

Supplemental figures

Legends:

Figure S1: (A) Per sample mutation statistics: number of coding small variants (SNVs and small indels), number of non-coding small variants, ratio of coding vs. non-coding small variants, number of structural variants, and aberrant genome fraction. The number of noncoding small variants and the ratio of coding vs. non-coding small variants are not shown for two samples without matched normal control. (B) Scatterplot of the number of SVs vs. small variant count. (C - E) Integrative plots of recurrent CNAs for FL (C), FL-DLBCL (D), DLBCL (E). Non-diploid cases were excluded.

Figure S2: (A) Additional rainfall plots with increasing mutational load from top to bottom. Genomic regions affected by Psichales appear as vertical streaks, but less focalized than the kataegis-like clusters or rainfalls. As in Figure 1, special genomic regions including those targeted by hallmark translocations are highlighted. (B) Overlaid violin plots and boxplots displaying distributions of different per sample quantities across gcBCL subgroups: (i) number of SNVs; (ii) number of kataegis clusters; (iii) number of affected Kataegis regions inside the IG loci; (iv) number of affected Kataegis regions outside the IG loci; (v) overall number of affected Kataegis regions; (vi) counts of SNVs inside kataegis clusters; (vii) ratio of affected kataegis regions inside the IG loci vs. all affected Kataegis regions (*ROIs_Ka_ratio*); (viii) counts of SNVs inside Psichales clusters; (ix) counts of SNVs inside affected Psichales regions; and (x) ratio of affected kataegis regions vs. affected Psichales regions. DLBCLs, FL/DLBCLs and FLs had comparable numbers of psichales clusters and affected psichales regions, but DLBCLs and FL/DLBCLs exhibited significantly higher numbers of SNVs in psichales clusters and regions than FLs, reflecting their overall higher mutational load. (C) Overlaid violin plots and boxplots displaying distributions of different per sample quantities across gcBCL subgroups: (i) counts of SNVs in kataegis clusters of type “other” (*SNVs_Mech_other*); (ii) counts of SNVs in CSR-like kataegis clusters (*SNVs_Mech_CSR*); (iii) counts of SNVs in SHM-like Kataegis clusters (*SNVs_Mech_SHM*); (iv) ratio of counts of SNVs in Kataegis clusters of type “other” vs. counts of SNVs in all kataegis clusters (*SNVs_Mech_other_ratio*); (v) ratio of counts of SNVs in CSR-like kataegis clusters vs. counts of SNVs in all kataegis clusters (*SNVs_Mech_CSR_ratio*); (vi) ratio of counts of SNVs in SHM-like Kataegis clusters vs. counts of SNVs in all kataegis clusters (*SNVs_Mech_SHM_ratio*). (D) The first four items (i) – (iv) are analogous to the first, second, fifth and sixth item of (B), but for DLBCL subtypes. (v) counts of SNVs inside kataegis clusters overlapping with kataegis regions; (vi) ratio of SNVs in kataegis cluster overlapping with kataegis regions vs. SNVs in all kataegis clusters; (vii) same ratio as in (vi), but restricted to the IG loci. (E) The first two items are analogous to the second and third item in (C), but for DLBCL subtypes; (iii) ratio of counts of SNVs in CSR-like kataegis clusters vs. counts of SNVs in SHM-like kataegis clusters. (F) Clustering of the kataegis clusters according to their contributions from CSR-like and SHM-like mutational processes with contributions of SHM and CSR as axes (analogous to Figure 2B, but without Kataegis cluster located inside IG loci). Kataegis clusters dominated by a CSR-like pattern are colored in orange, clusters dominated by a SHM-like pattern are colored in green and clusters dominated by neither pattern (other) are colored in purple.

Figure S3: Kataegis clusters and kataegis regions displayed as oncoprint, analogous to Figure 2C. The x-axis encodes samples, the y-axis the kataegis regions, which are ordered by recurrency of affection. The oncoprint carries 14 layers of annotation: i) the contributions per kataegis region of the three different categories SHM-like (green), CSR-like (orange) and other (purple) as stacked barplots; ii) a categorical annotation whether the kataegis region is located in VDJ genes (blue) or other parts of the IG loci (turquoise), whether it was described to be a target of aberrant SHM before (purple), whether the closest gene was associated with lymphomagenesis in some fashion (red), whether the closest gene is a cancer gene (orange) or all others (yellow); iii) barplots indicating the total number of SNVs in the respective kataegis region cohort-wide; iv) the density of SNVs per genomic window (SNVs per bp); v) genomic size of the respective kataegis region in bp; vi) distance to the closest transcription start site; vii) replication timing (encoded by RepliSeq scores) in lymphoblastoid cell lines; viii) integrated subject score of the respective kataegis region; ix) the integrated object score of the respective kataegis region; x) distribution of cellular fractions per kataegis regions; xi) distribution of distances to the closet breakpoint; xii) fractions of the different functional Annovar annotations; xiii) fractions of the different chromatin states from CG B cells; xiv) fractions of SNVs matching the motif DGYW in the respective kataegis regions.

Figure S4: Distribution of replication timing of SNVs in psichales clusters vs. SNVs in background as illustrated by combined boxplots and violin plots. (A) Replication time averaged over cell lines vs. psichales in the different entities (three breast cancer cohorts (BRCA), one cohort of chronic lymphocytic leukemia (CLL), one external DLBCL cohort, two cohorts of ovarian cancer (OV) and one sarcoma (SA) cohort (Nik-Zainal et al., Nature 2016; Nik-Zainal, Cell 2012; Puente et al., Nature 2015; Patch et al., Nature 2015; results shown here are in part based upon data generated by the TCGA research network: <http://cancergenome.nih.gov/>). (B) Replication time averaged over cell lines vs. psichales across all entities and samples. (C) Comparison of replication time in the breast cancer cell line MCF7 (left) against replication time averaged over cell lines (right) vs. psichales in the breast cancer entities. (D) Comparison of replication time in the lymphoblastic cell lines GM06990, GM12801, GM12812, GM12813 and GM12878 against replication time averaged over cell lines (right) vs. psichales in the B-cell derived entities. The label "TRUE" represents the psichales SNVs, the label "FALSE" the non-psichales SNVs. (E) Distribution of replication timing of SNVs in kataegis clusters vs. SNVs in background as illustrated by combined boxplots and violin plots - in contrast to psichales, kataegis is not enriched in late replicating genomic regions.

Figure S5: Psichales clusters and psichales regions displayed as oncoprint, analogous to Figure S3. 18,794 psichales clusters were detected in 219 genomes. Using a recurrence threshold of three affected cases, 252 recurrent psichales regions were identified. The x-axis encodes samples, the y-axis the psichales regions, which are ordered by recurrency of affection. The oncoprint carries 12 layers of annotation analogous to Figure S3, but without colour coding of the individual psichales clusters and without integrated subject or object scores (all not defined).

Figures S6/S7: Genomic windows in which the IGH switch regions were defined (in the order of their genomic coordinates): **(A)** IGHA2 (switch $\alpha 2$), **(B)** IGHE (switch ϵ), **(C)** IGHG4 (switch $\gamma 4$), **(D)** IGHG2 (switch $\gamma 2$), **(E)** IGHA1 (switch $\alpha 1$), **(F)** IGHG1 (switch $\gamma 1$), **(G)** IGHG3 (switch $\gamma 3$), **(H)** IGHM (switch μ). Each subplot is composed of several tracks, showing from bottom to top (if a track is empty it is left out in the respective subplot): (i) transcripts from ENSEMBL; (ii) repeats from repeat masker (with colour coding Alu – light blue, L1 – blue, L2 – light green, L3 – green, L4 – pale red, LINE – light red, LTR – ocre, MER – orange, MIR – pale purple, MLT – purple, RNA – pale yellow, SINE – brown, rich – grey50, oligomer – grey30, monomer – black, other – grey80, simple

- red); (iii) merged repeats whenever overlapping (irrespective of repeat class); (iv – vi) tracks showing alignments with contigs from the IMGT data base where *oGL* = „original global-local“ (global-local alignments of contigs known to originate from one single switch region), and *oL* = „original local“ (local Smith-Waterman alignments of contigs known to originate from one single switch region) and *jL* = „junction local“ (local alignments of contigs known to originate from recombinations); (vii) a track for Ka-ROIs (regions of interest of recurrent Kataegis-like clusters); (viii – xii) tracks showing hallmark events (translocations between IGH and *MYC*, *BCL2* or *BCL6*) in the different subgroups, where the subgroups „BL_solid“, „BL_leukemia“ and „BL_pleura“ are grouped together and designated „BL“ and the subgroups „DHL“ and „B-NOS“ are grouped together and designated „other“. All subgroup specific tracks use a common colour coding for the hallmark events: *MYC* - blue, *BCL2* - dark green, *BCL6* – brown; (xiii) track showing non-hallmark translocations; and (xiv) track showing coordinates of breakpoints of rearrangements within the IGH locus.

Figure S8: Mutational mechanisms in gcBCL subgroups, comparing DLBCL, FL-DLBCL and FL (left part of the figure, items A, B, E, F, I, J, M and N) and ABC-DLBCL vs. GCB-DLBCL (right part of the figure, items C, D, G, H, K, L, O, P and R). (A) Absolute counts and (B) normalized fractions of CSR-like kataegis clusters among all kataegis clusters in gcBCL subgroups. (E) Absolute counts and (F) normalized fractions of SHM-like kataegis clusters among all kataegis clusters in gcBCL subgroups. (I) Absolute counts and (J) normalized fractions of kataegis clusters of type “other” among all kataegis clusters in gcBCL subgroups. (C) Absolute counts and (D) normalized fractions of CSR-like kataegis clusters among all kataegis clusters in ABC-DLBCL vs. GCB-DLBCL. (G) Absolute counts and (H) normalized fractions of SHM-like kataegis clusters among all kataegis clusters in ABC-DLBCL vs. GCB-DLBCL. (K) Absolute counts and (L) normalized fractions of kataegis clusters of type “other” among all kataegis clusters in ABC-DLBCL vs. GCB-DLBCL. (M) Absolute counts of HbP instances per sample in gcBCL subgroups. (N) Number of HbP instances per sample normalized to the square of the number of kataegis clusters per sample in gcBCL subgroups. (O, P) Analogous to (M) and (N) but for ABC-DLBCL vs. GCB-DLBCL. (Q) Dotplot with linear fit between (on the x-axis) the number of kataegis clusters per sample and (on the y-axis) the square root of the number of HbP instances per sample. The linear fit in this display illustrates a quadratic relationship, i.e. $y = c \cdot x^2$. (R) Enrichment and depletion patterns of mutational signatures in ABC-DLBCL vs. GCB-DLBCL. Abbreviations: HbP: hypermutation by proxy.

Figure S9: Hypermutation by proxy: overview. Pairwise subject-object relationships of the kataegis regions as determined by the custom method introduced in this work. All kataegis regions identified in gcBCL are included in the plot. Subject regions are displayed as rows and object regions as columns. Note that this visualization is not symmetric.

Figure S10: Hypermutation by proxy: examples. (A-I) HbP at the *BCL6* gene cluster. (A) Oncoprint showing the samples affected by rKa-clusters (recurrent Kataegis-like clusters) in *BCL6* and the neighboring genes involved in HbP (*RP11-132N15.3*, *RP11-44H4.1*, *LPP*, *ST6GAL1*). Whereas three different Ka-ROIs are located within *LPP*, the other genes are affected by only one Ka-ROI each, leading to a total of six Ka-ROIs involved in the *BCL6* gene cluster. Four out of six Ka-ROIs overlap with chromatin segments which are annotated as enhancers. (B - E) Scatterplots of the expression of the genes *RP11-132N15.3*, *RP11-44H4.1*, *LPP*, *ST6GAL1* as compared to *BCL6* per subgroup. **(F-I) HbP around *PAX5*.** (F) Oncoprint showing the samples affected by rKa-clusters (recurrent Kataegis-like clusters) in *PAX5* and the neighboring genes involved in HbP (*ZCCHC7*, *GRHPR* and *RP11-397D12.4*). Whereas two different Ka-ROIs are located within *ZCCHC7*, the other genes are affected by only one Ka-ROI each, leading to a total of four Ka-ROIs involved

in genes in the vicinity of *PAX5*. Three out of four Ka-ROIs overlap with chromatin segments which are annotated as enhancers. (G - I) Scatterplots of the expression of the genes *ZCCHC7*, *GRHPR* and *RP11-397D12.4* as compared to *PAX5* per subgroup.

Figure S11: Unsupervised analysis of mutational signatures with NMF. (A) Riverplot graph, in which nodes correspond to mutational signatures and edges represent pairwise cosine similarities between mutational signatures. The vertical direction encodes different factorization ranks, increasing from top to bottom. The horizontal direction shows the different mutational signatures extracted at the respective factorization rank. The horizontal position of one node (corresponding to one mutational signature) is chosen such that crossing edge weights are minimized. Mutational signatures are colored and labeled according to their cosine similarity to known mutational signatures (from <http://cancer.sanger.ac.uk/cosmic/signatures>). At the optimal factorization rank $k = 11$ (cf. subplot (B)), three new signatures can be found. These new signatures are taken into account in the coloring and labeling of the entire graph. (B) Quality metrics for the choice of the factorization rank: norm of the residuals (“FrobError” for Frobenius reconstruction error, to be minimized), cophenetic correlation coefficient (“copheneticCoeff”, to be maximized) and the mean Amari distance (“meanAmariDist”, to be minimized). Combining these three measures, an optimal factorization rank of 11 is evaluated. (C) Determination of the cutoffs for the stratified analysis of mutational signatures with replication timing (encoded by RepliSeq score) as stratification axis. The vertical red lines show the chosen cutoffs in the density function of the RepliSeq scores of all SNVs of our cohort.

