**Characterizing selection signatures in coding and noncoding regions of 14,886 cancer genomes**
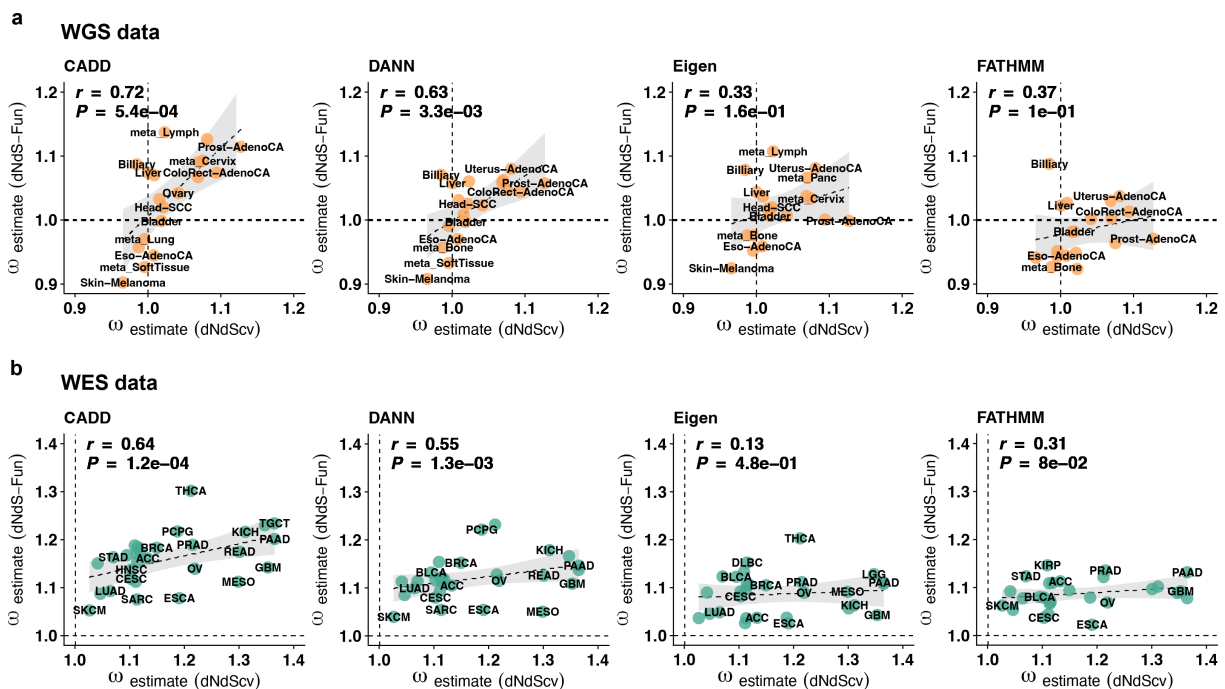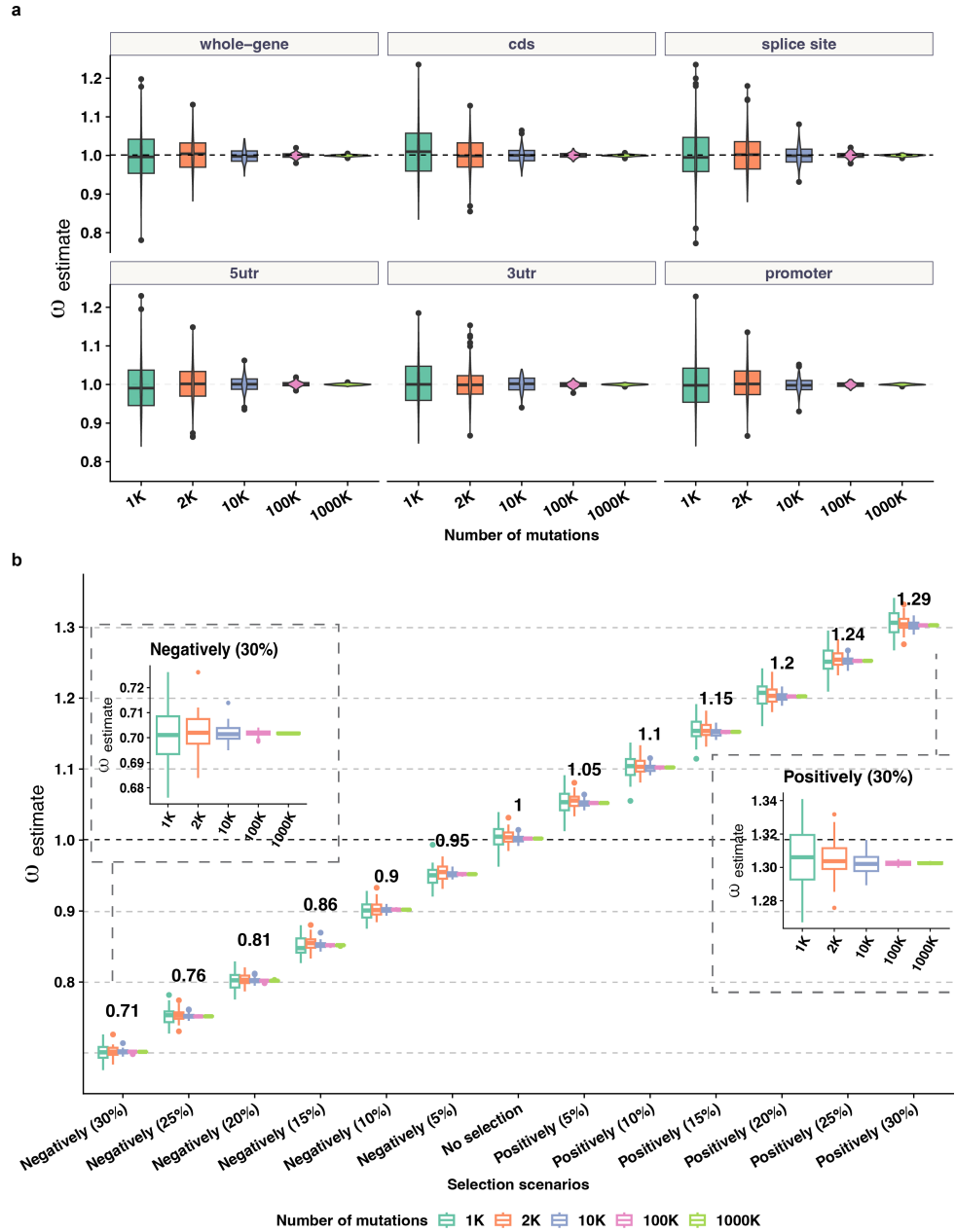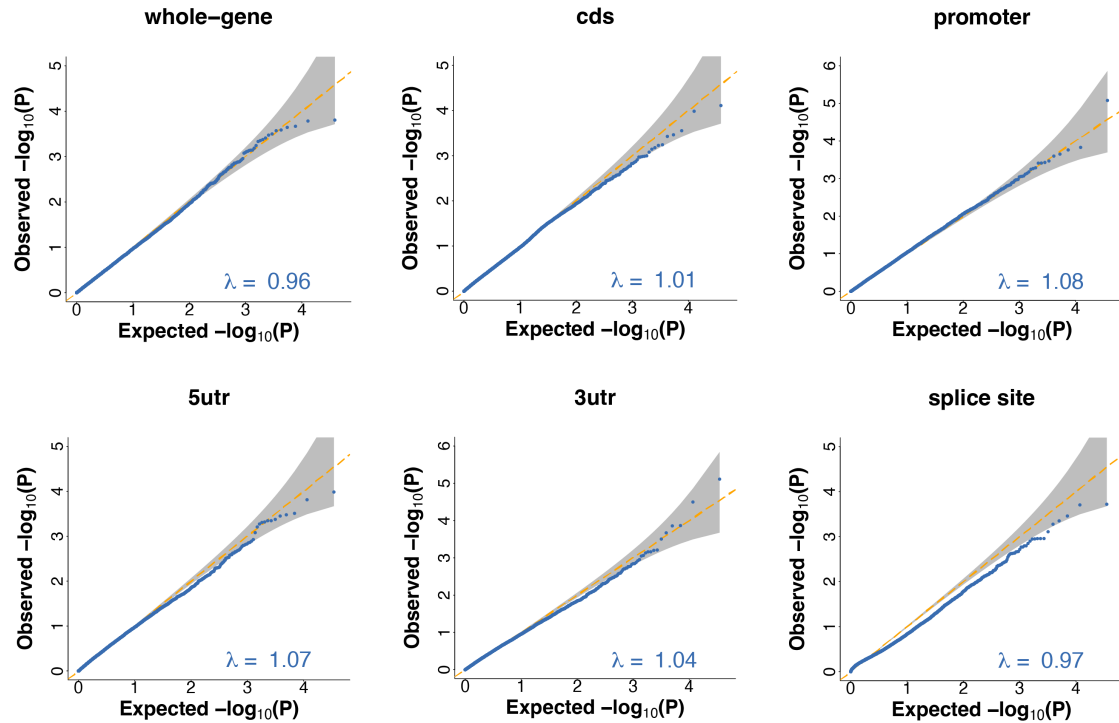
Zheng *et al.*

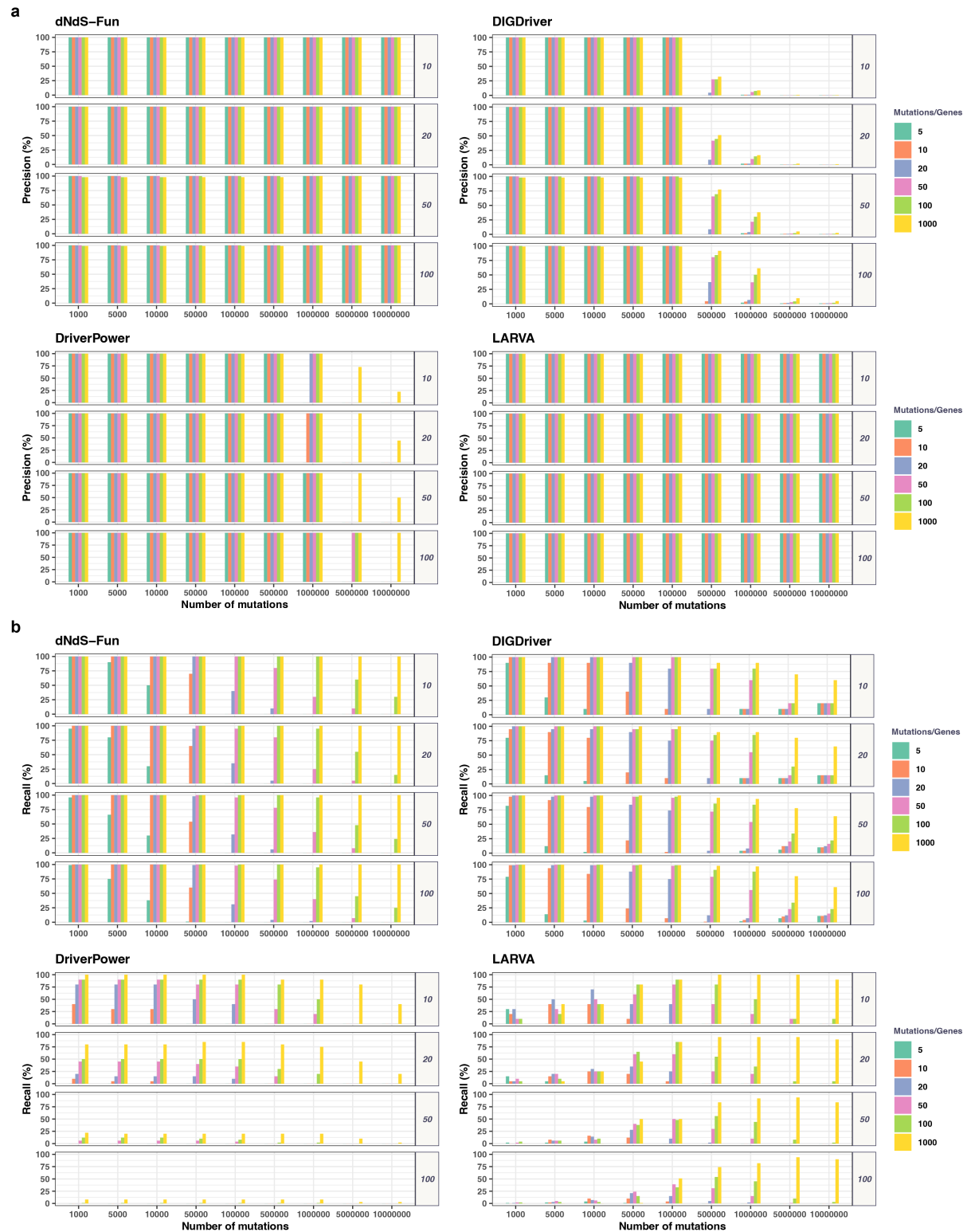**Supplementary Figure 1. Comparison of $\omega$ estimates between dNdS-Fun and dNdScv.** The plot illustrates the correlation between $\omega$ estimates obtained from dNdS-Fun and dNdScv, using only coding mutations from the PCAWG dataset (panel a) and the TCGA dataset (panel b). The dNdS-Fun models are constructed using CADD, DANN, Eigen, and FATHMM scores and are displayed sequentially from left to right. The Pearson correlation coefficient (r) and *P*-value are shown in the top left corner of each plot. The grey shaded area represents the 95% CI.

