

Supplementary Information

Genetic variants affecting equivalent protein family positions reflect human diversity

Francesco Raimondi (1,2), Matthew J. Betts (1,2), Qianhao Lu (1,2), Asuka Inoue (3,4), J. Silvio Gutkind (5), Robert B. Russell (1,2)*

1. CellNetworks, Bioquant, Heidelberg University, Im Neuenheimer Feld 267, 69120 Heidelberg, Germany
2. Biochemie Zentrum Heidelberg (BZH), Heidelberg University, Im Neuenheimer Feld 328, 69120, Heidelberg, Germany
3. Graduate School of Pharmaceutical Science, Tohoku University, Sendai, Miyagi, Japan
4. Japan Science and Technology Agency (JST), Precursory Research for Embryonic Science and Technology (PRESTO), Kawaguchi, Saitama, Japan.
5. Moores Cancer Center, University of San Diego, San Diego, USA

* Corresponding author:

Email: bq_rrussell@bioquant.uni-heidelberg.de

Tel: +49 6221 54 51 362

Supplemental Figures & Tables

Figure S1

GO biological processes linked to genes in families most affected by variants; percentages of genes are shown in the chart or next to the names.

Figure S2

a) Allelic counts versus position within the structure **b)** cartoon of a GPCR representative structure (PDB code 4eiy) showing the top 10 positions most affected by variants, shown as spheres centered on C β or C α atoms. The size of the spheres is proportional to the variant counts as for Figure 3a.

Figure S3

Tree of OR3A1 dRy alleles generated from the PAML analysis. The branches corresponding to dRy minor alleles are highlighted in orange.

Figure S4

Same representation as Fig S3 for OR8D2.

Figure S5

Hierarchical clustering of 2,504 individuals (Y-axis) based on distances between fingerprints of OR dRy alleles within ORs (X-axis). Darker red dots indicate homozygous variants. Left and bottom dendrograms indicate similarity of groups of individuals/receptors based on OR dRy fingerprints. Dendrogram color scheme reflect clustering based on fingerprint similarity. Left color band between dendrogram and main matrix indicate super population membership of each individual. The matrix at the top shows associations between odorants and olfactory receptors(35).

Figure S6

Anosmia estimation based on observed loss-of-function variants. For each allele in each of 2504 individuals, we summed log(EC50) values for all the odorant-receptor pairs, unless a deleterious variant (either dRy, Stop gain or Frameshift) was observed. We

considered the overall response of an individual to a given odorant as the average of the log(EC50) values. We considered the average of all the individual responses to a given odorant for the members of the same population and we normalized each ligand response relative to the maximum value.

Figure S7

a) Genotypic diversity versus the number of animals in 50 simulations (ten for each of 0, 25, 50, 75 or 100% odorant driven movement) for eight different random starting environments. Note how the higher odorant driven simulations (yellow, orange, red) tend to be towards the bottom of the plots indicating the highest diversity in genotype. **b)** Genotypic diversity versus the number of iterations (time step) in 50 simulations (ten for each of 0, 25, 50, 75 or 100% odorant driven movement) for eight different random starting environments. Note how the higher odorant driven simulations (yellow, orange, red) tend to be towards the bottom of the plots indicating the highest diversity in genotype.

Table S1

Protein family domains most enriched in missense variants.

Table S2

ORs from Swissprot/Uniprot annotated with family and pseudogene classification (from HORDE) and with deleterious mutations information (dRy, Stop gains and frameshifts)

Table S3

Protein family domain positions most enriched in missense variants.

Table S4

Assessing the ability of mutations from canonical GPCR residues to predict odorant binding from OR/odorant screens (19, 34, 35)

Table S5

dRy missense variants and associated Minor Allele Frequency (MAF) and zigosity for **a)** non-OR GPCRs and **b)** OR GPCRs.

Table S6

a) Total counts of variants and total number of positions for different Arginine positions within humans. Ratios (rightmost column) and enrichments (bottom) are also shown. **b)** Amino acids mutation enrichments and log odds in OR and non-OR genes relative to other genes. **c)** Non-synonymous/synonymous variants ratio for variants at dRy as well as other positions for OR/non-OR. **d)** Non-synonymous/synonymous variant ratio for all human genes .

Table S7

PAML analysis of OR alleles from the 1000 Genomes project

Table S8

dRy, stop-gain and frameshift variants within ORs .

Table S9

OR dRy arginine missense variant enrichment within each continental group.

Table S10

a) Enrichment of OR dRy variants within clusters of individuals defined through variant profile similarity and **b)** ethnic group composition of each cluster.

Table S11

Estimation of Anosmia by combining information on OR loss-of-function alleles in each of 2505 individuals in the 1000 genomes dataset with information on odorant/OR binding data taken from Mainland et al.

Movie S1

Examples of simulating animal behavior in a 2D matrix with scattered odorants and where animals contain 100 receptor genes that determine attraction or repulsion to one or more odorants. More details are given in the movie itself and in Online Methods.

