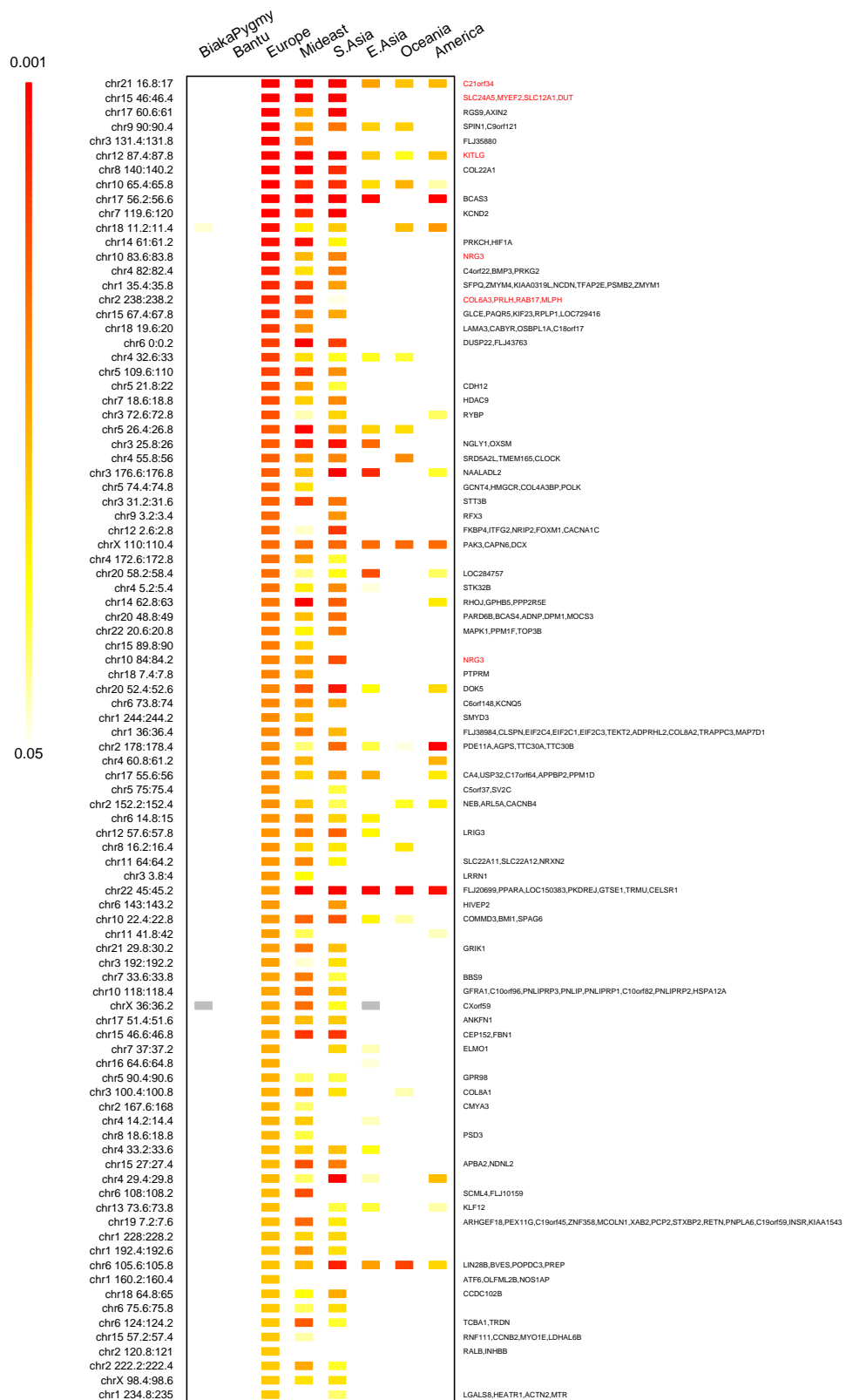| Position | Genes |
|---|---|
| chr21 16.8:17 | C21orf34 |
| chr15 46:46.4 | SLC24A5,MYEF2,SLC12A1,DUT |
| chr17 60.6:61 | RGS9,AXIN2 |
| chr9 90:90.4 | SPIN1,C9orf121 |
| chr3 131.4:131.8 | FLJ35880 |
| chr12 87.4:87.8 | KITLG |
| chr8 140:140.2 | COL22A1 |
| chr10 65.4:65.8 | |
| chr17 56.2:56.6 | BCAS3 |
| chr7 119.6:120 | KCND2 |
| chr18 11.2:11.4 | |
| chr14 61:61.2 | PRKCH,HIF1A |
| chr10 83.6:83.8 | NRG3 |
| chr4 82:82.4 | C4orf22,BMP3,PRKG2 |
| chr1 35.4:35.8 | SFPQ,ZMYM4,KIAA0319L,NCDN,TFAP2E,PSMB2,ZMYM1 |
| chr2 238:238.2 | COL6A3,PRLH,RAB17,MLPH |
| chr15 67.4:67.8 | GLCE,PAQR5,KIF23,RPLP1,LOC729416 |