**Supplementary Figure 2. Estimation of $\omega$ under different dichotomization thresholds.** a) Comparison of $\omega$ estimates between dNdS-Fun with different dichotomization thresholds (i.e., 5:5, 6:4, 7:3, and 8:2) and dNdScv in coding regions (top panel), and between dNdS-Fun with a threshold of 5:5 and dNdS-Fun with other thresholds in noncoding regions (bottom panel). The dichotomization thresholds—5:5, 6:4, 7:3, and 8:2—represent the ratios of less functional to more functional genomic sites. b) Estimates of $\omega$ across different cancer types using dNdS-Fun with different thresholds. c) Standard errors of the $\omega$ estimate across various cancer types using dNdS-Fun with different thresholds.

**Supplementary Figure 3. Simulations under the null and selection models.** a) Estimates of $\omega$ from dNdS-Fun using data simulated under the null model. Each box plot represents the results from 300 simulation replicates, categorized by the number of mutations. b) Estimates of $\omega$ from dNdS-Fun using data simulated under various selection models. Each boxplot shows the results from 50 simulation replicates. The median estimate of $\omega$ for each scenario is displayed above the corresponding box plots, with values in the parentheses indicating the proportion of mutations under selection.
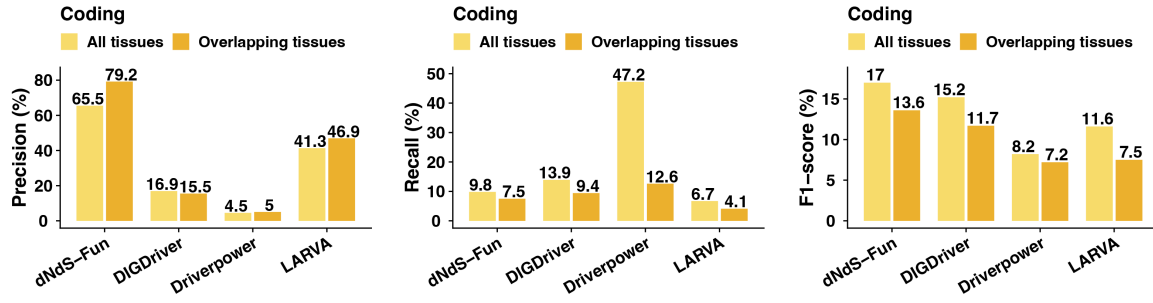
**Supplementary Figure 4. Quantile-Quantile (QQ) plots of the dNdS-Fun test-statistics across genes or genomic elements in simulations under the null model.** The QQ plots compared the distribution of observed *P*-values for individual genes or genomic elements against those expected under the null hypothesis. Each dot represents an individual gene or genomic element. The grey shaded area denotes the 95% CI. $\lambda$ represents the median of the observed test statistics divided by its expected value under the null hypothesis.