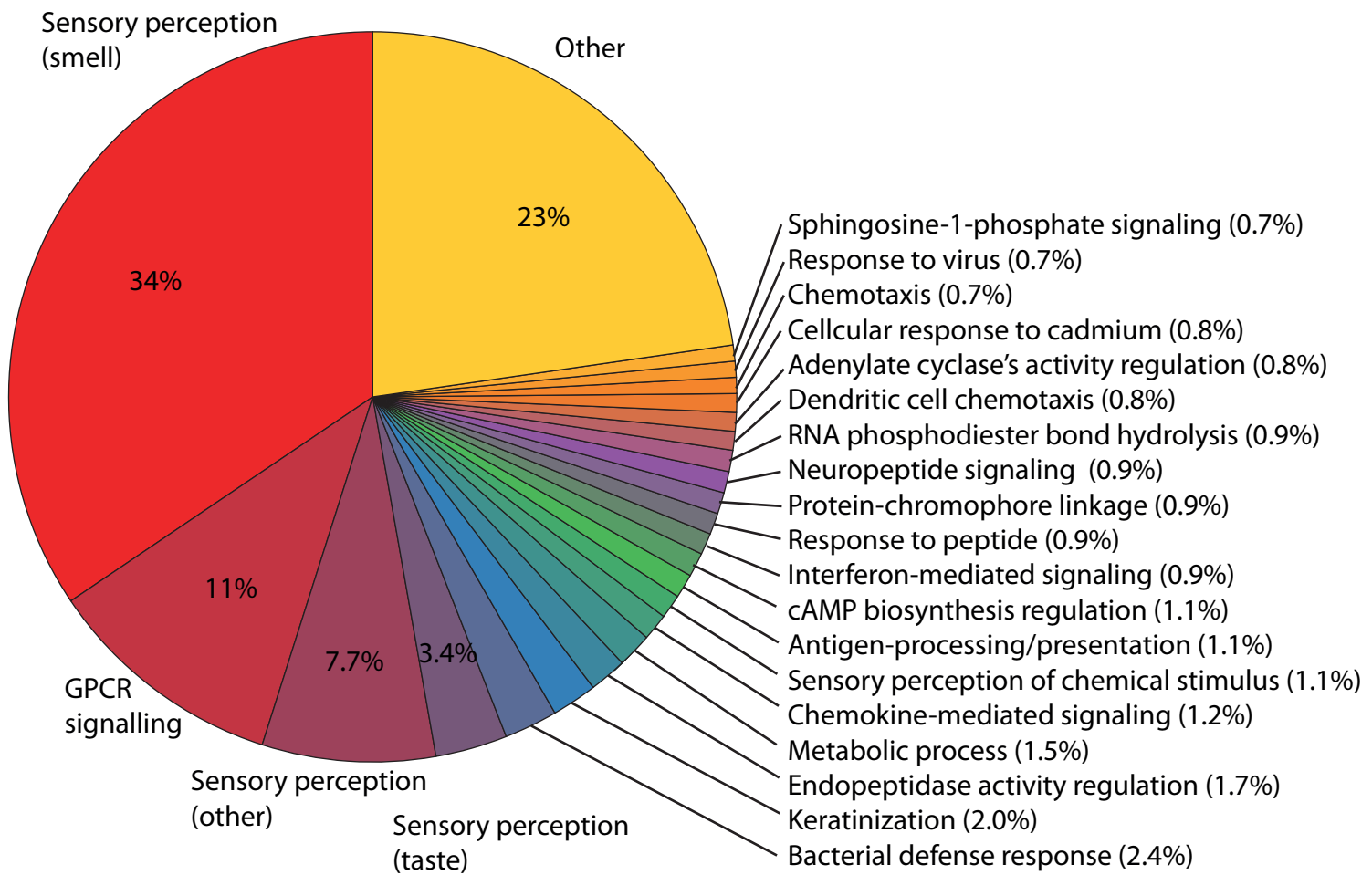


Figure S1