Figure S12: Unsupervised analysis of mutational signatures with NMF. (A) Exposures to the identified mutational signatures for 219 lymphomas (179 cases with matched normal from this study and additional 40 pediatric Burkitt lymphomas (BL), 39 from [Lopez et al., manuscript accepted at Nature Communications] and one adult BL). (B) Overlaid violin plots and boxplots of the distributions of cancer cell fractions (CCF) of SNVs in IG const genes, in IG switch regions (where the CSR-like profile stems from) and in IG VDJ genes (where the SHM-like profile stems from). (C) Normalized exposures to the identified mutational signatures. (D - G) More enrichment and depletion patterns as indicated in the titles of the subplots. (D) Stratified analysis of mutational signatures showed an enrichment of AC1 (spontaneous deamination, $p_{KW} = 8.8 \times 10^{-23}$, $p_{Nem} = 5.8 \times 10^{-13}$) and AC2 (APOBEC, $p_{KW} = 2.2 \times 10^{-3}$, $p_{Nem} = 9.7 \times 10^{-3}$) in early clonal evolution. L1 showed a trend towards enrichment in early ($p_{KW} = 9.6 \times 10^{-2}$) and AC9 in late ($p_{KW} = 1.6 \times 10^{-2}$, $p_{Nem} = 0.2$) clonal evolution. (G) L1 and L2 were enriched in kataegis clusters, with L1 being enriched in CSR-like ($p_{KW} = 9.8 \times 10^{-21}$, $p_{Nem} = 1.7 \times 10^{-5}$) and L2 in SHM-like ($p_{KW} = 4.4 \times 10^{-35}$, $p_{Nem} = 5.4 \times 10^{-4}$) kataegis clusters as compared to the non-clustered SNV stratum. Abbreviations: CCF – cancer cell fraction, id. – intermutation distance. Error bars represent standard error of the mean (SEM).

Figure S13: Occurrences of the most frequently mutated genes. (A) In DLBCL vs. FL and (B) in ABC-DLBCL vs. GCB-DLBCL.

Figure S14: Extended SNV statistics for the driver genes. To account for the broad range of mutation frequencies in the different driver genes, the visualization is split into three groups with different axis scales. From left to right, the panels show: i) the number of samples with any SNVs in the respective gene or in a 2.5 kb window around the gene; ii) the number of kataegis clusters overlapping with the respective gene in the gcBCL cohort; iii) the distribution of SNV counts per sample; iv) the distribution of the distance of each SNV to the TSS (of the canonical transcript); v) fractions of SNVs belonging to SHM-like kataegis clusters (green), CSR-like kataegis clusters (orange), kataegis clusters of type “other” (purple) and no kataegis clusters at all (yellow); vi) the

distribution of the position of the SNVs relative to the (canonical) TSS of the respective genes: outside of a 2.5 kb window 5' of the gene (light blue), inside a 2.5 kb window 5' of the gene (dark blue), within 2.5 kb from TSS in the gene body (dark green), more than 2.5 kb away of the TSS in the gene body, and 3' of the gene (light red); and vii) the ratio of the SNV density in the window 2.5 kb around the (canonical) TSS divided by the SNV density in the remaining gene.

Figure S15: NMF-based consensus clustering of alterations in driver genes identified in DLBCL cases of our gcBCL cohort. (A) Four clusters were identified and labelled "MYD88-like" (red), "BCL2-like" (green), "BCL6-like" (light blue) and "TP53-like" (purple) and annotated at the bottom of the heatmap. The cell-of-origin attribution of the DLBCL cases is also annotated below the heatmap. All alterations are encoded as red squares in affected genes. Inside the clusters, affected genes are ordered by occurrence of alterations. (B) Assessment of SNV load for the cases belonging to the different clusters. (C) Fraction of the genome with aberrant copy number state per consensus cluster.

Figure S16: NMF-based consensus clustering of alterations in driver genes identified in the gcBCL cohort across sub-entities. (A) Nine clusters were identified and labelled "TP53-like" (red), "B2M-like" (orange), "PIM1-like" (light green), "BCL6-like" (green), "PAX5-like" (turquoise), "BCL2-like" (light blue), "CSMD1-like" (blue), "SOCS1"-like (purple), and "MYD88-like" (magenta), and annotated at the bottom of the heatmap. The cell-of-origin attribution as well as the subgroup/entity of the samples are also annotated below the heatmap. The three hypermutated samples 4145528, 4109808 and 4163639 were excluded from the whole cohort of 181 gcBCL. All alterations are encoded as red squares in affected genes. Inside the clusters, affected genes are ordered by occurrence of alterations. (B) Mutational load in the different consensus clusters. (C) Fraction of the genome with aberrant copy number state.

Figure S17: Distributions of cancer cell fraction per driver gene.

Figure S18: Tandem RNA chimeras between S1PR2 and DNMT1 in the DLBCL cell line SU-DHL-10. A variety of splice variants were identified by PCR.

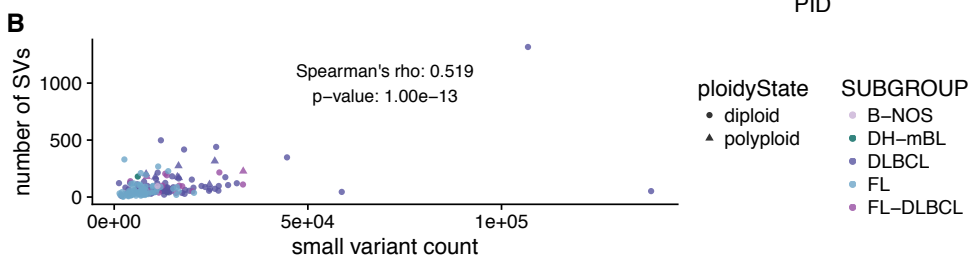
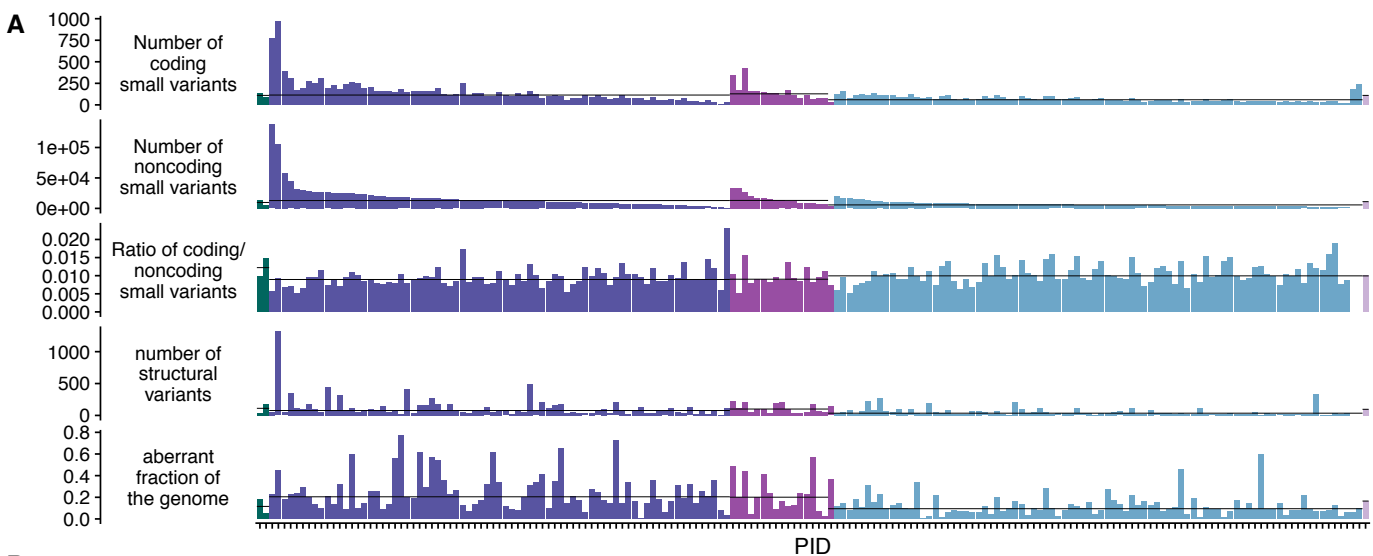
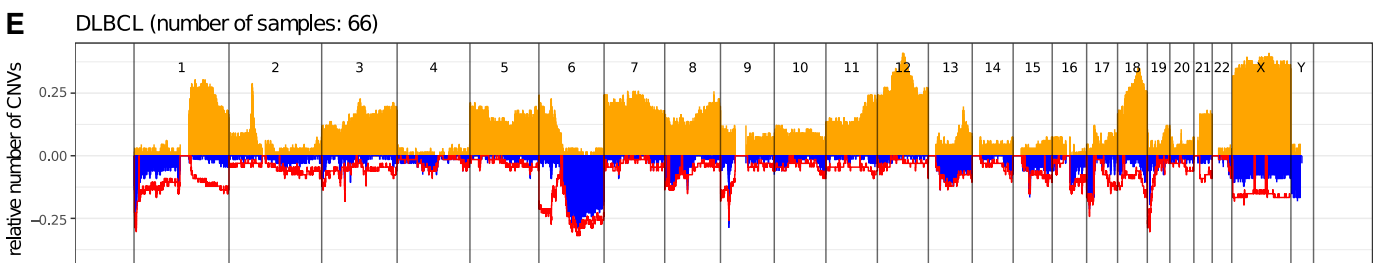
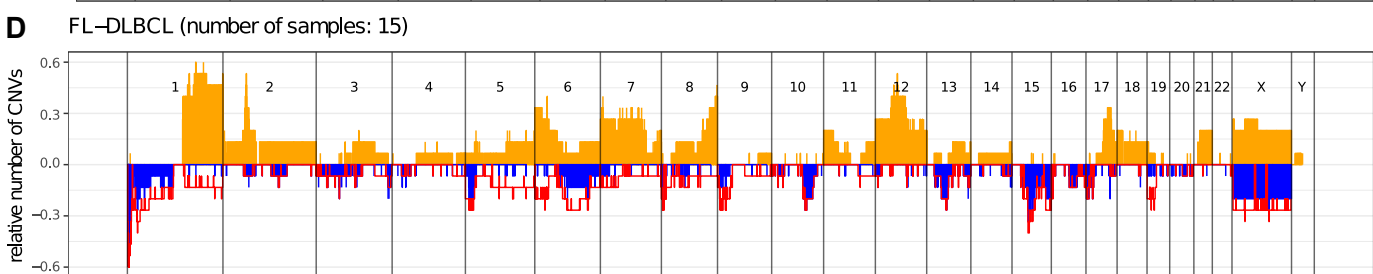
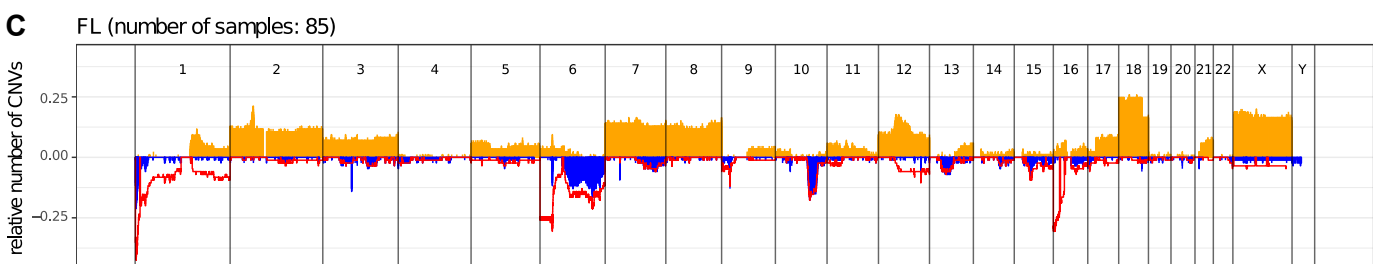