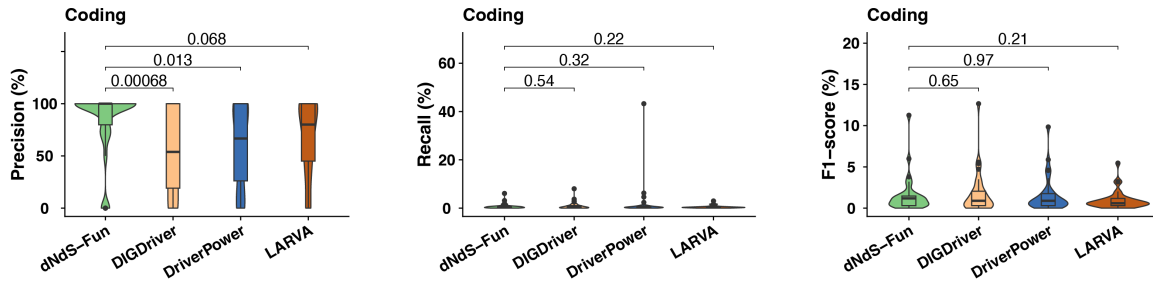
**Supplementary Figure 5. Comparison of the performance of dNdS-Fun with DIGDriver, DriverPower, and LARVA using simulated data.** The simulations varied in the number of driver

genes, driver mutations per gene, and total number of mutations to assess each method's performance. Panel a displays the precision of