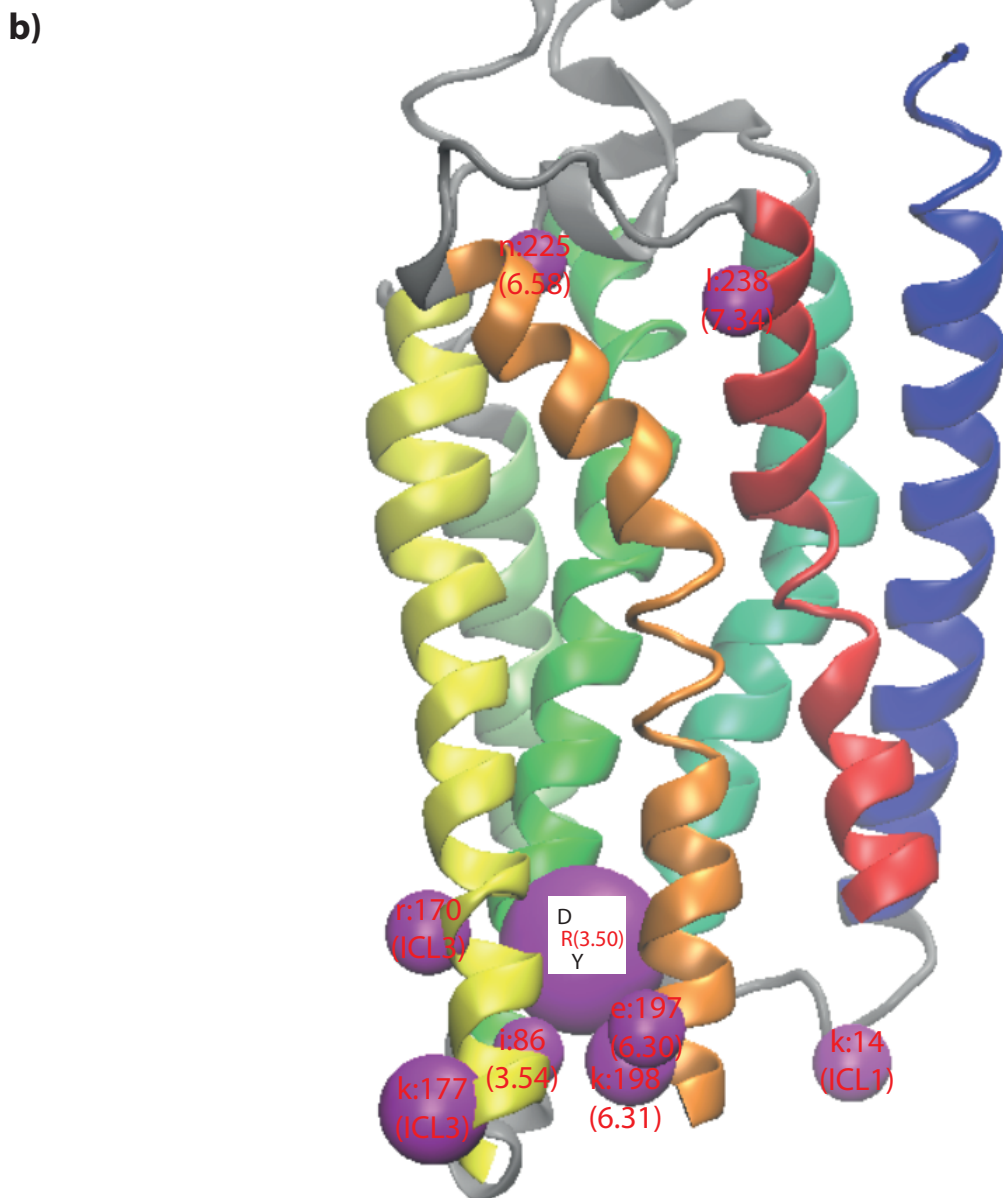
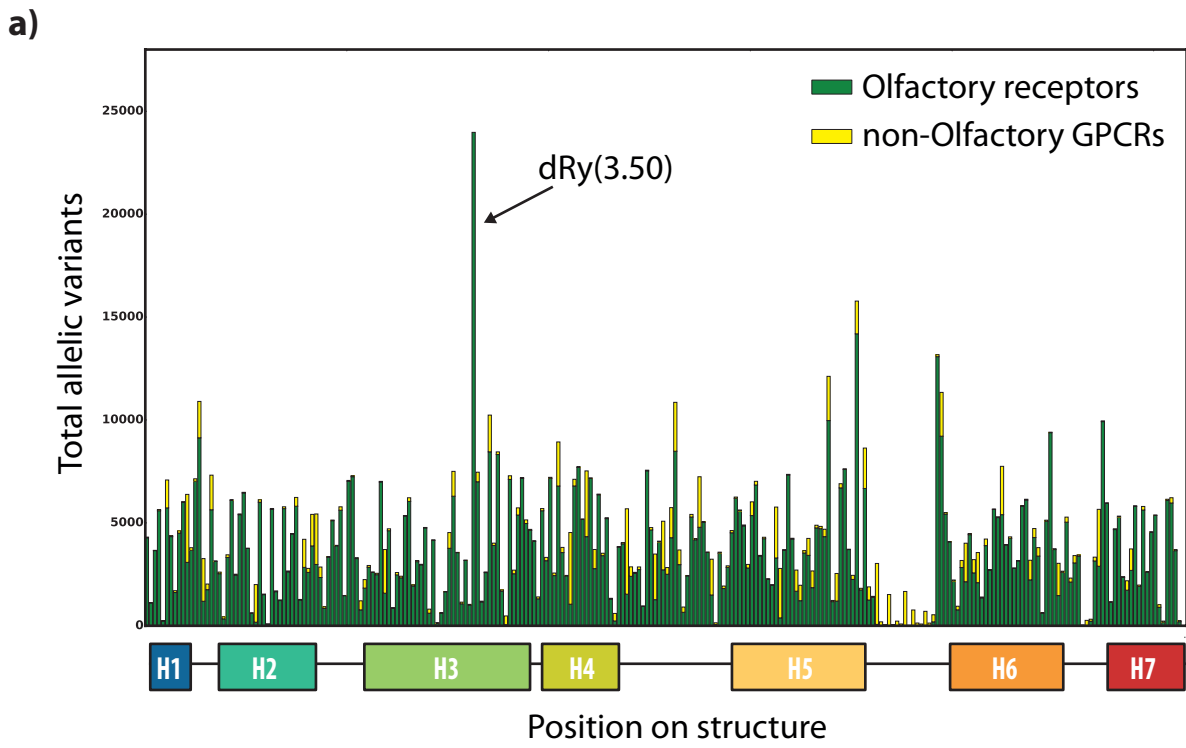