| chr18 19.6:20 | LAMA3,CABYR,OSBPL1A,C18orf17 |
| chr6 0:0.2 | DUSP22,FLJ43763 |
| chr4 32.6:33 | |
| chr5 109.6:110 | |
| chr5 21.8:22 | CDH12 |
| chr7 18.6:18.8 | HDAC9 |
| chr3 72.6:72.8 | RYBP |
| chr5 26.4:26.8 | |
| chr3 25.8:26 | NGLY1,OXSM |
| chr4 55.8:56 | SRD5A2L,TMEM165,CLOCK |
| chr3 176.6:176.8 | NAALADL2 |
| chr5 74.4:74.8 | GCNT4,HMGCR,COL4A3BP,POLK |
| chr3 31.2:31.6 | STT3B |
| chr9 3.2:3.4 | RFX3 |
| chr12 2.6:2.8 | FKBP4,ITFG2,NRIP2,FOXM1,CACNA1C |
| chrX 110:110.4 | PAK3,CAPN6,DCX |
| chr4 172.6:172.8 | |
| chr20 58.2:58.4 | LOC284757 |
| chr4 5.2:5.4 | STK32B |
| chr14 62.8:63 | RHOJ,GPHB5,PPP2R5E |
| chr20 48.8:49 | PARD6B,BCAS4,ADNP,DPM1,MOCS3 |
| chr22 20.6:20.8 | MAPK1,PPM1F,TOP3B |
| chr15 89.8:90 | |
| chr10 84:84.2 | NRG3 |
| chr18 7.4:7.8 | PTPRM |
| chr20 52.4:52.6 | DOK5 |
| chr6 73.8:74 | C6orf148,KCNQ5 |
| chr1 244:244.2 | SMYD3 |
| chr1 36:36.4 | FLJ38984,CLSPN,EIF2C4,EIF2C1,EIF2C3,TEKT2,ADPRHL2,COL8A2,TRAPPC3,MAP7D1 |
| chr2 178:178.4 | PDE11A,AGPS,TTC30A,TTC30B |
| chr4 60.8:61.2 | |
| chr17 55.6:56 | CA4,USP32,C17orf64,APPBP2,PPM1D |
| chr5 75:75.4 | C5orf37,SV2C |
| chr2 152.2:152.4 | NEB,ARL5A,CACNB4 |
| chr6 14.8:15 | |
| chr12 57.6:57.8 | LRIG3 |
| chr8 16.2:16.4 | |
| chr11 64:64.2 | SLC22A11,SLC22A12,NRXN2 |