Fig. S1



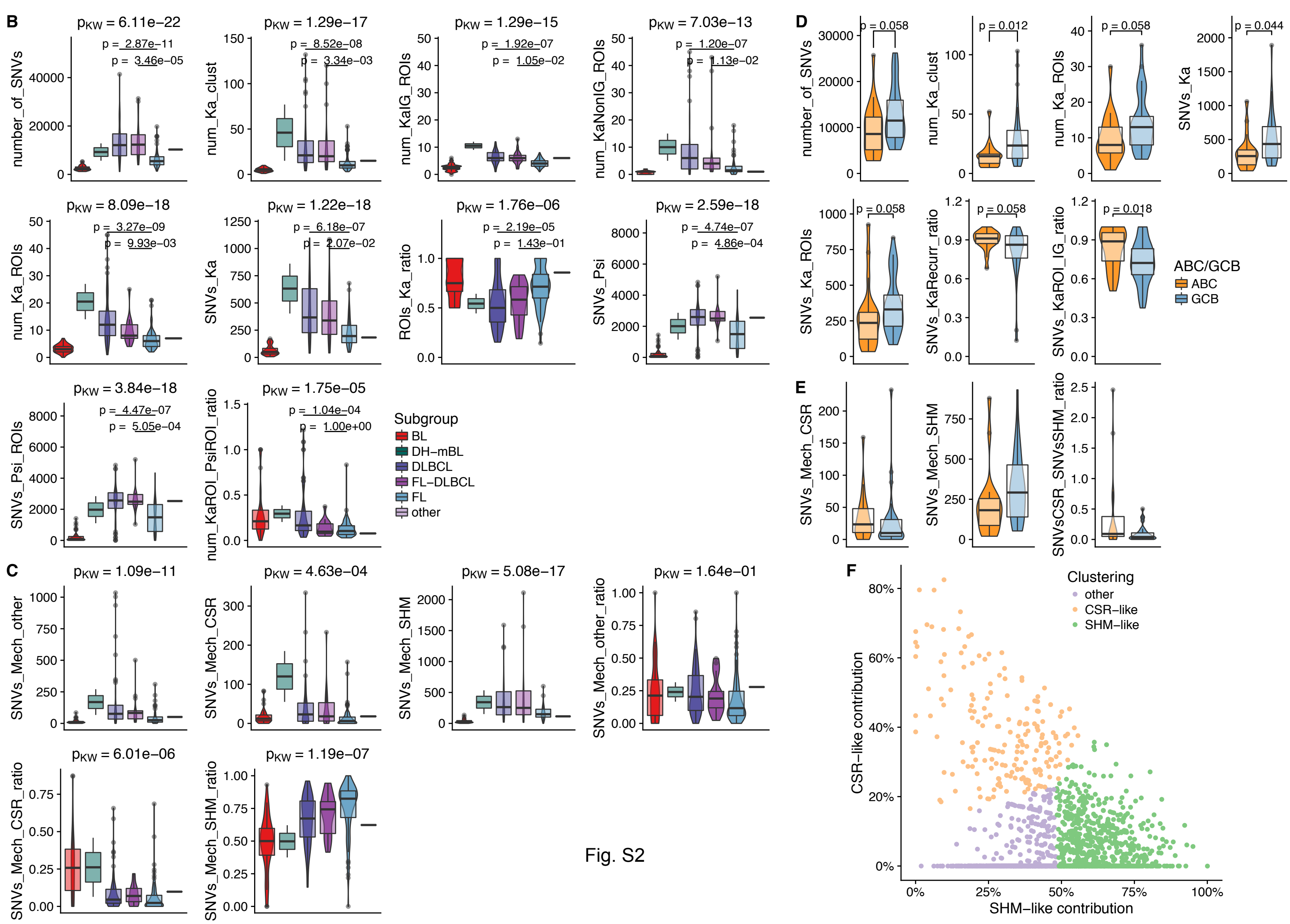
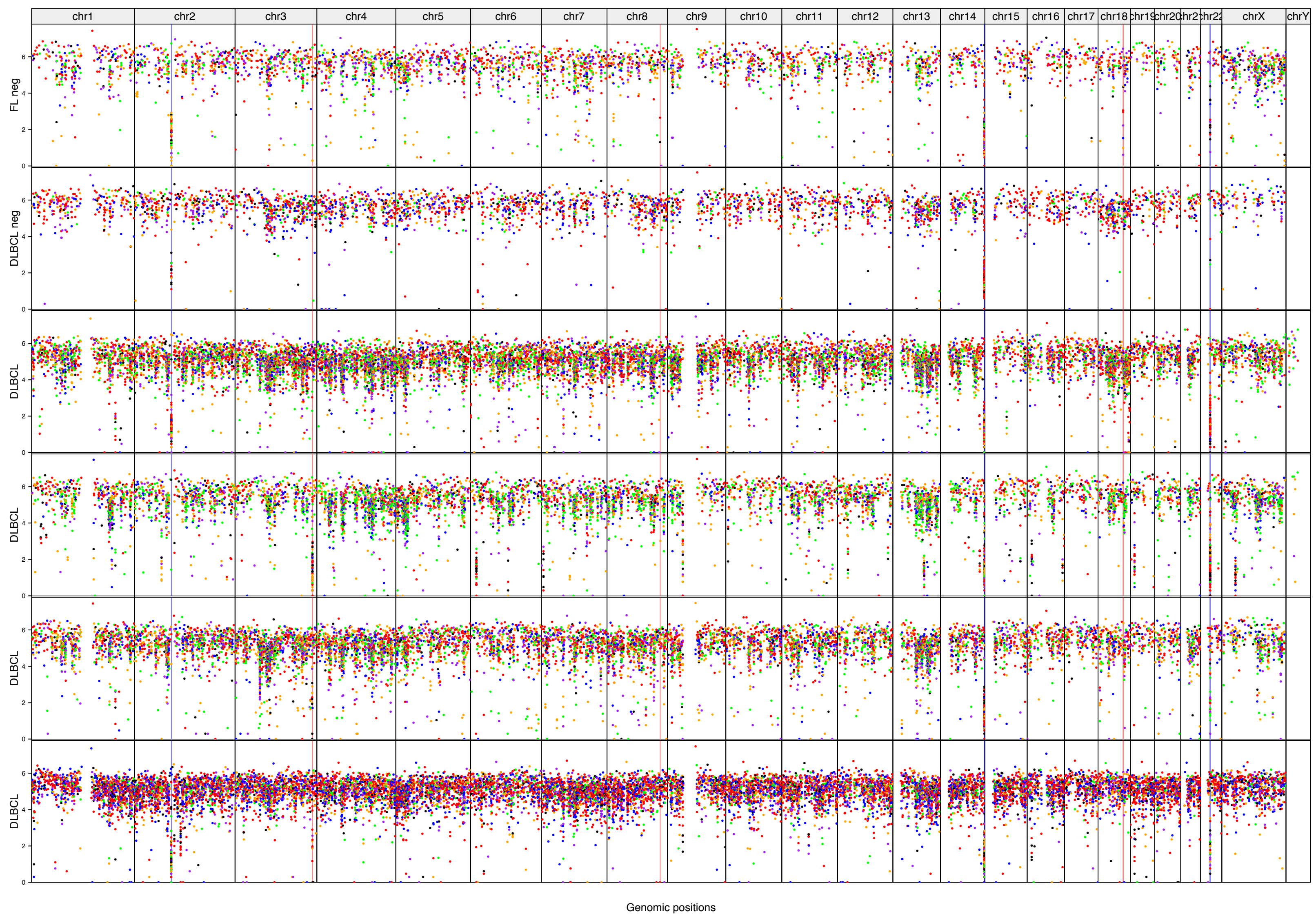


Fig. S2



Fig. S3

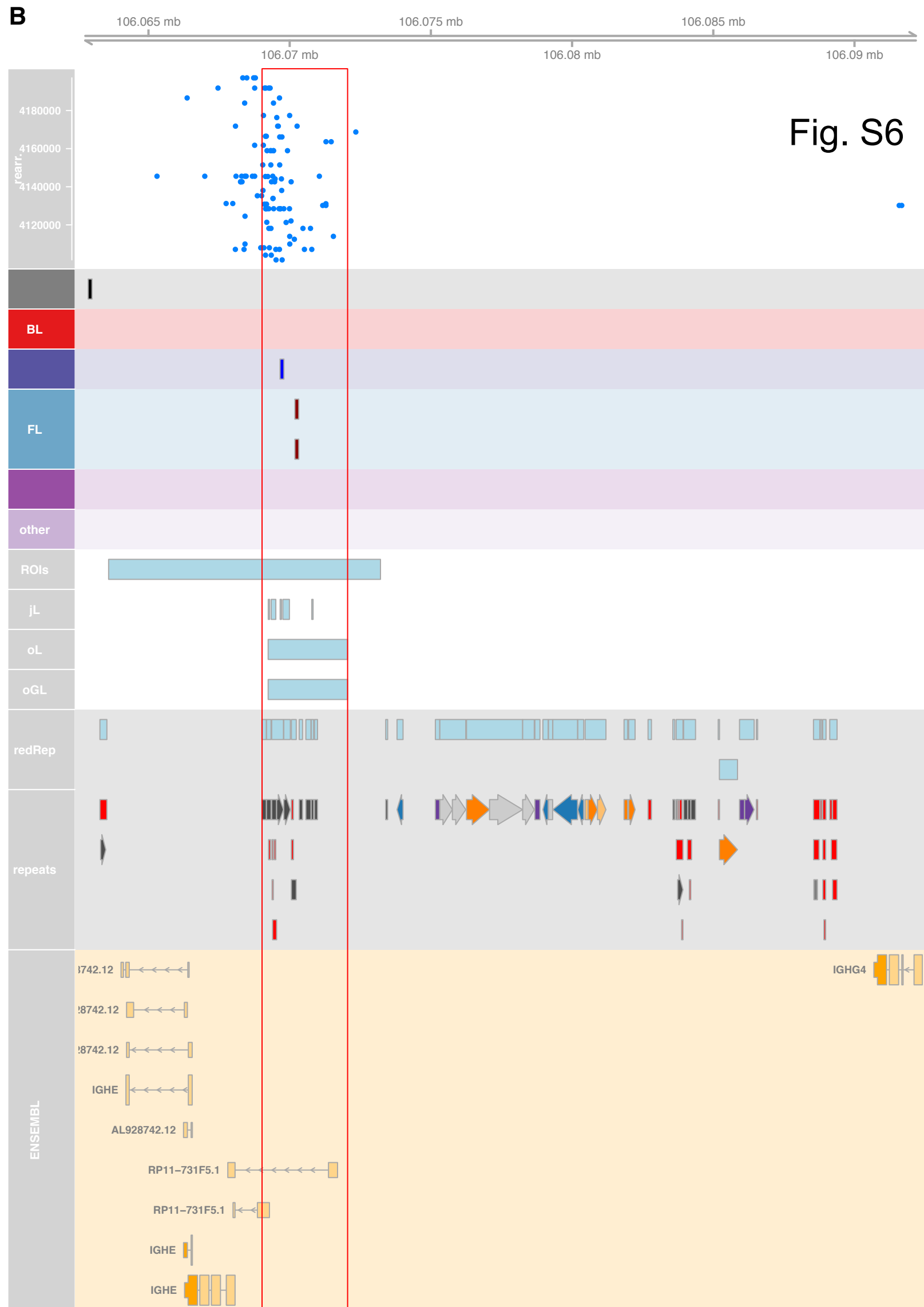
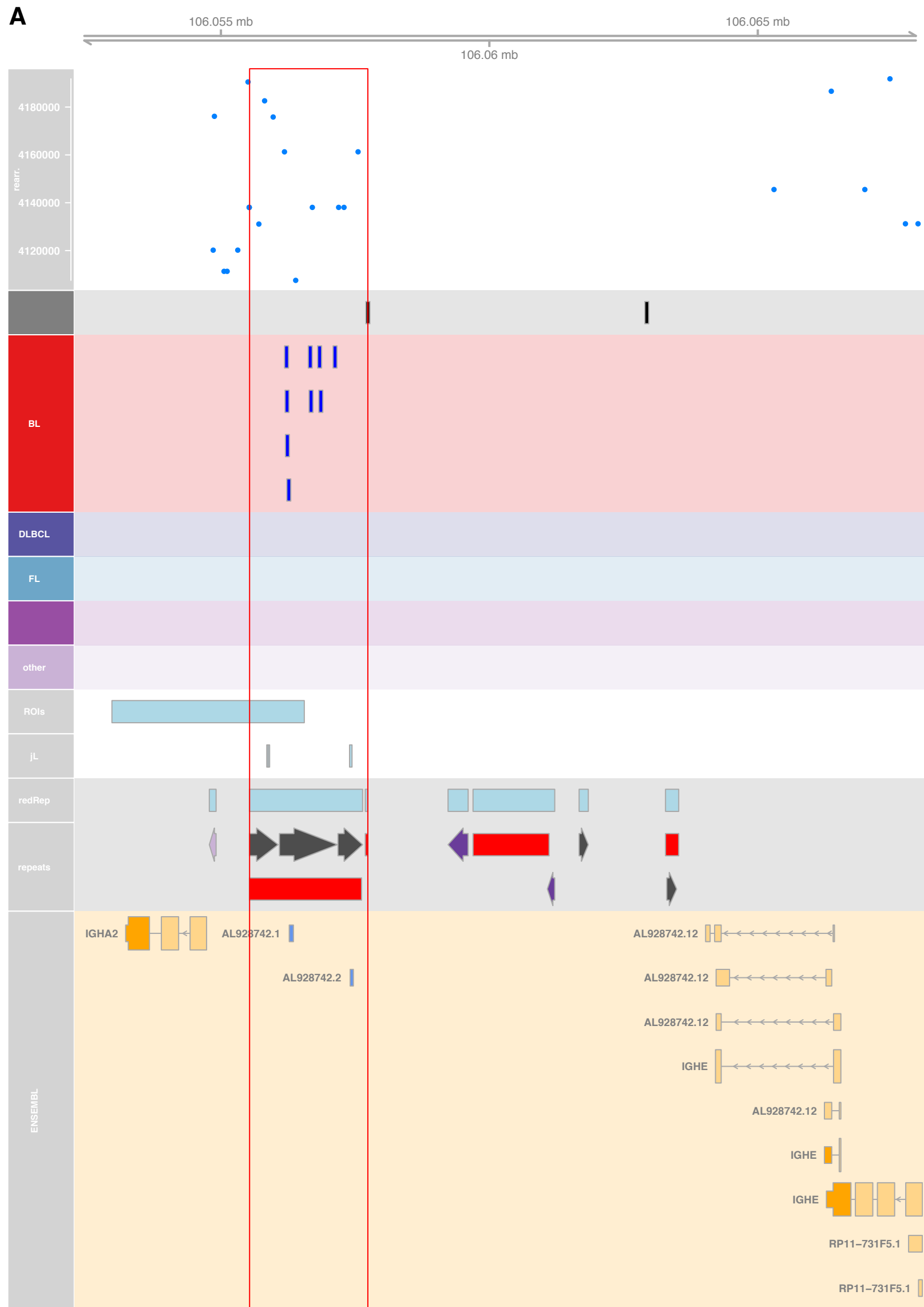
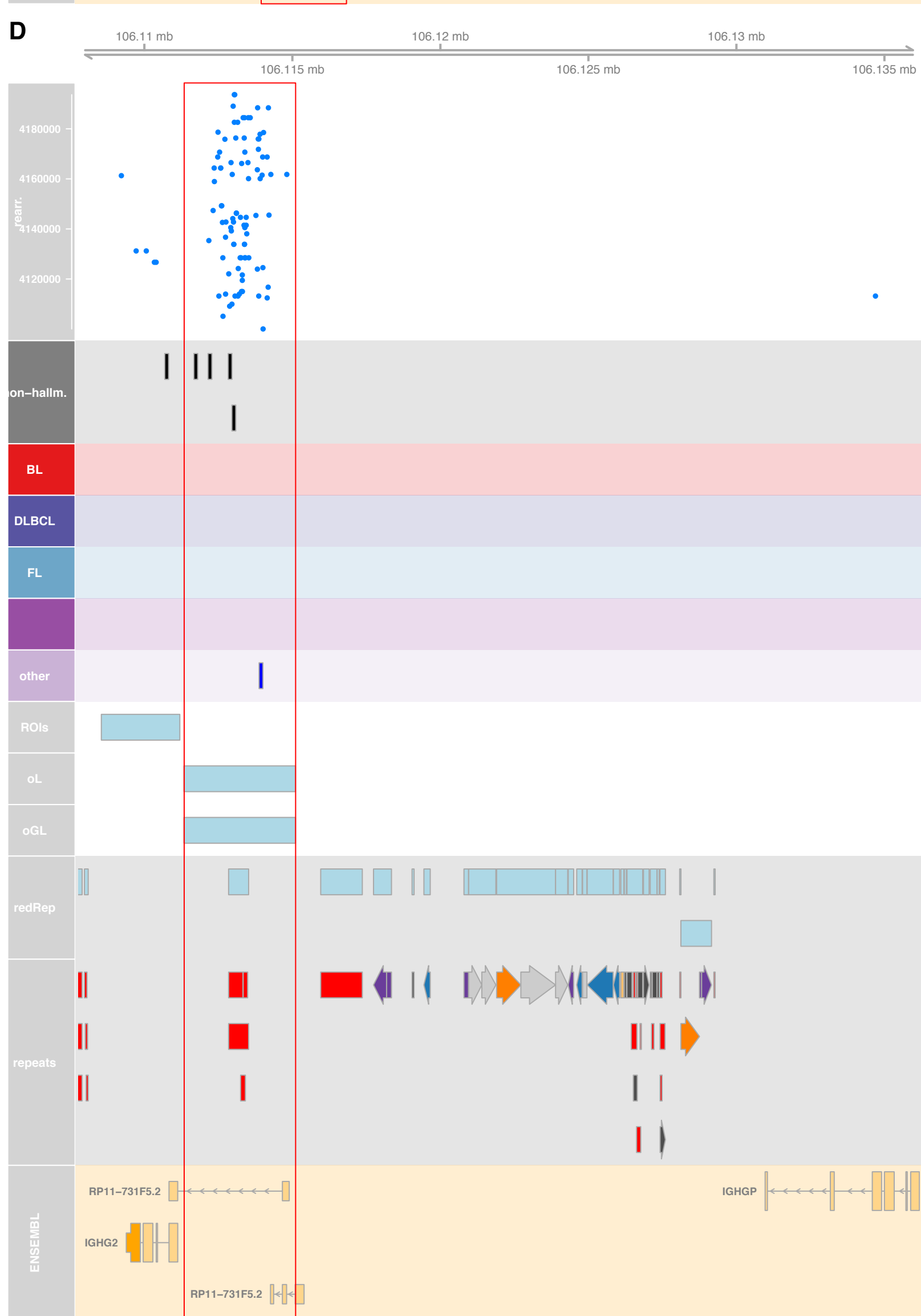
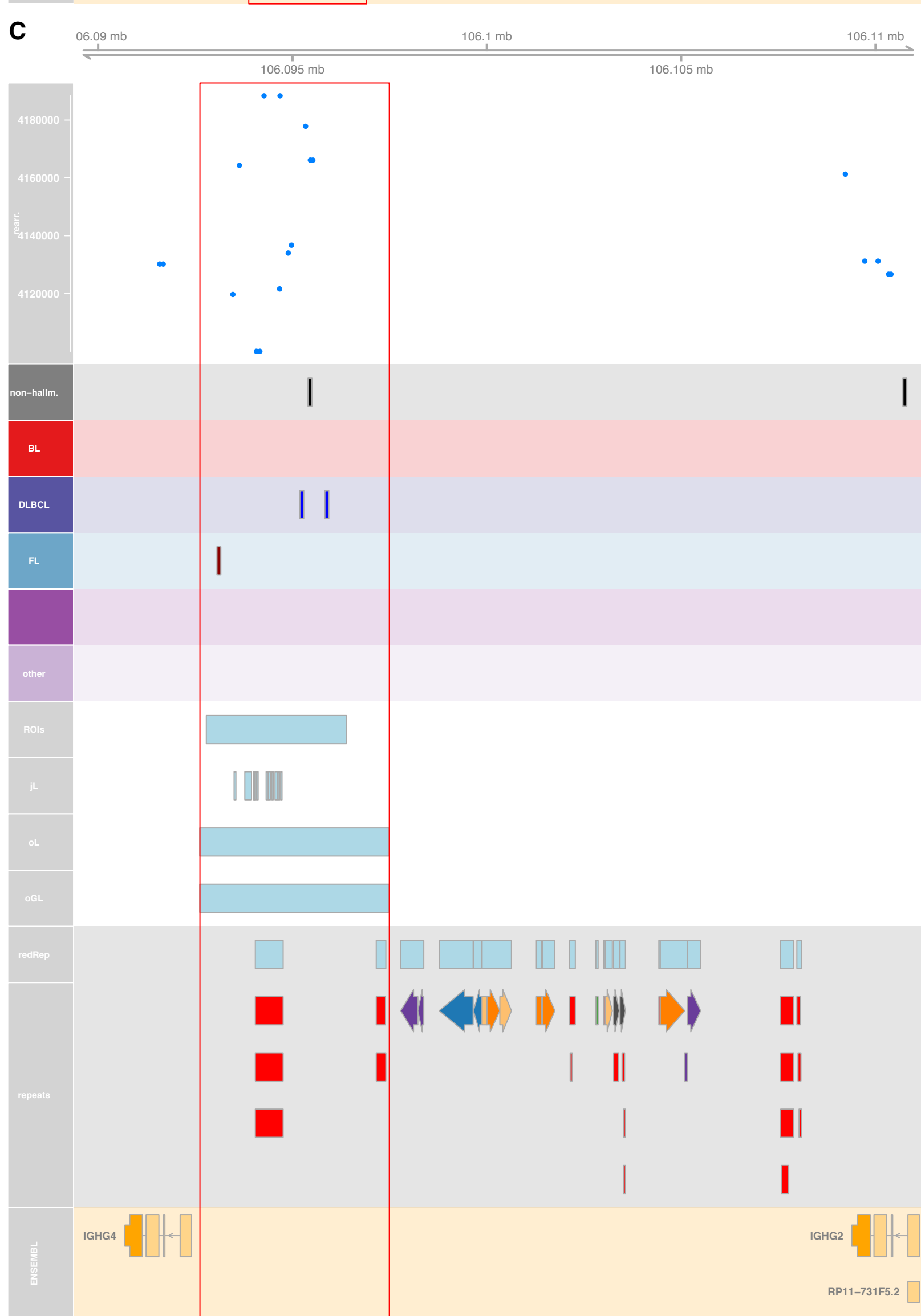