Figure S2

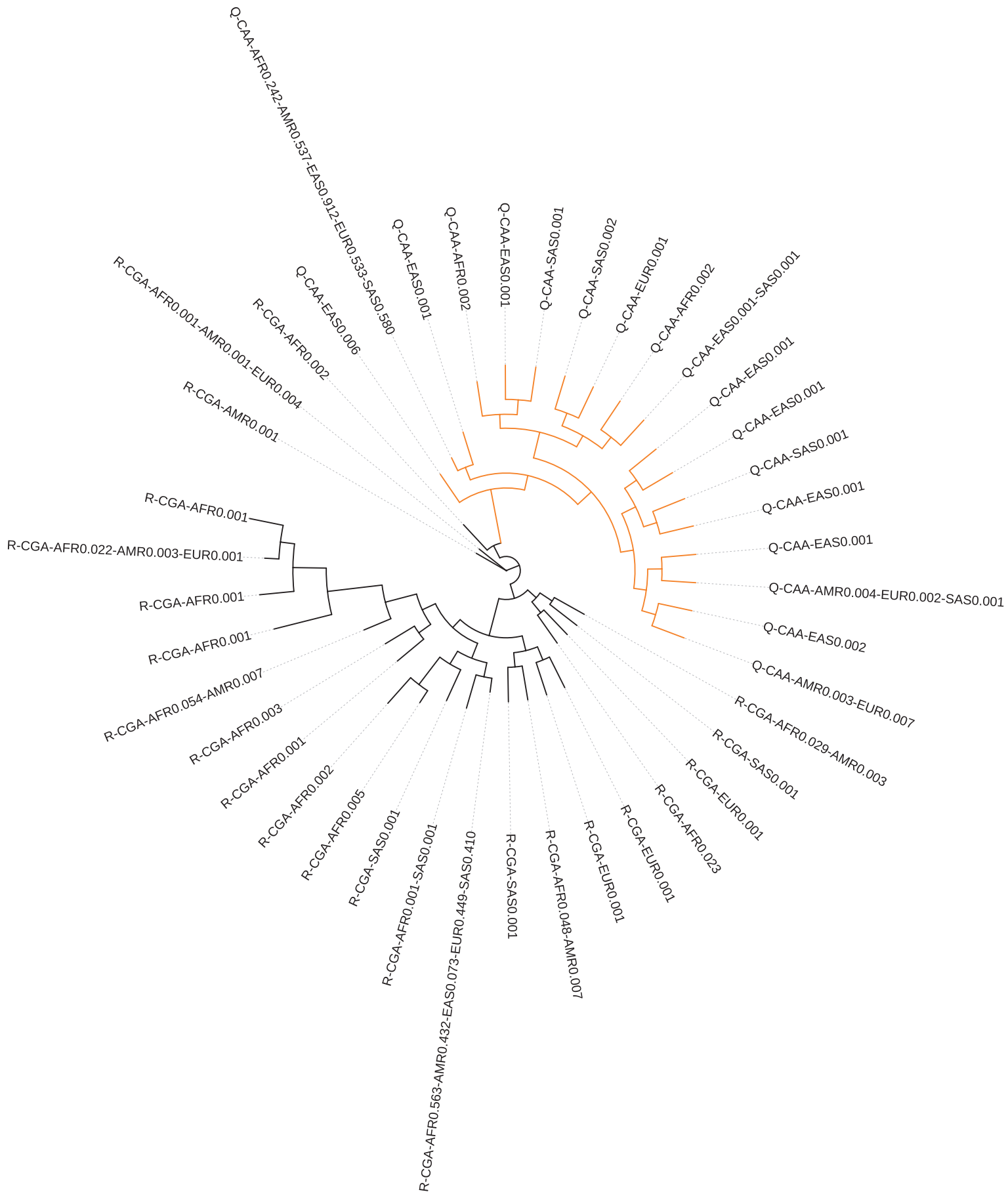


Figure S3