| chr3 3.8:4 | LRRN1 |
| chr22 45:45.2 | FLJ20699,PPARA,LOC150383,PKDREJ,GTSE1,TRMU,CELSR1 |
| chr6 143:143.2 | HIVEP2 |
| chr10 22.4:22.8 | COMMD3,BMI1,SPAG6 |
| chr11 41.8:42 | |
| chr21 29.8:30.2 | GRIK1 |
| chr3 192:192.2 | |
| chr7 33.6:33.8 | BBS9 |
| chr10 118:118.4 | GFRA1,C10orf96,PNLIPRP3,PNLIP,PNLIPRP1,C10orf82,PNLIPRP2,HSPA12A |
| chrX 36:36.2 | CXorf59 |
| chr17 51.4:51.6 | ANKFN1 |
| chr15 46.6:46.8 | CEP152,FBN1 |
| chr7 37:37.2 | ELMO1 |
| chr16 64.6:64.8 | |
| chr5 90.4:90.6 | GPR98 |
| chr3 100.4:100.8 | COL8A1 |
| chr2 167.6:168 | CMYA3 |
| chr4 14.2:14.4 | |
| chr8 18.6:18.8 | PSD3 |
| chr4 33.2:33.6 | |
| chr15 27:27.4 | APBA2,NDNL2 |