Fig. S6



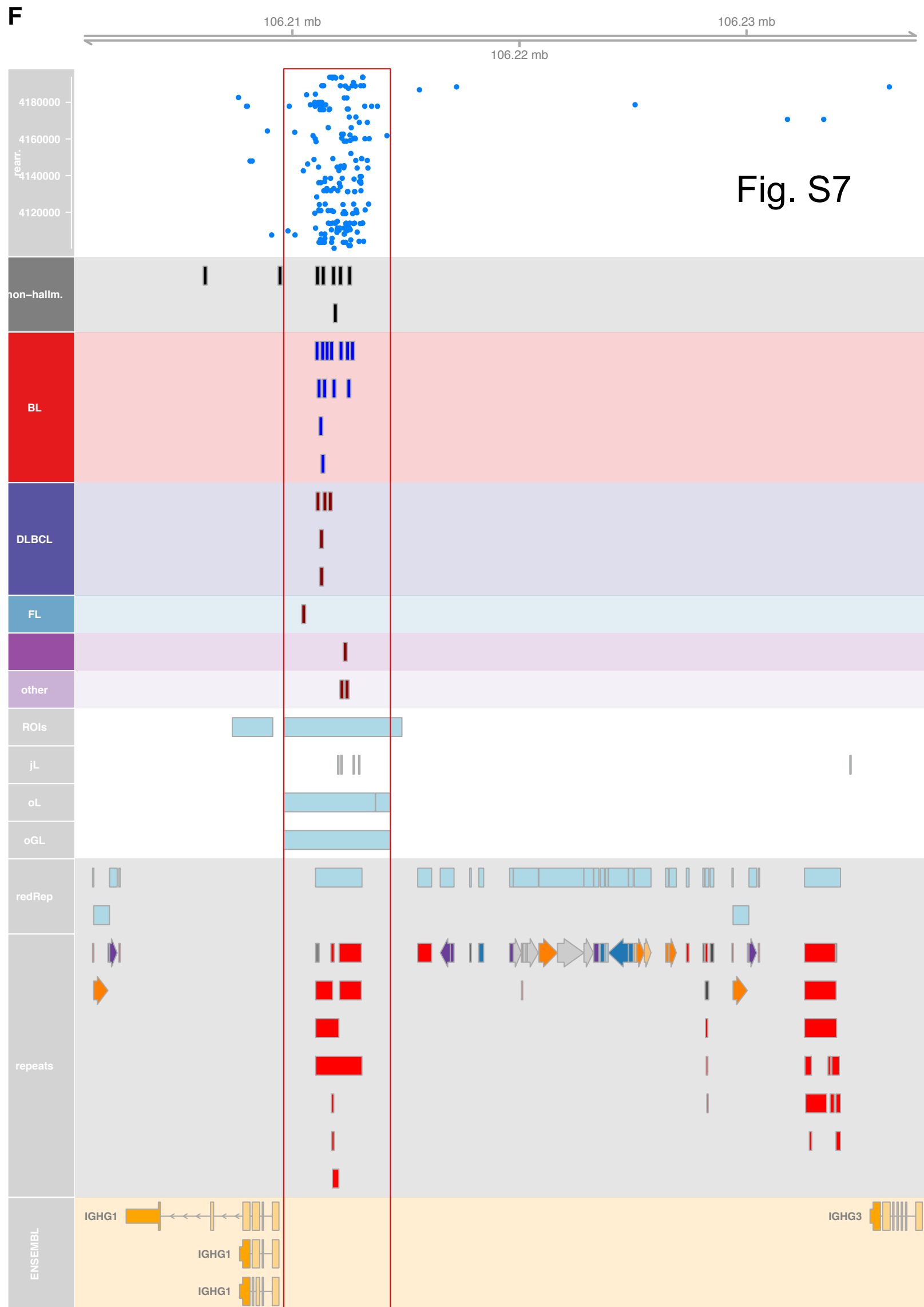
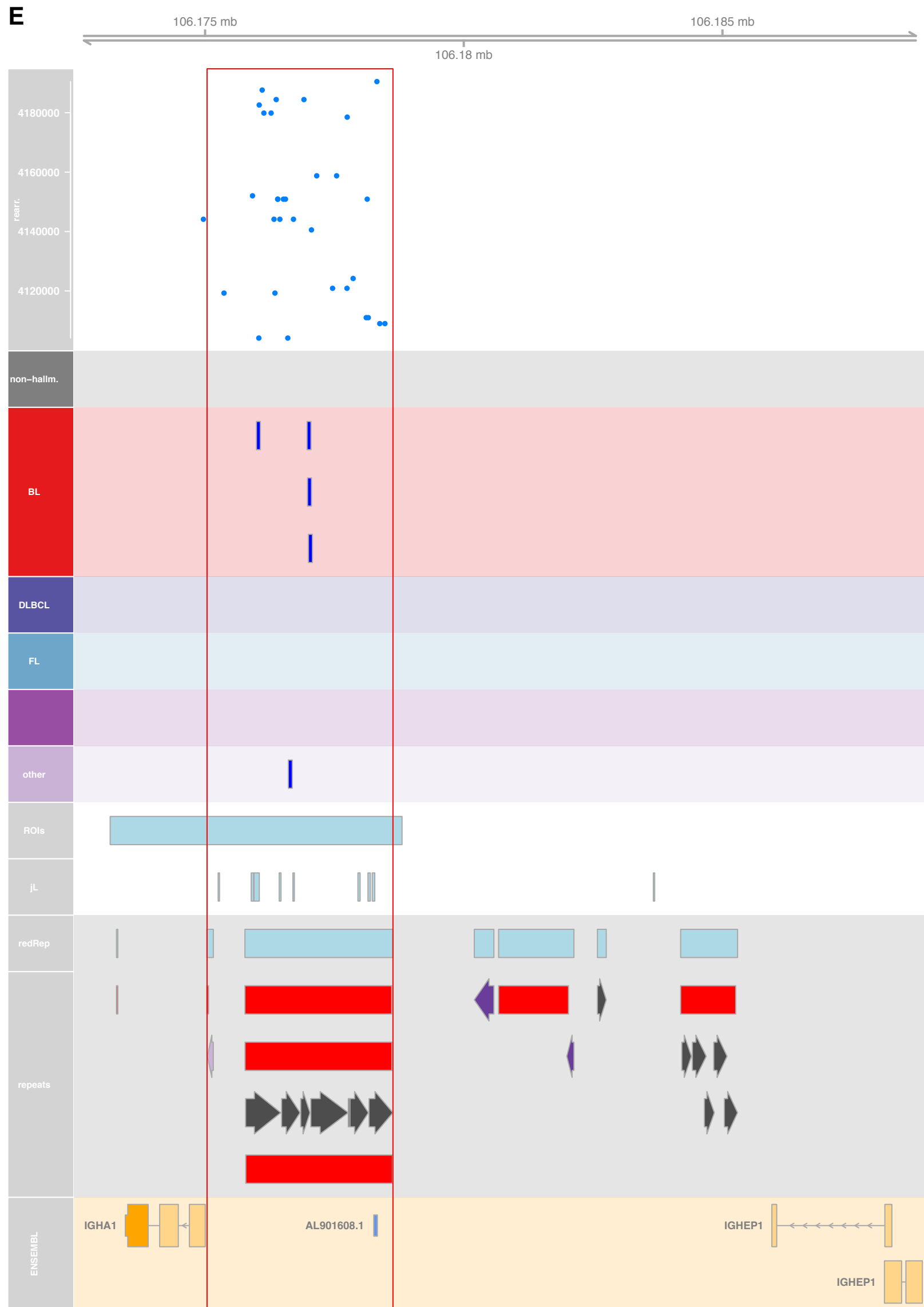
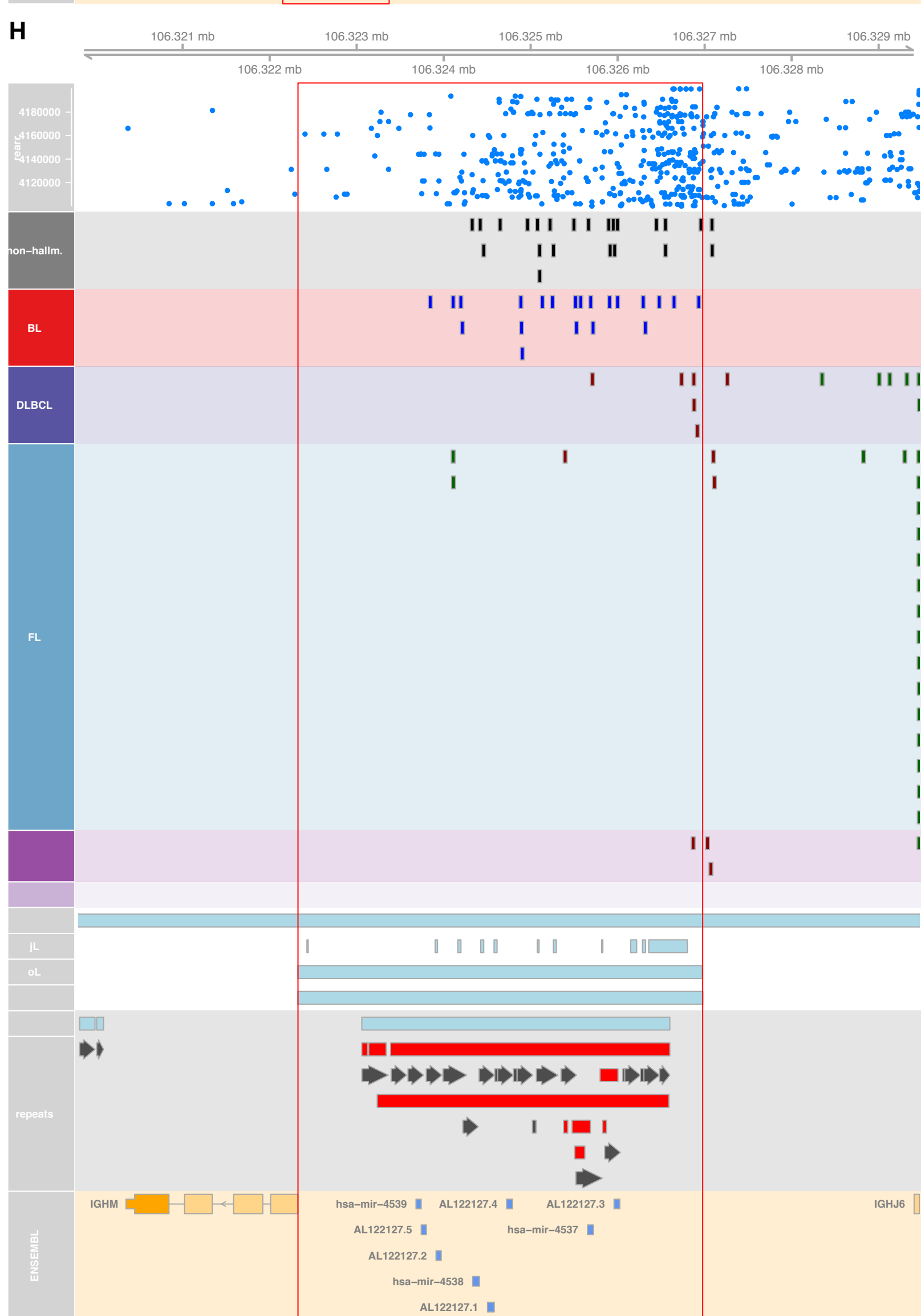
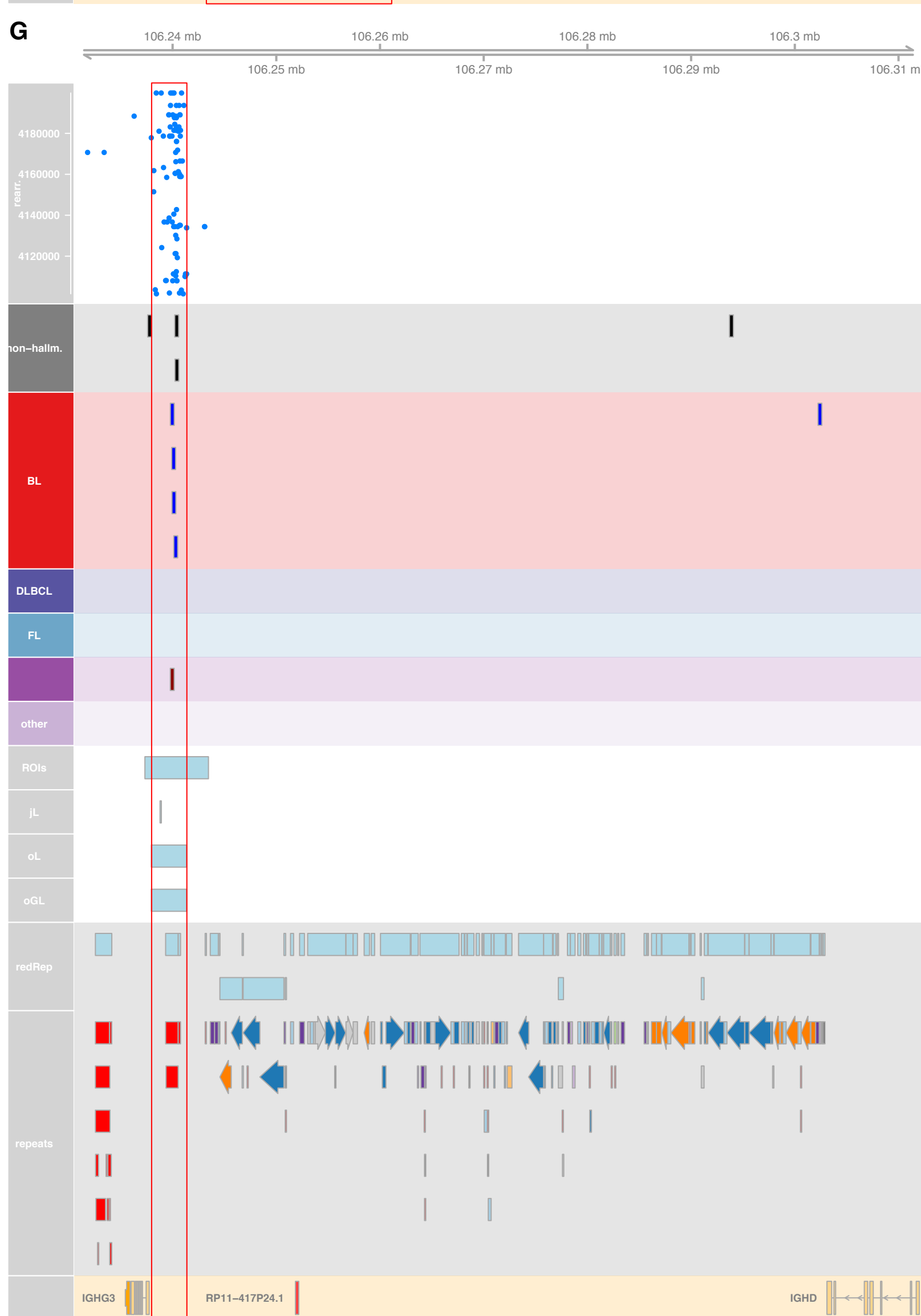