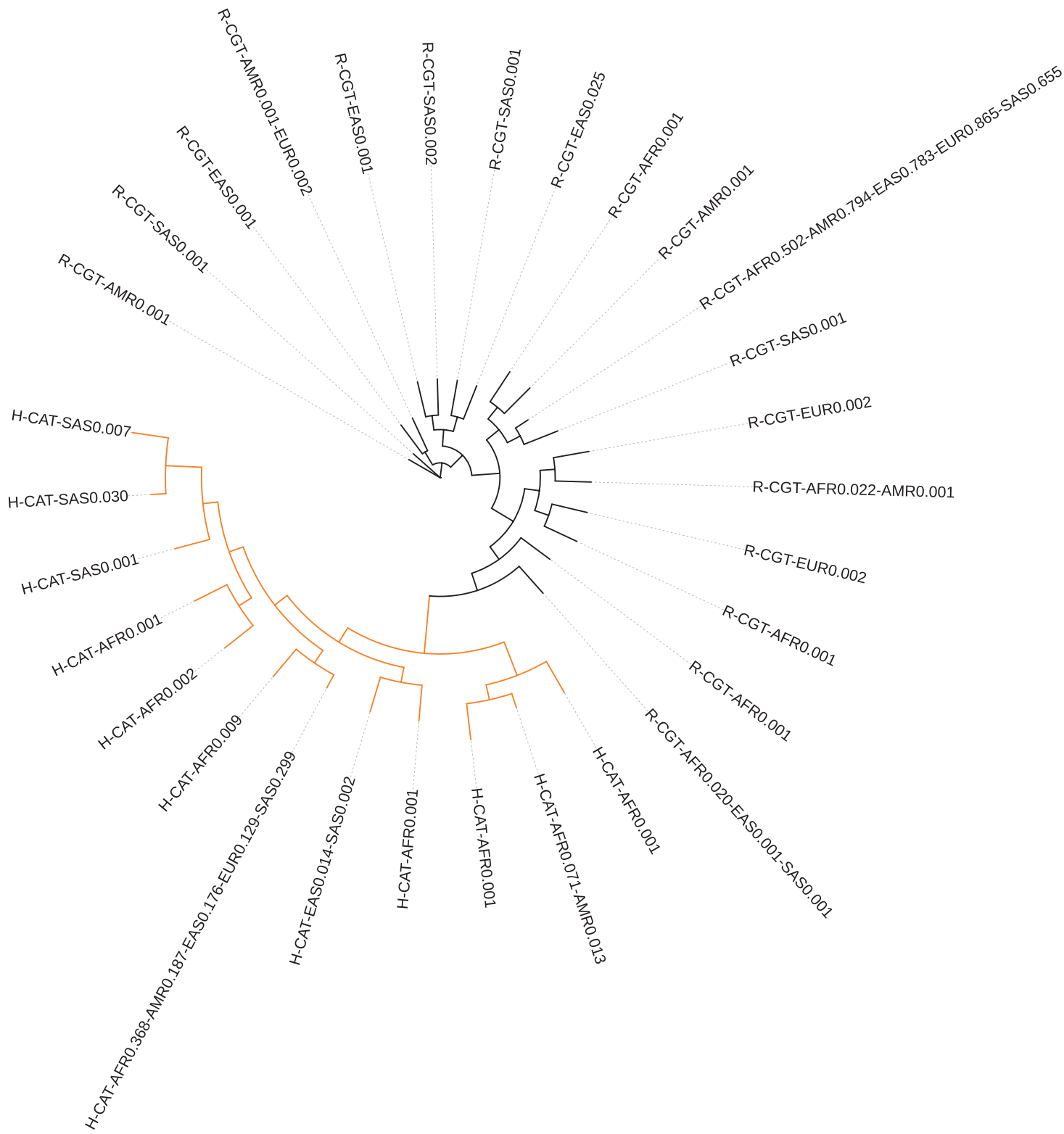


Figure S4

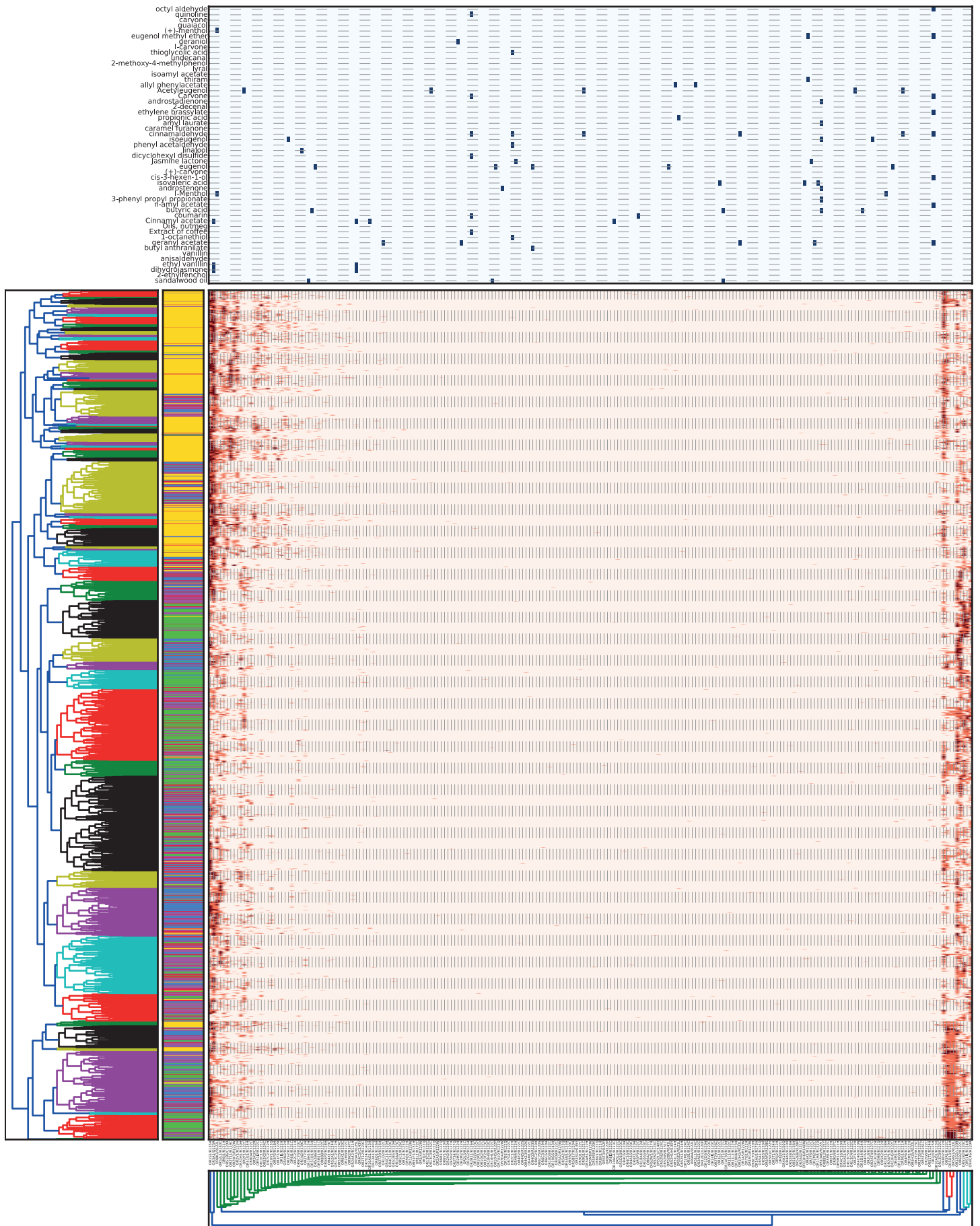


Figure S5