| chr4 29.4:29.8 | |
| chr6 108:108.2 | SCML4,FLJ10159 |
| chr13 73.6:73.8 | KLF12 |
| chr19 7.2:7.6 | ARHGEF18,PEX11G,C19orf45,ZNF358,MCOLN1,XAB2,PCP2,STXBP2,RETN,PNPLA6,C19orf59,INSR,KIAA1543 |
| chr1 228:228.2 | |
| chr1 192.4:192.6 | |
| chr6 105.6:105.8 | LIN28B,BVES,POPDC3,PREP |
| chr1 160.2:160.4 | ATF6,OLFML2B,NOS1AP |
| chr18 64.8:65 | CCDC102B |
| chr6 75.6:75.8 | |
| chr6 124:124.2 | TCBA1,TRDN |
| chr15 57.2:57.4 | RNF111,CCNB2,MYO1E,LDHAL6B |
| chr2 120.8:121 | RALB,INHBB |
| chr2 222.2:222.4 | |
| chrX 98.4:98.6 | |
| chr1 234.8:235 | LGALS8,HEATR1,ACTN2,MTR |

0.05

Figure 23: The top 1% of XP-EHH signals in Europe, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.
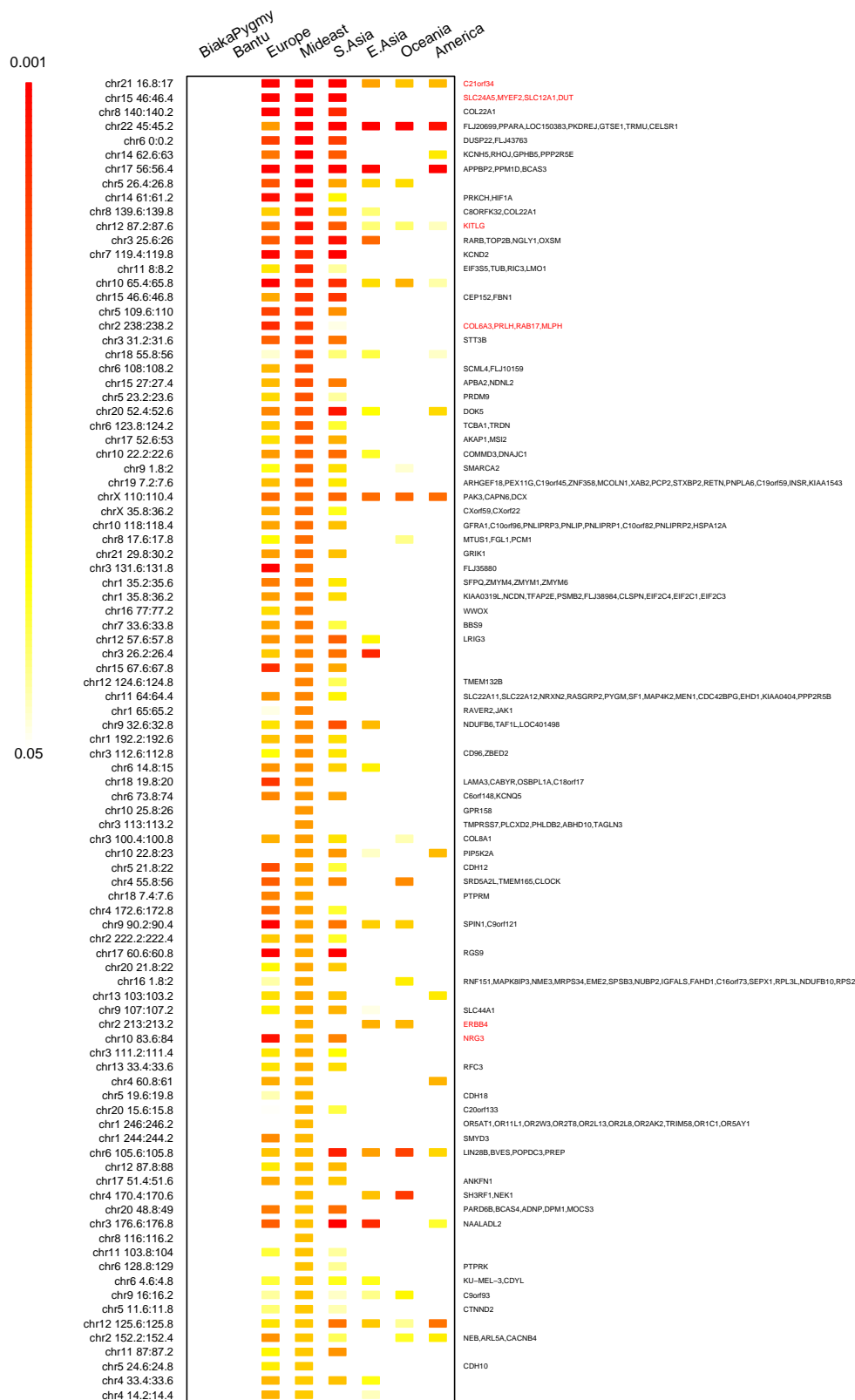
Figure 24: The top 1% of XP-EHH signals in the Middle East, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.
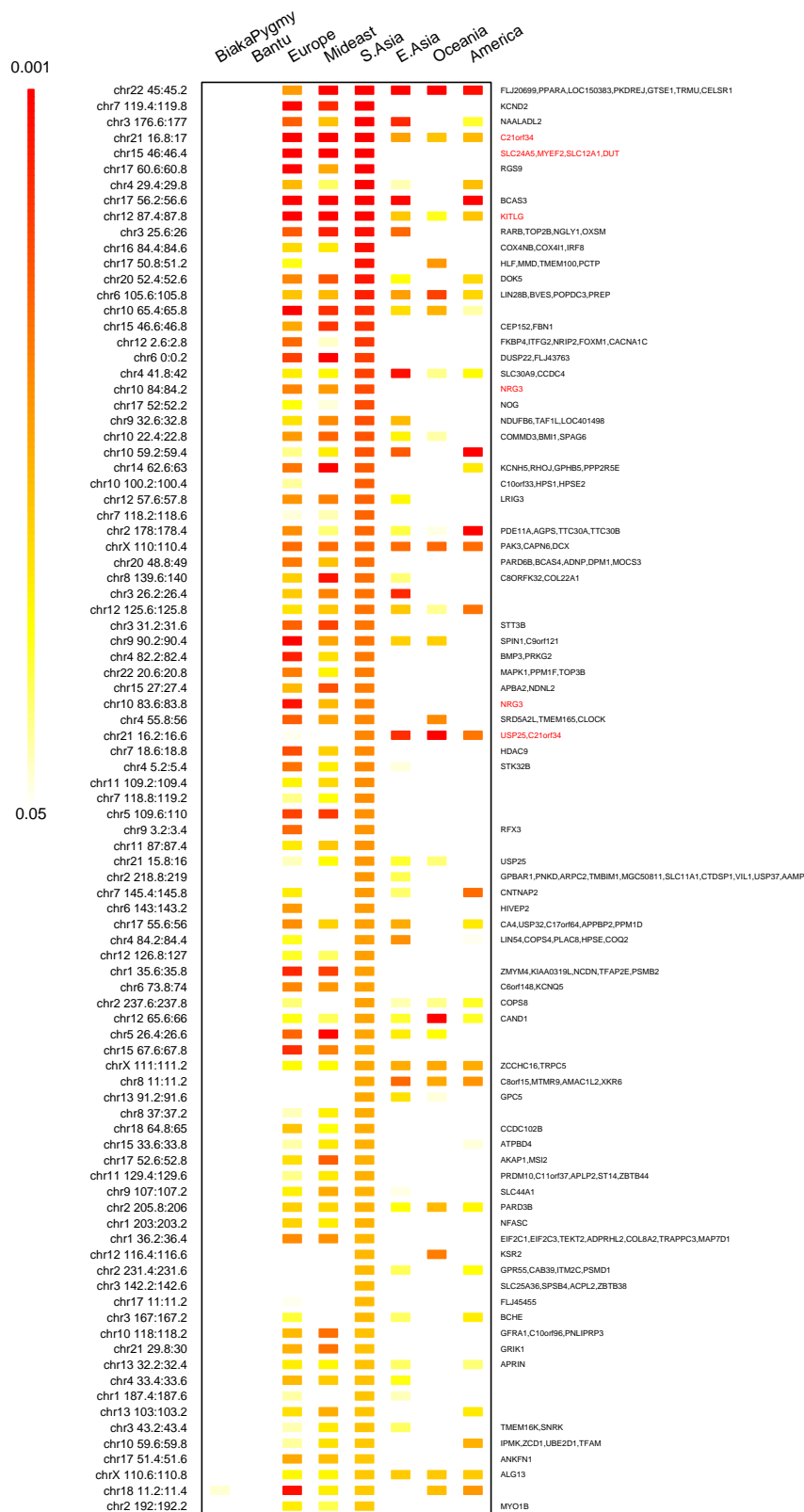
Figure 25: The top 1% of XP-EHH signals in South Asia, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.
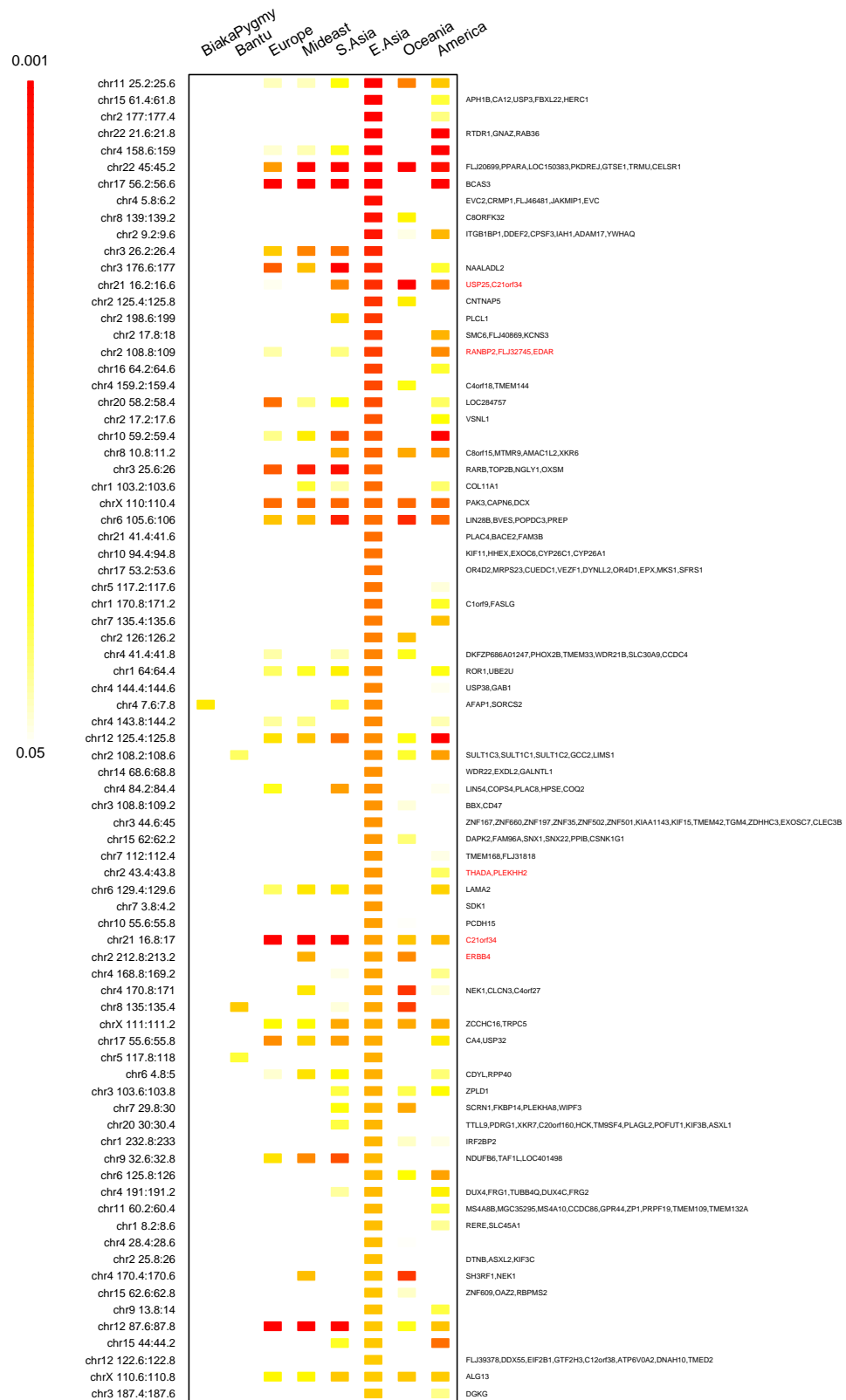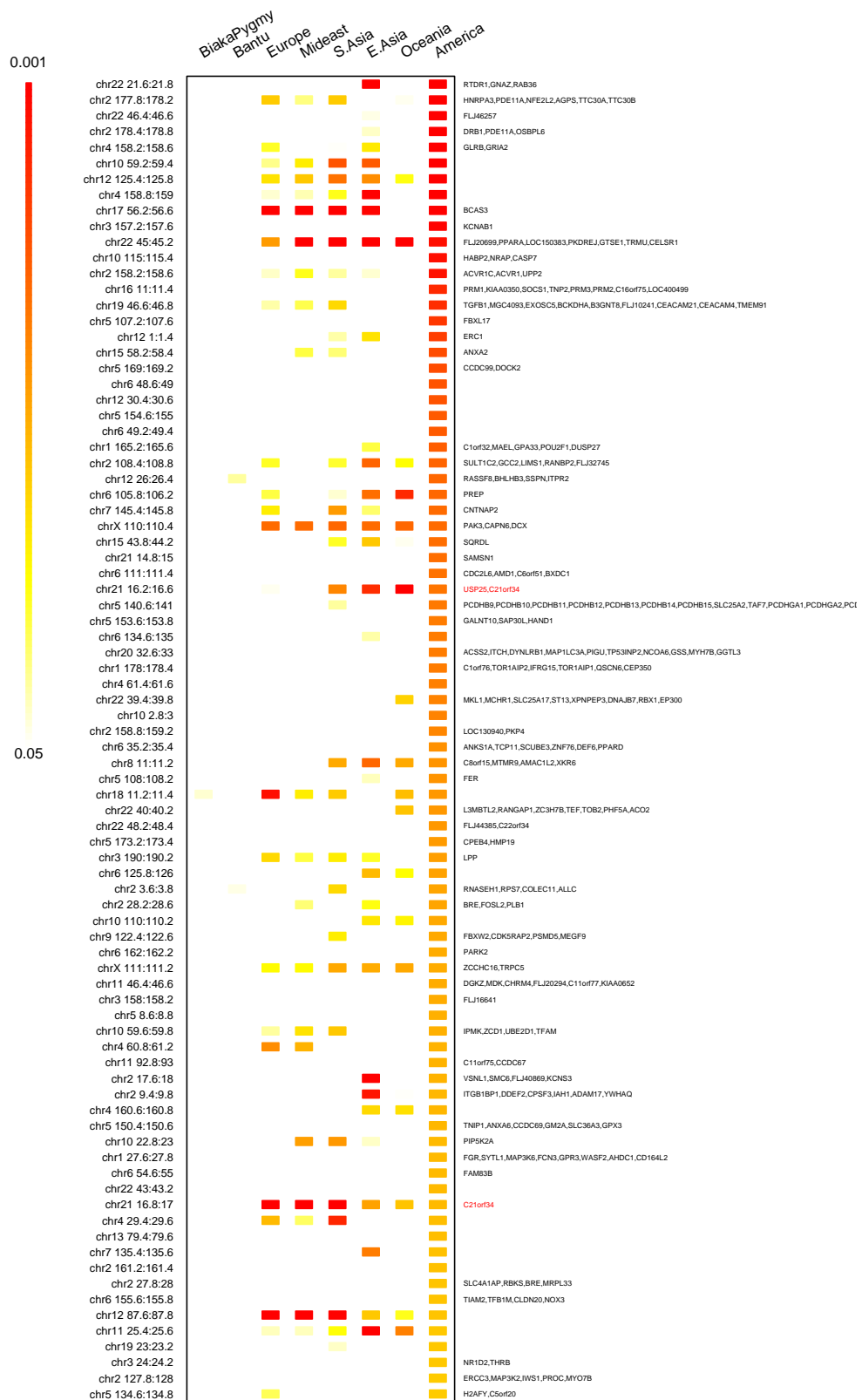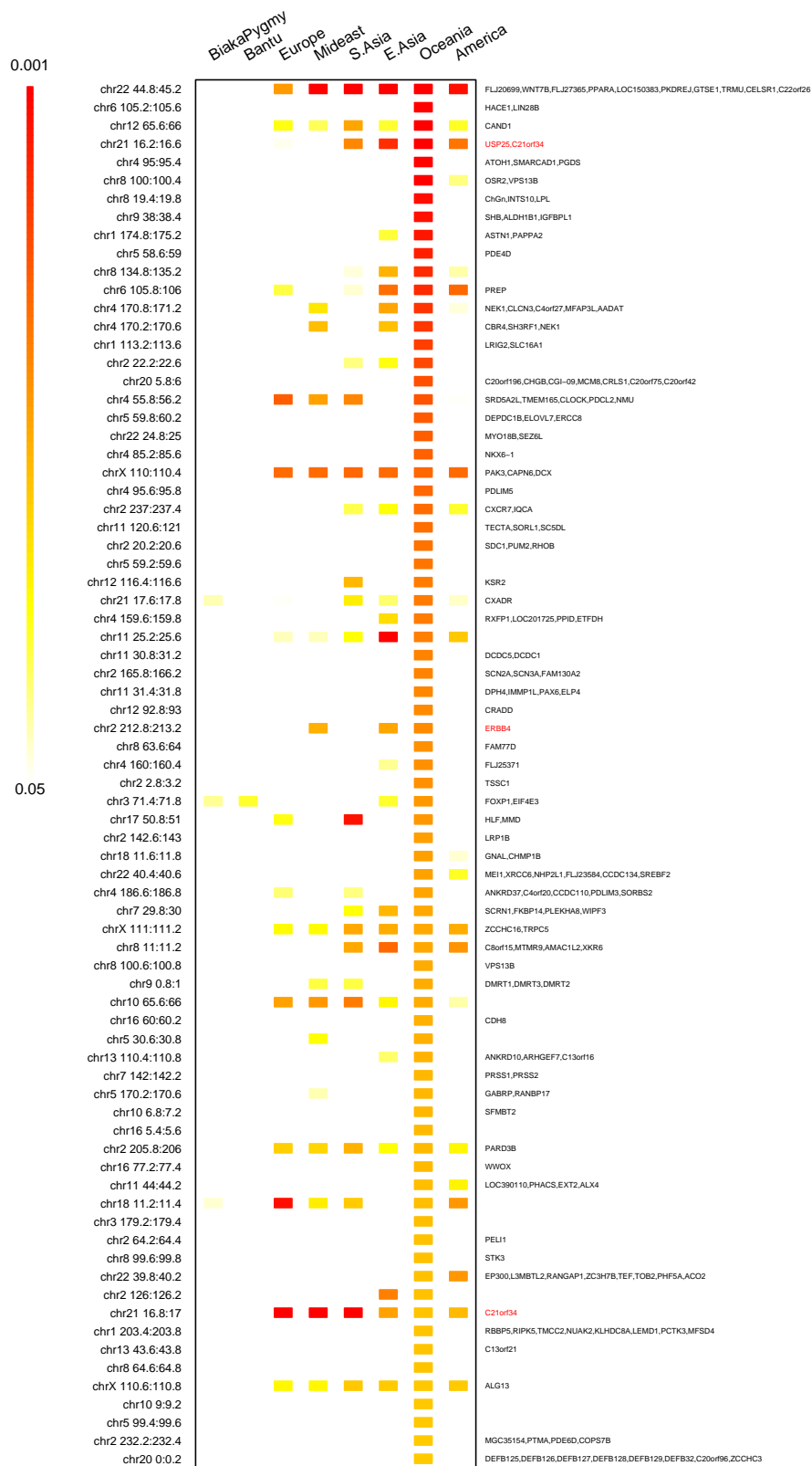
Figure 26: The top 1% of XP-EHH signals in East Asia, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 27: The top 1% of XP-EHH signals in the Americas, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

Figure 28: The top 1% of XP-EHH signals in Oceania, ordered from top to bottom in order of significance. See the caption of Figure 1 in the main text for details.

# References

[1] Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R (2005) Ascertainment bias in studies of human genome-wide polymorphism. Genome Res 15: 1496–1502.

[2] Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, et al. (2007) Genome-wide detection and characterization of positive selection in human populations. Nature 449: 913–918.

[3] Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ, et al. (2005) Calibrating a coalescent simulation of human genome sequence variation. Genome Res 15: 1576–1583.

[4] Tang K, Thornton K, Stoneking M (2007) A New Approach for Using Genome Scans to Detect Recent Positive Selection in the Human Genome. PLoS Biol 5: e171.

[5] Teshima KM, Coop G, Przeworski M (2006) How reliable are empirical genomic scans for selective sweeps? Genome Res 16: 702–712.

[6] Thornton KR, Jensen JD (2007) Controlling the false-positive rate in multilocus genome scans for selection. Genetics 175: 737–750.

[7] Voight BF, Kudaravalli S, Wen X, Pritchard JK (2006) A map of recent positive selection in the human genome. PLoS Biol 4: e72.