Fig. S7



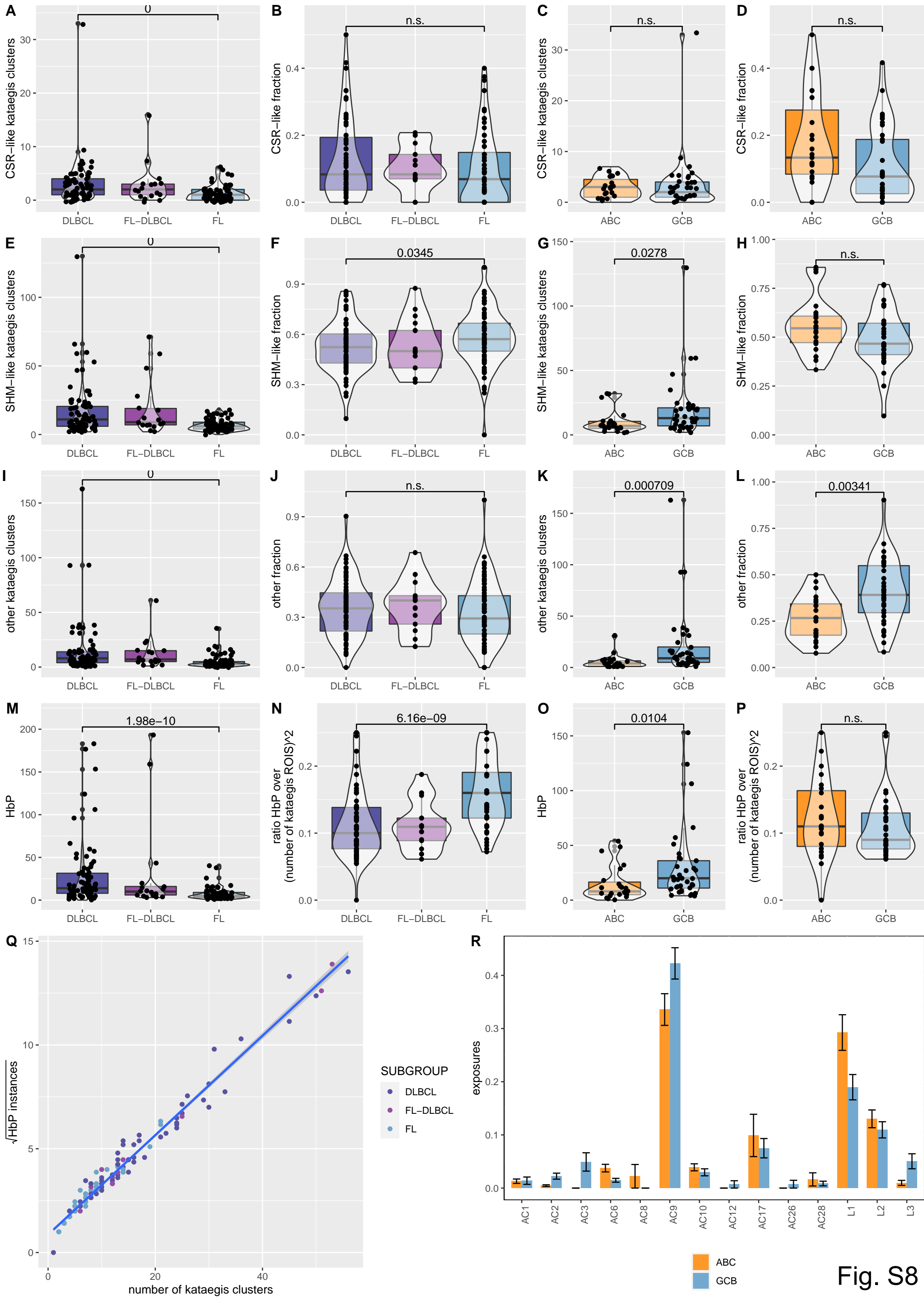


Fig. S8

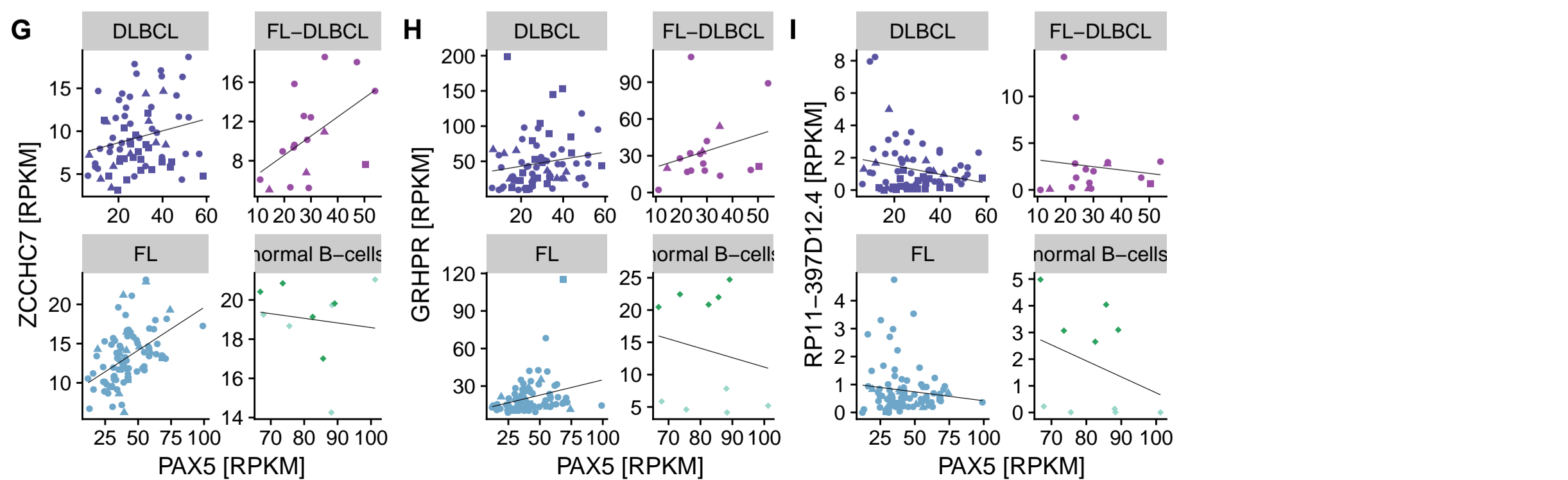
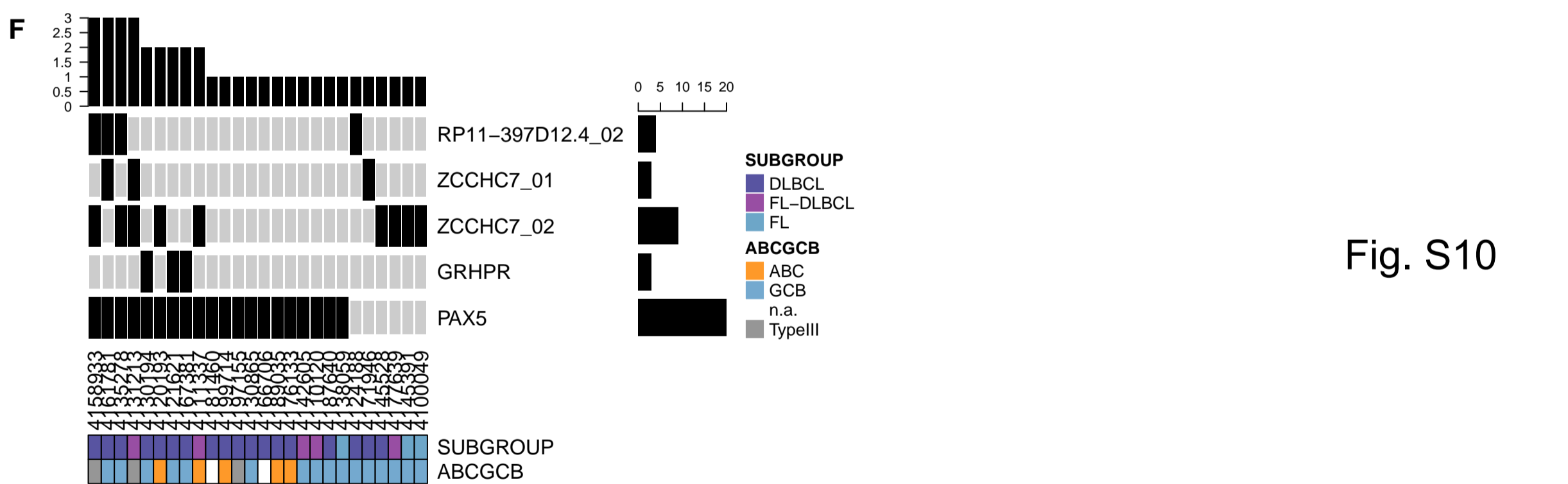
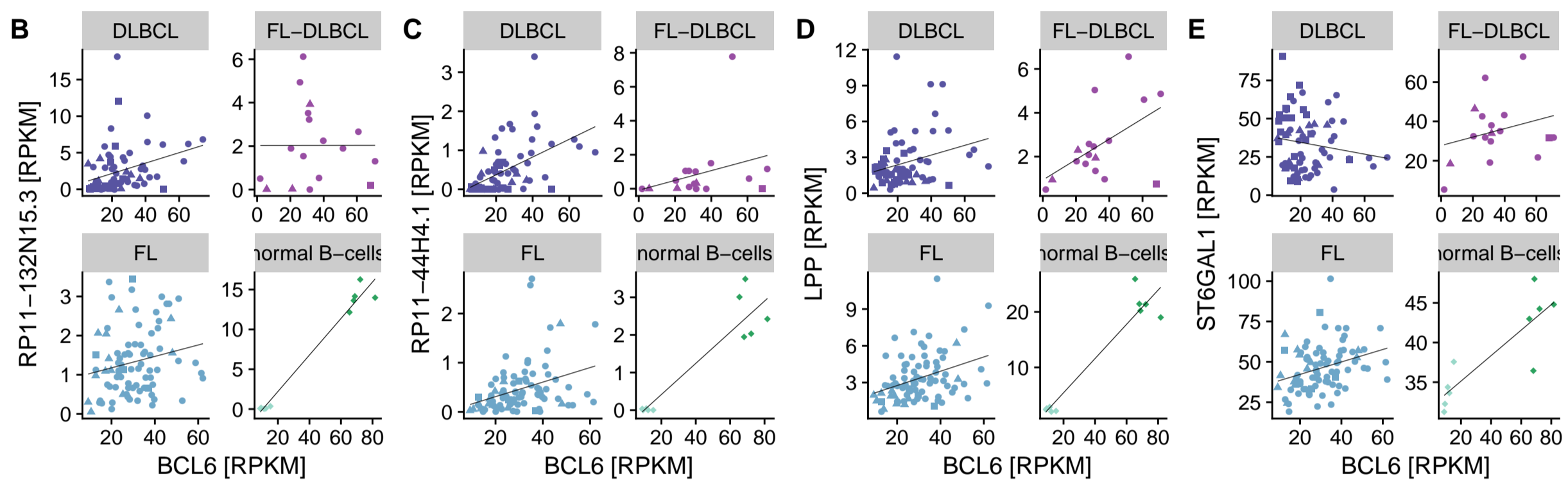