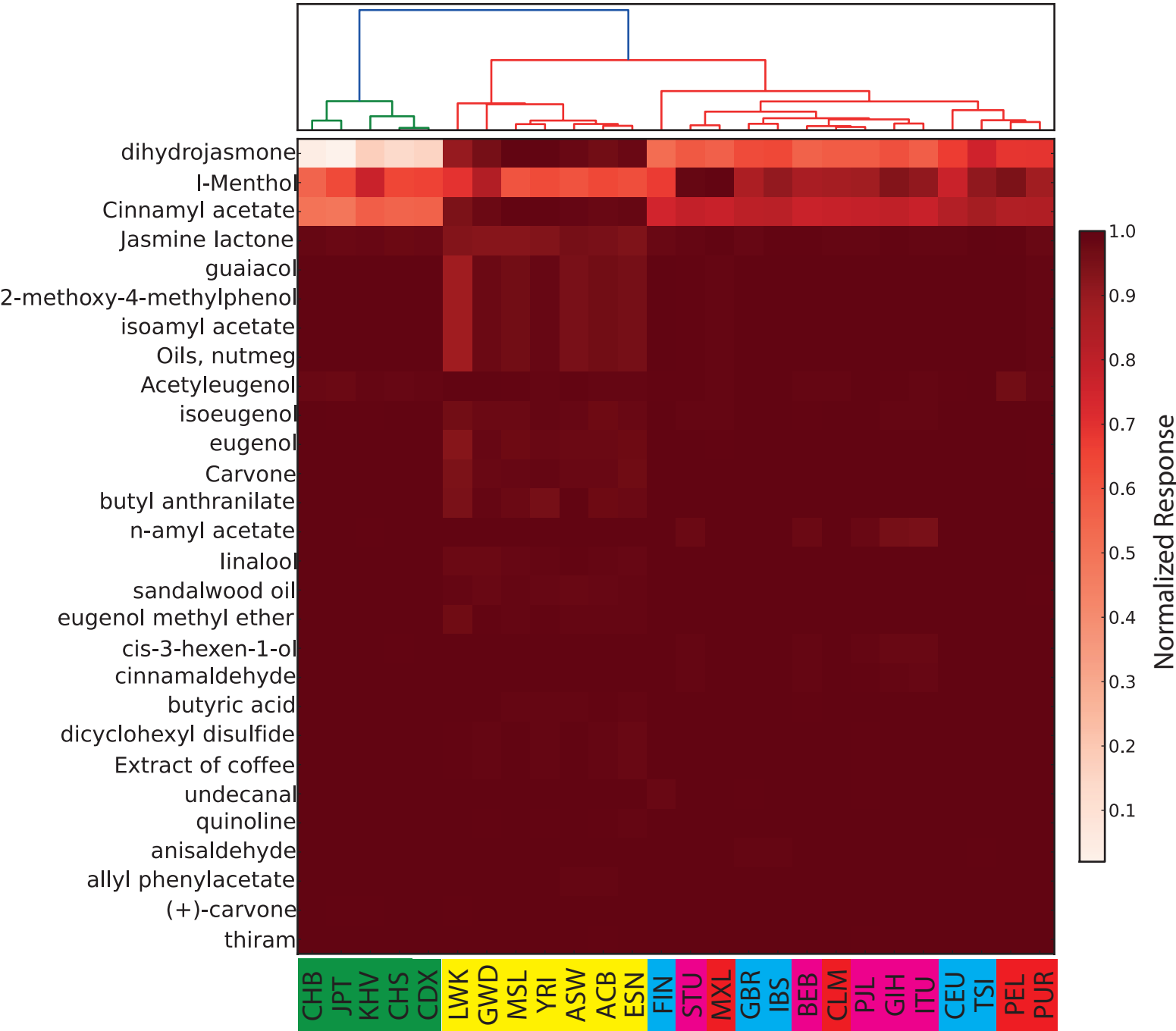


Figure S6

Distribution among super populations