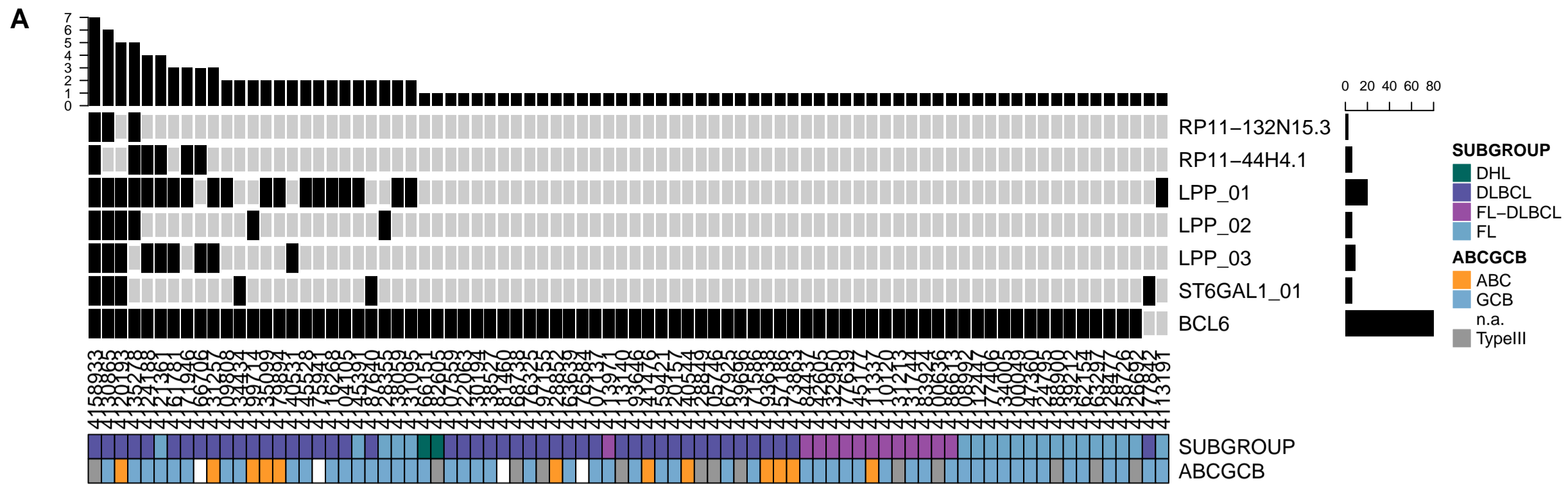
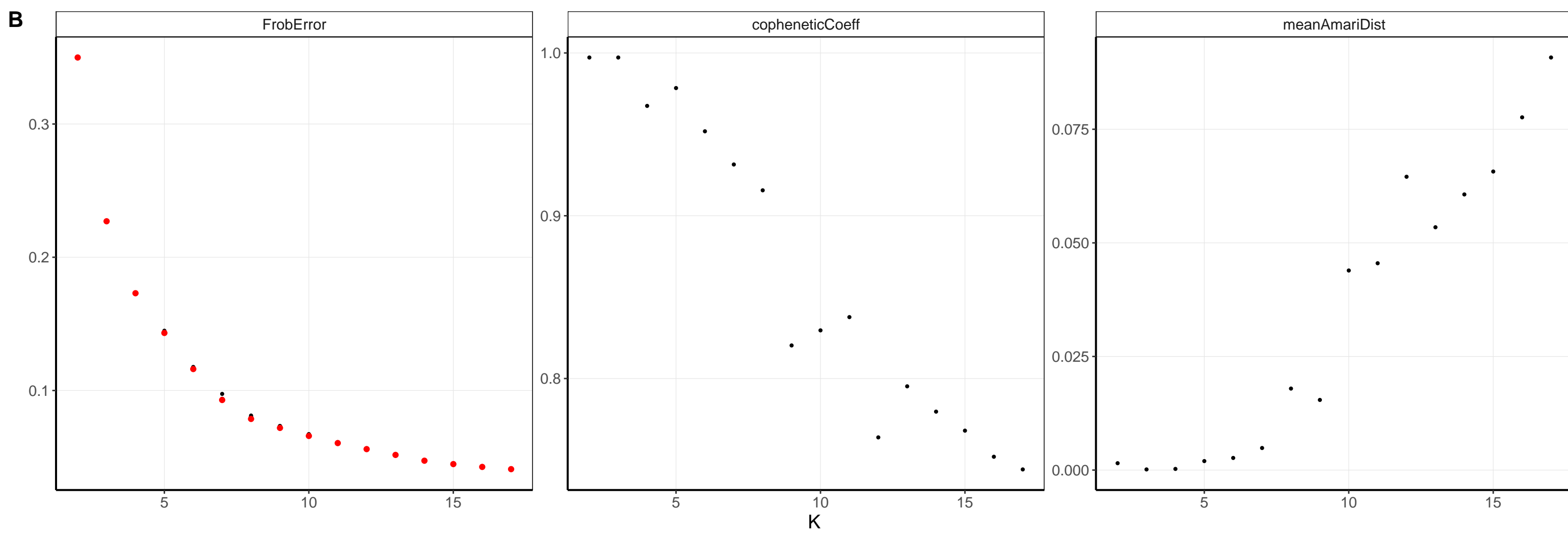
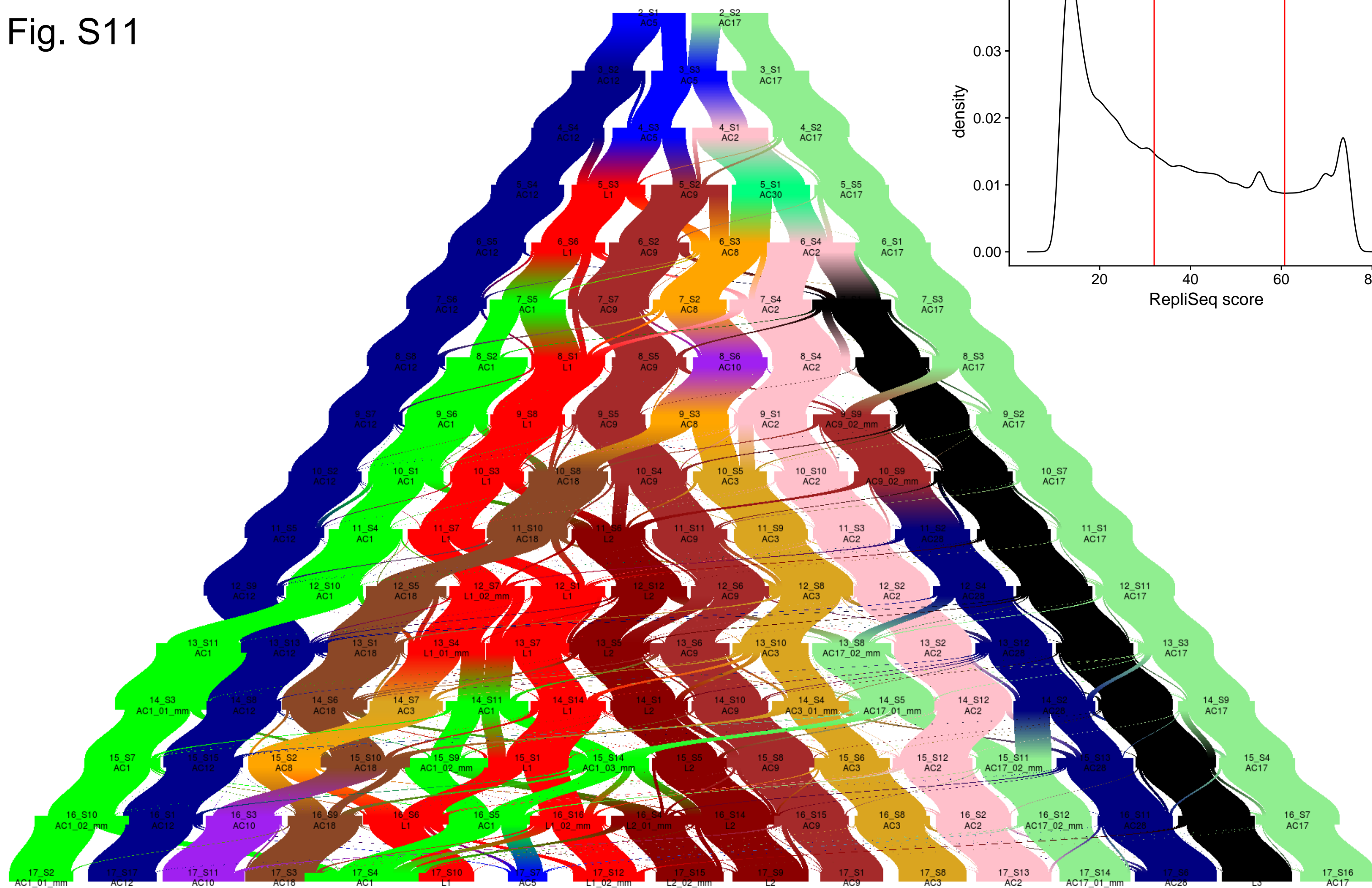
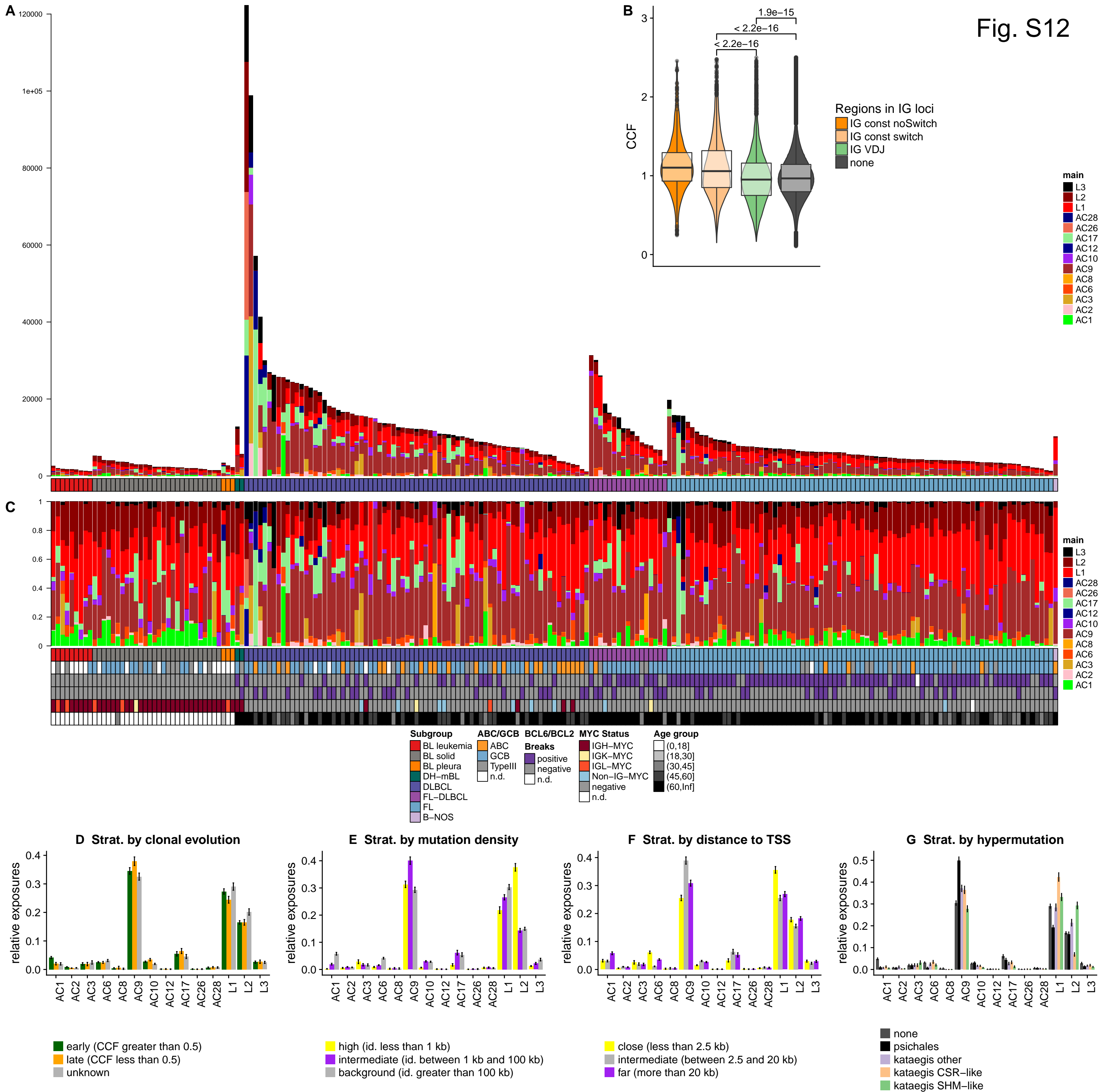


Fig. S10

Fig. S11





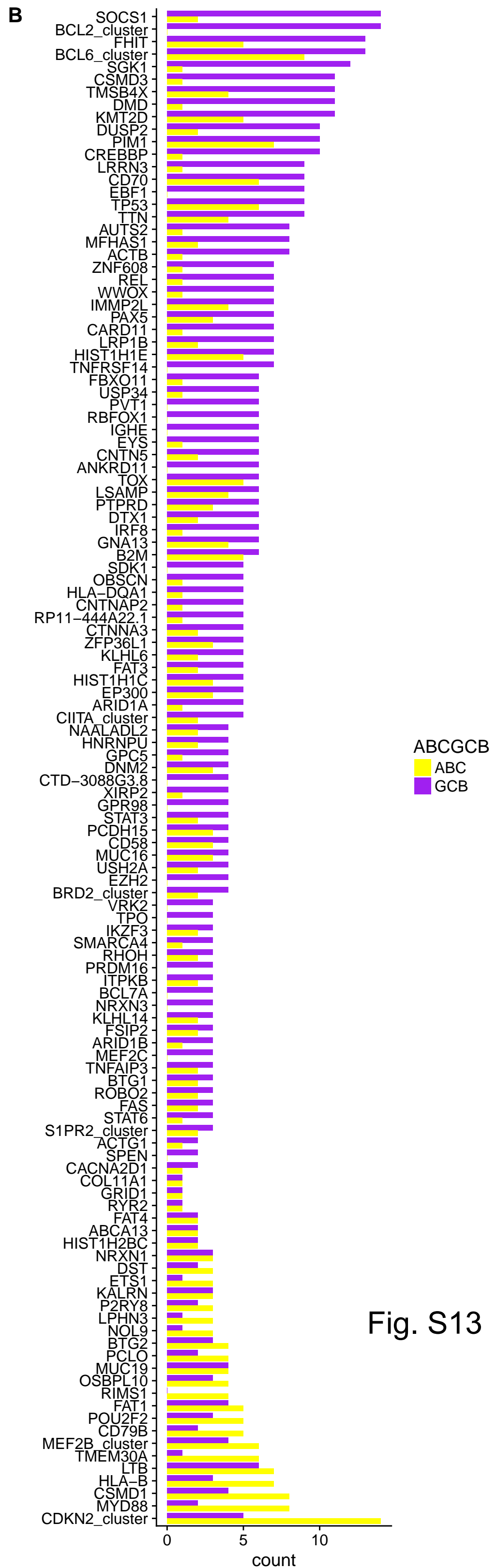
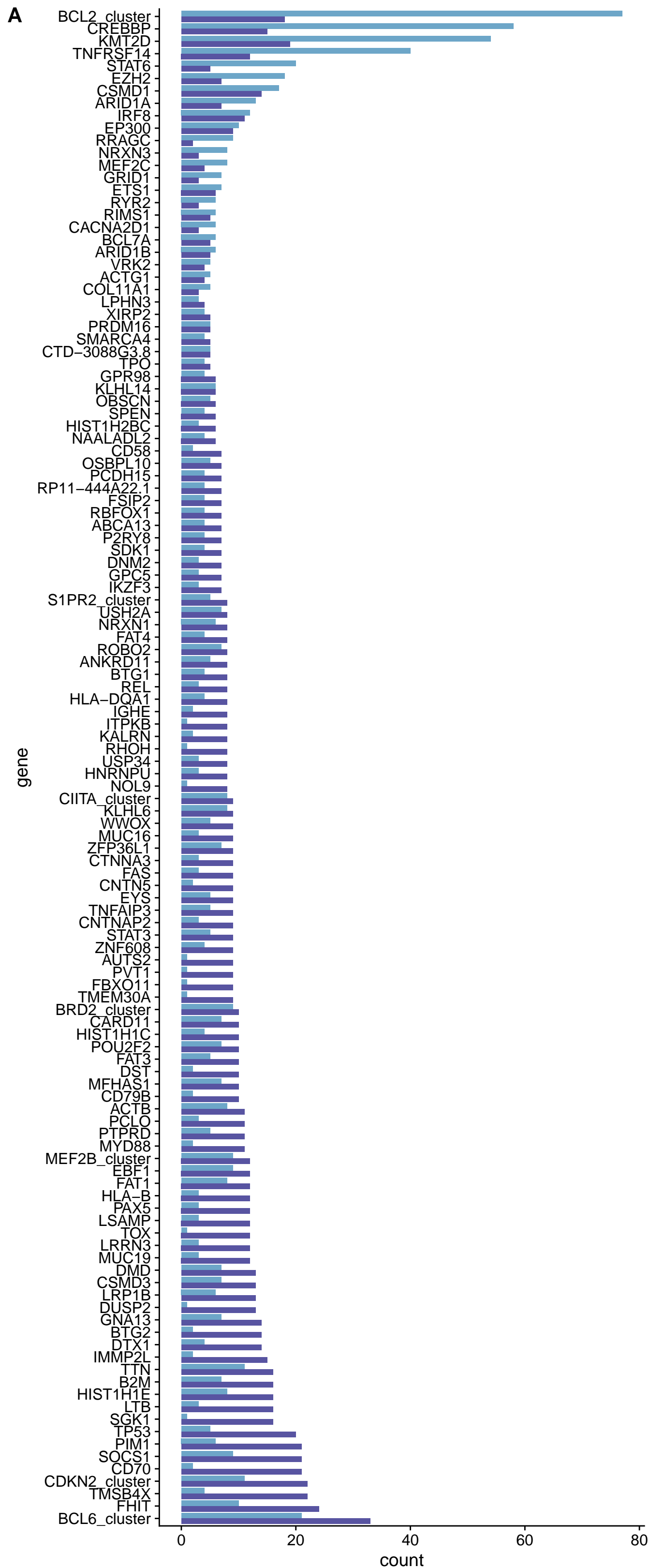