- Ad Mixed American
- African
- European
- East Asian
- South Asian

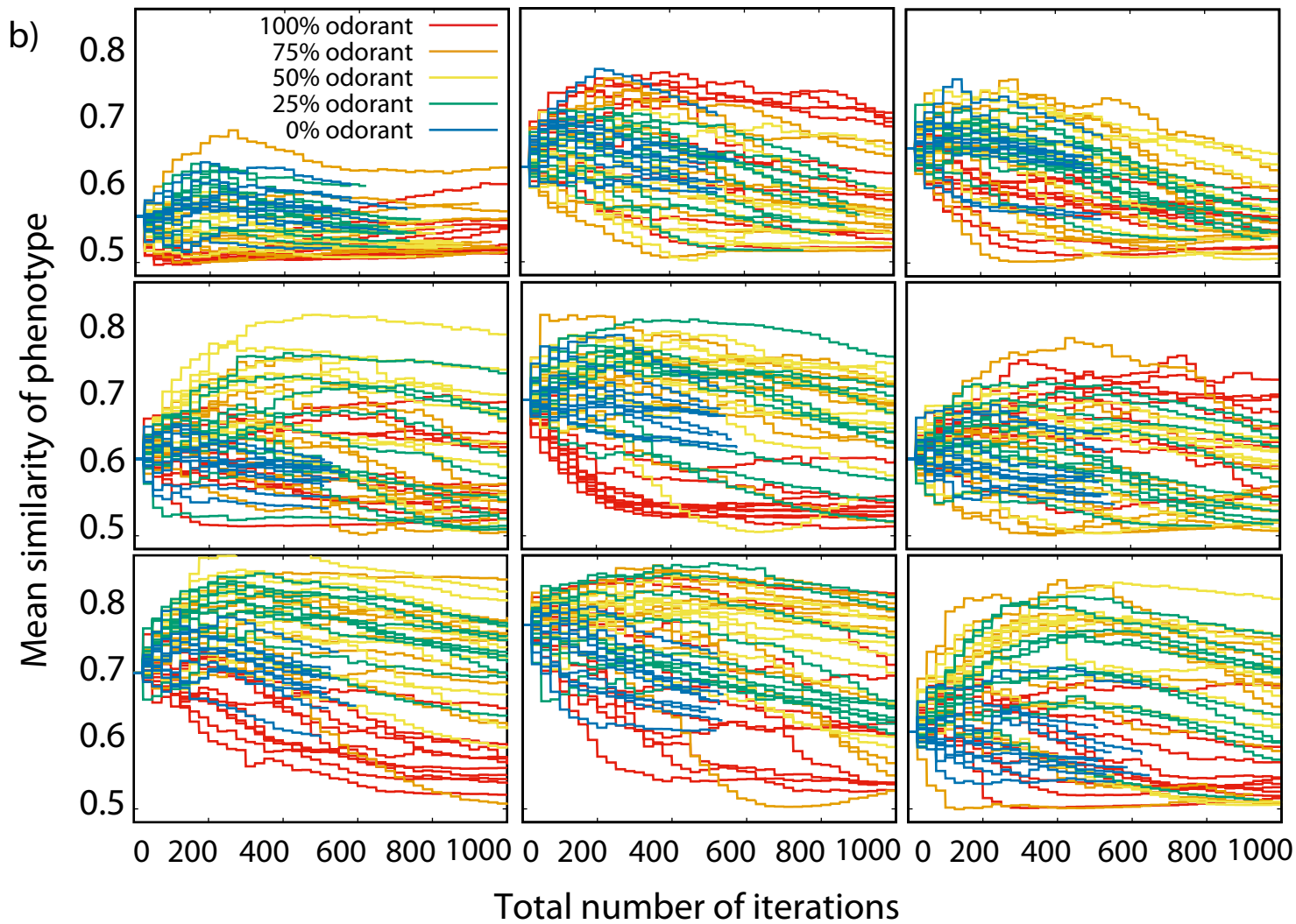
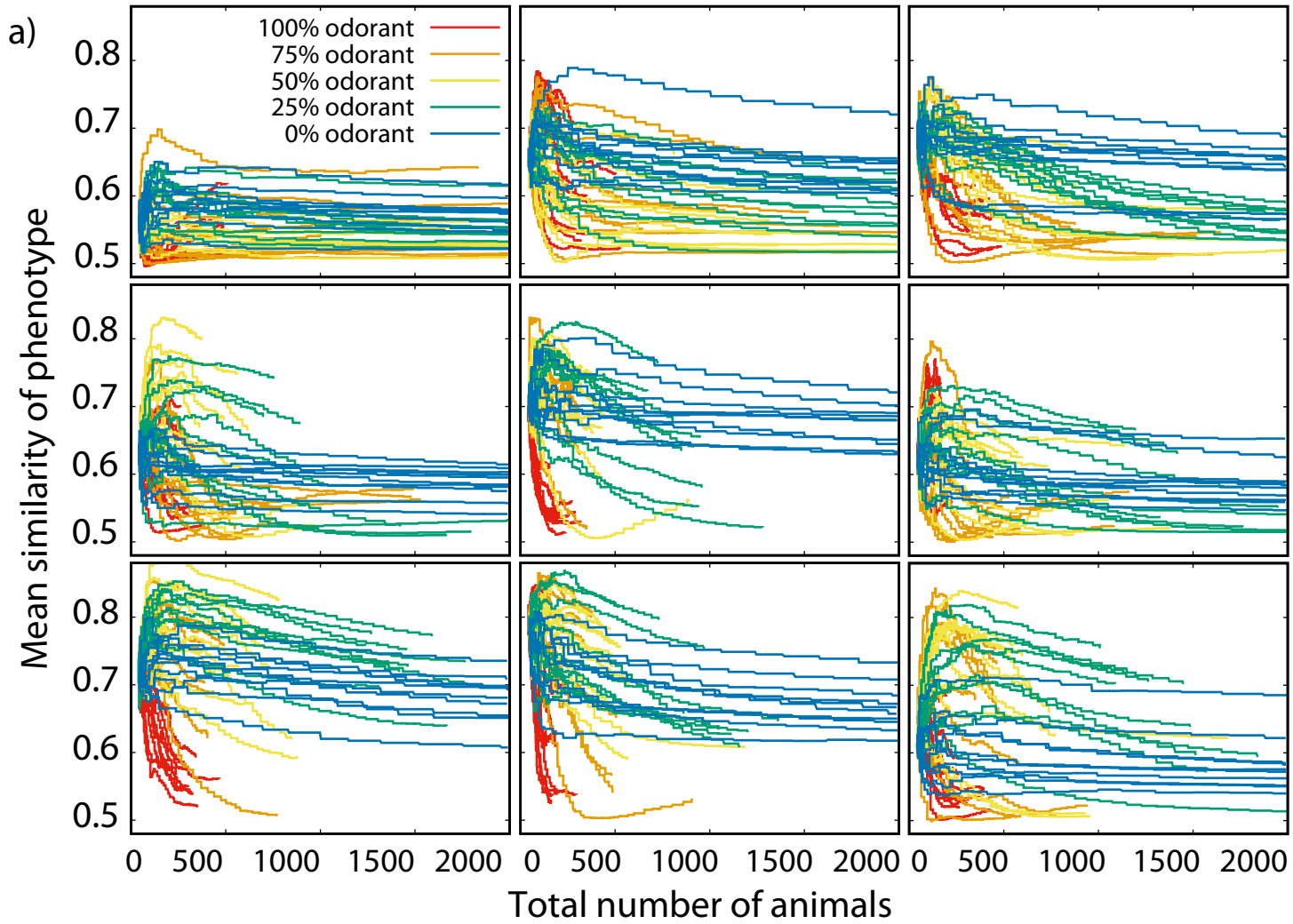


Figure S7