Fig. S13

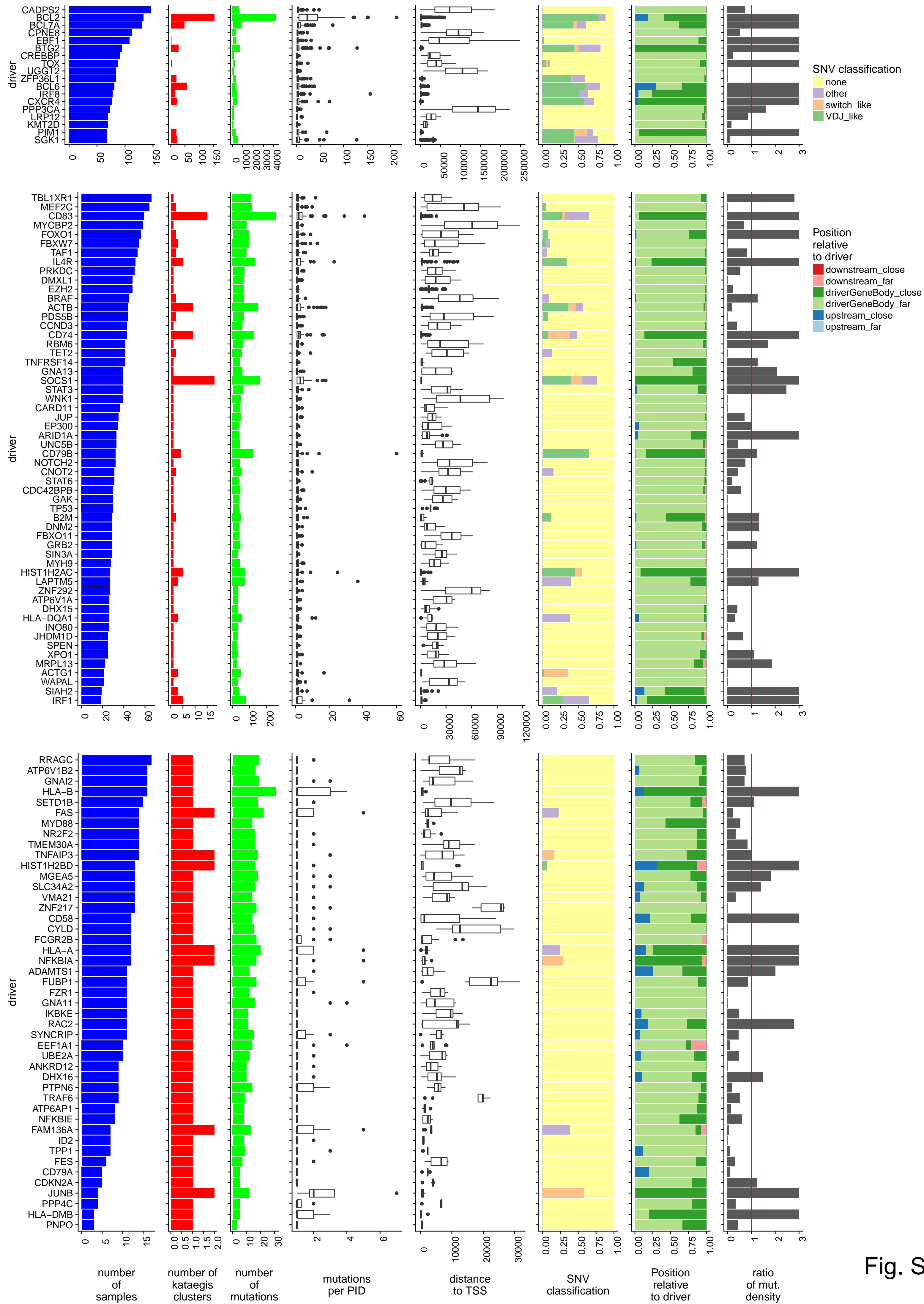


Fig. S14

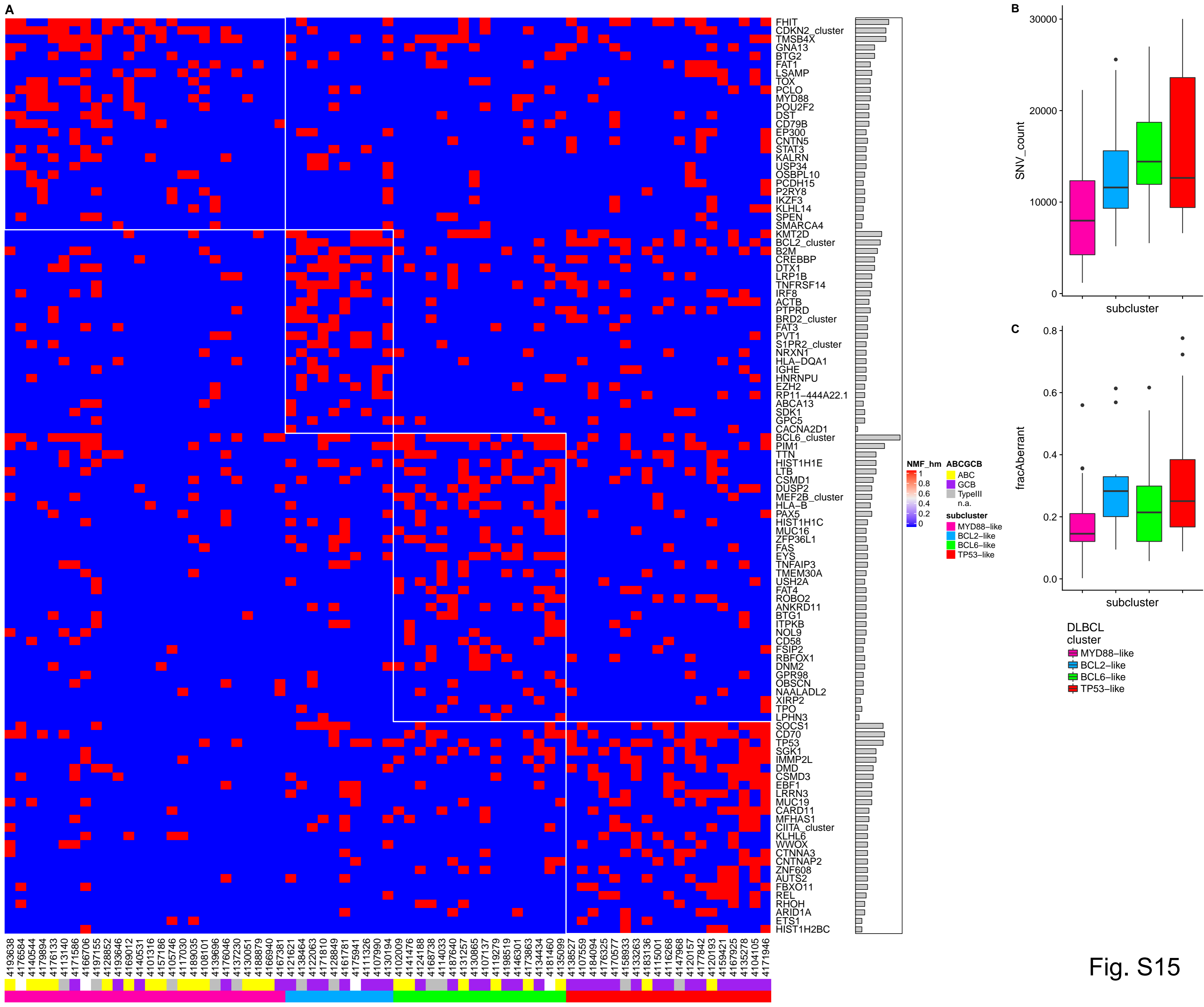


Fig. S15

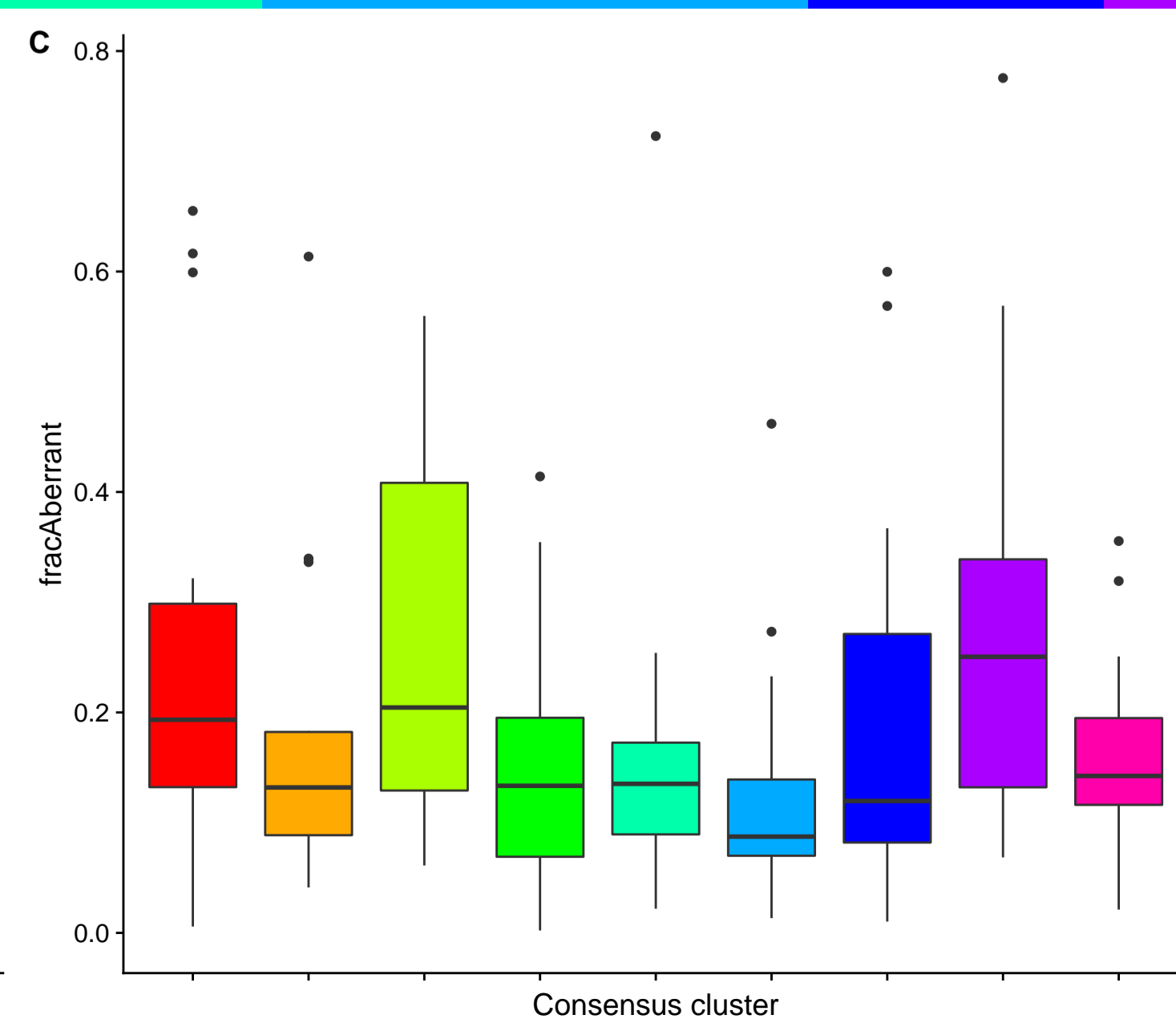
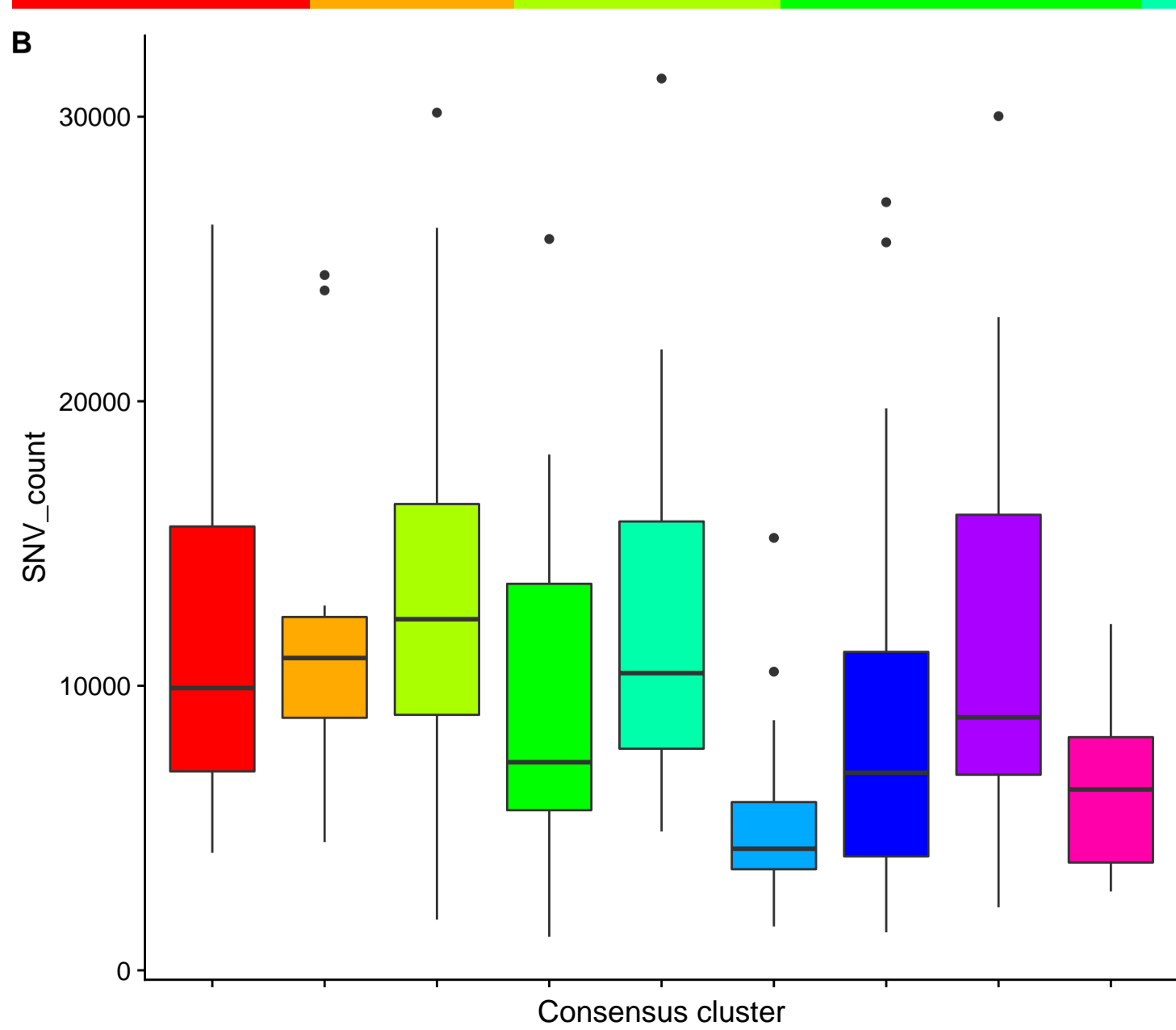
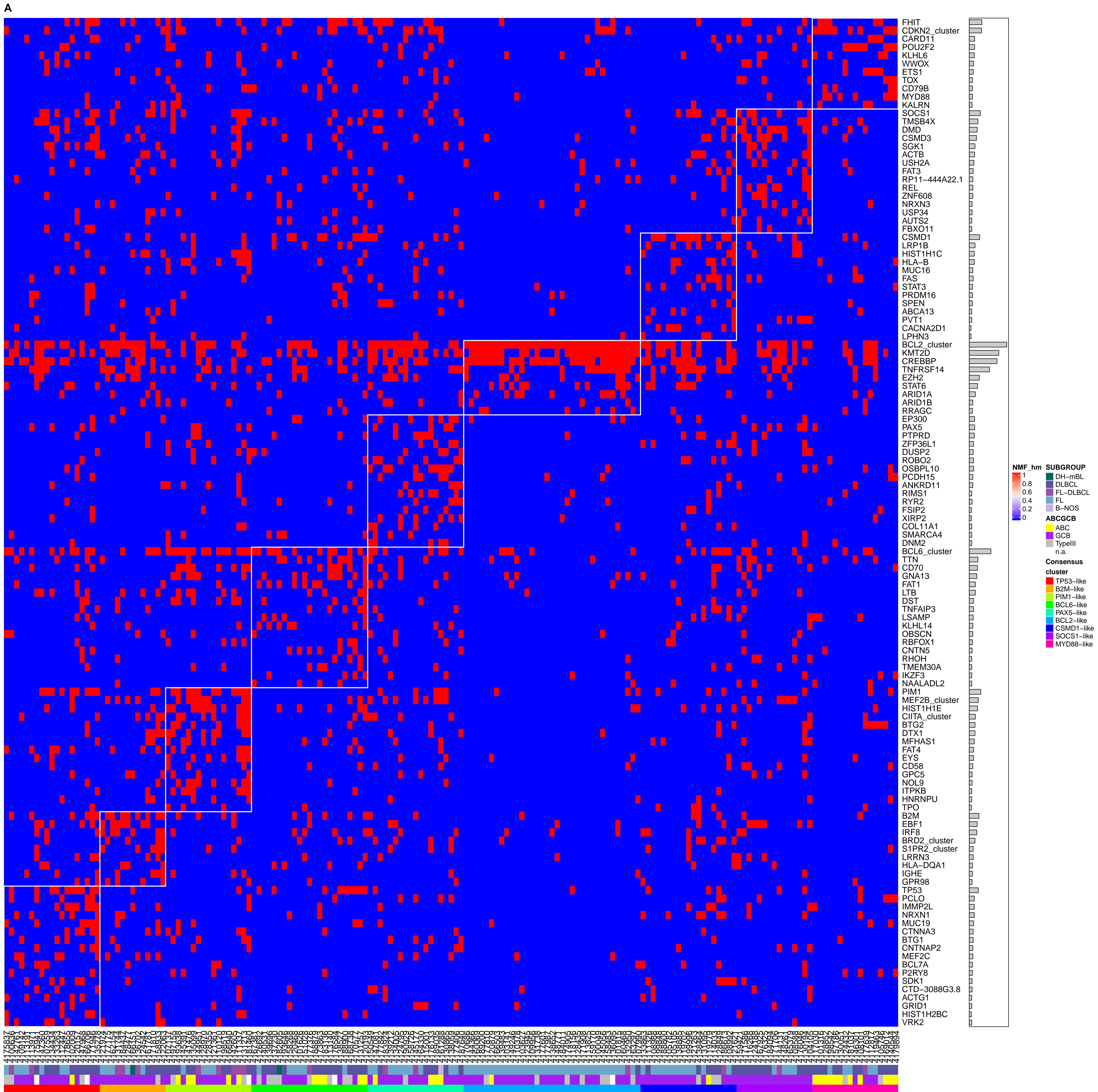


Fig. S16

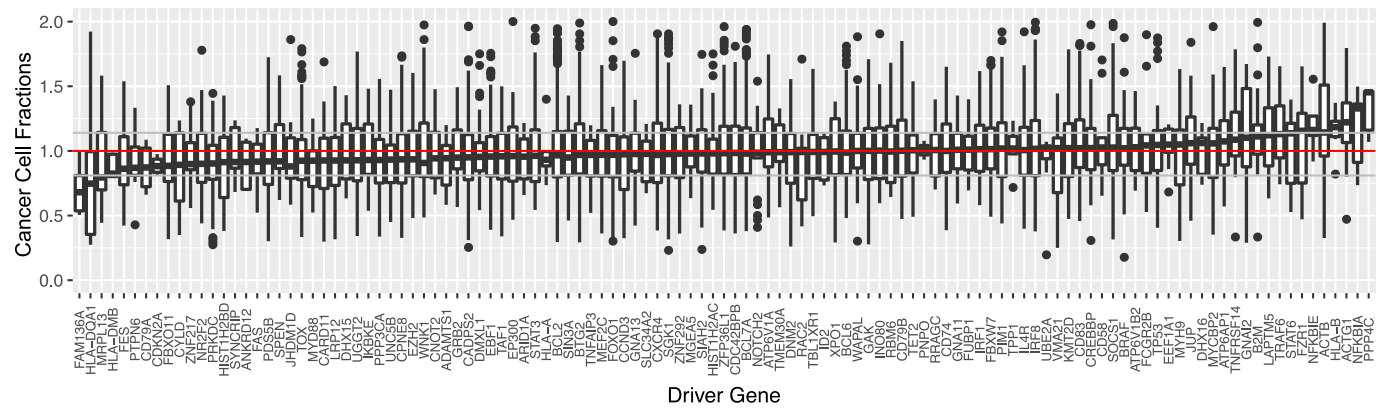


Fig. S17