each method, while panel b presents the recall, with both metrics plotted against the total number of mutations. The simulated driver genes serve as the gold standard for evaluating performance.
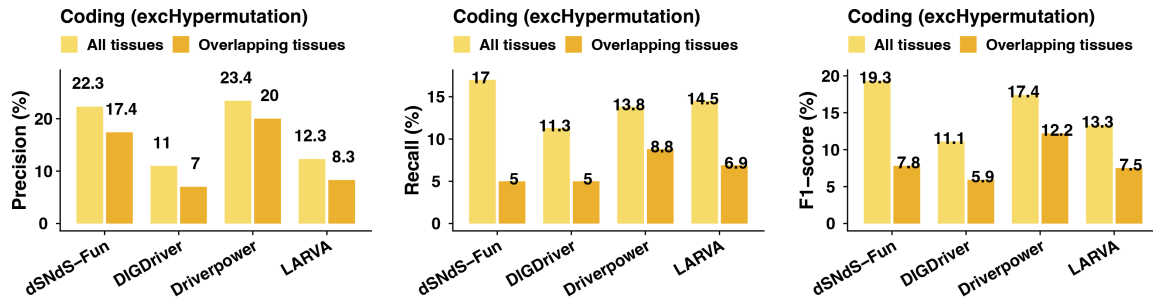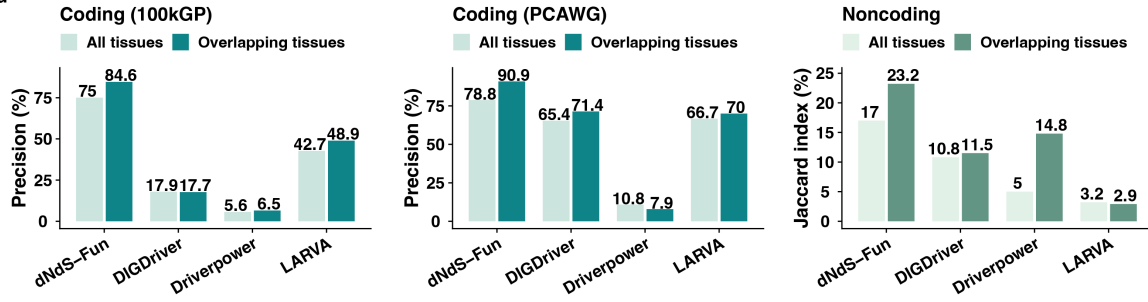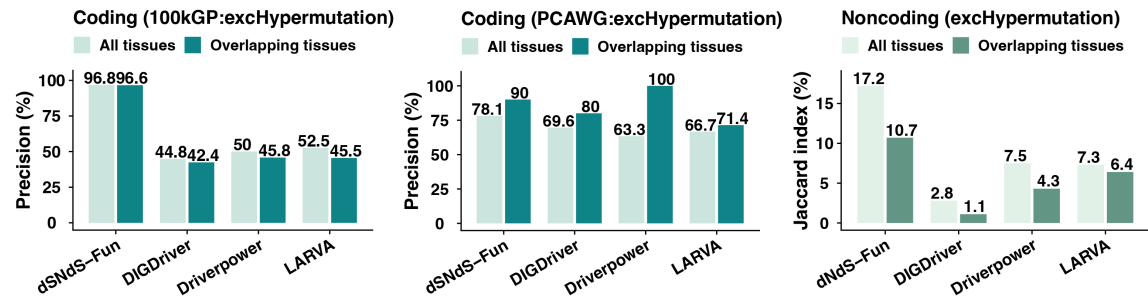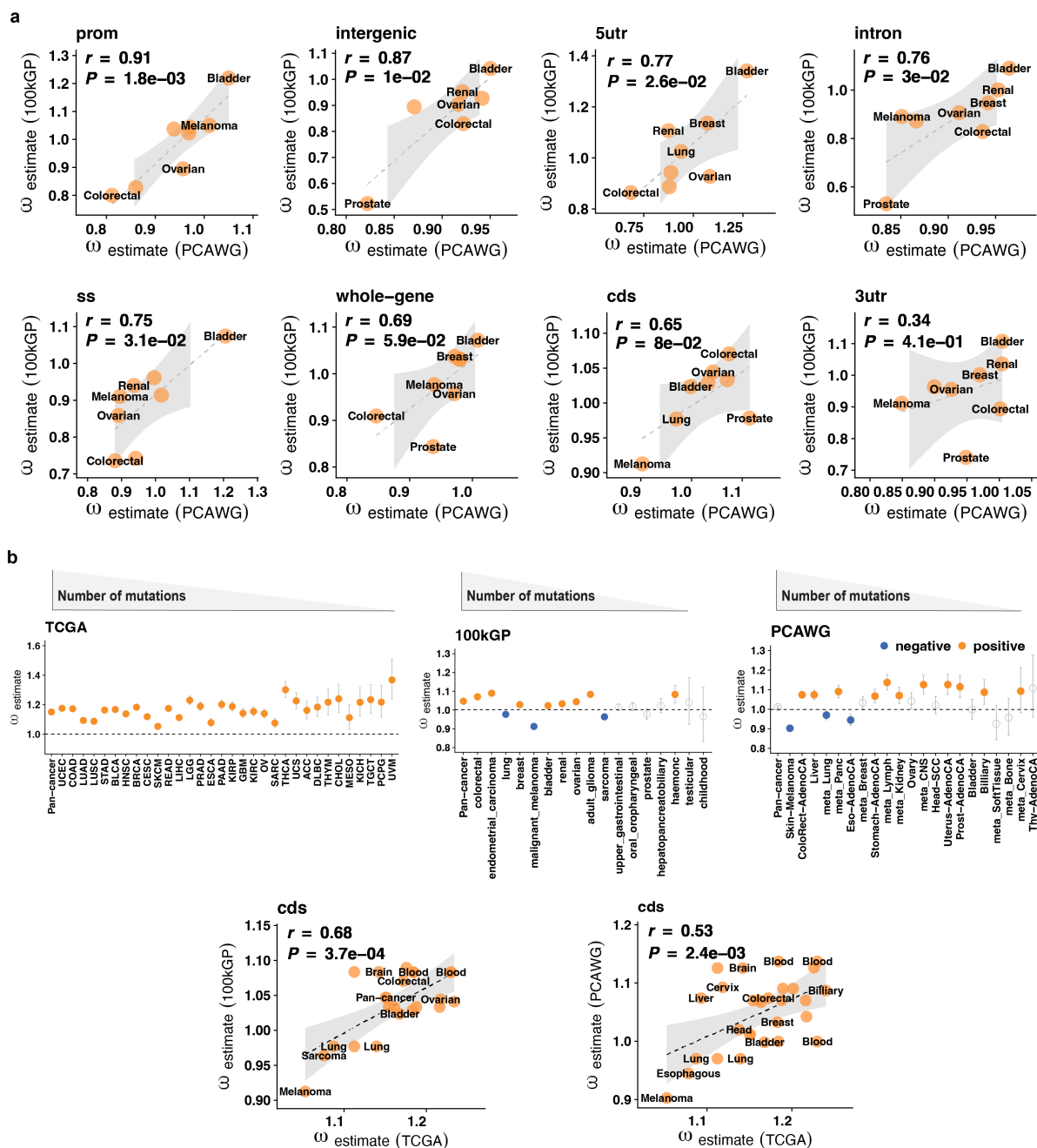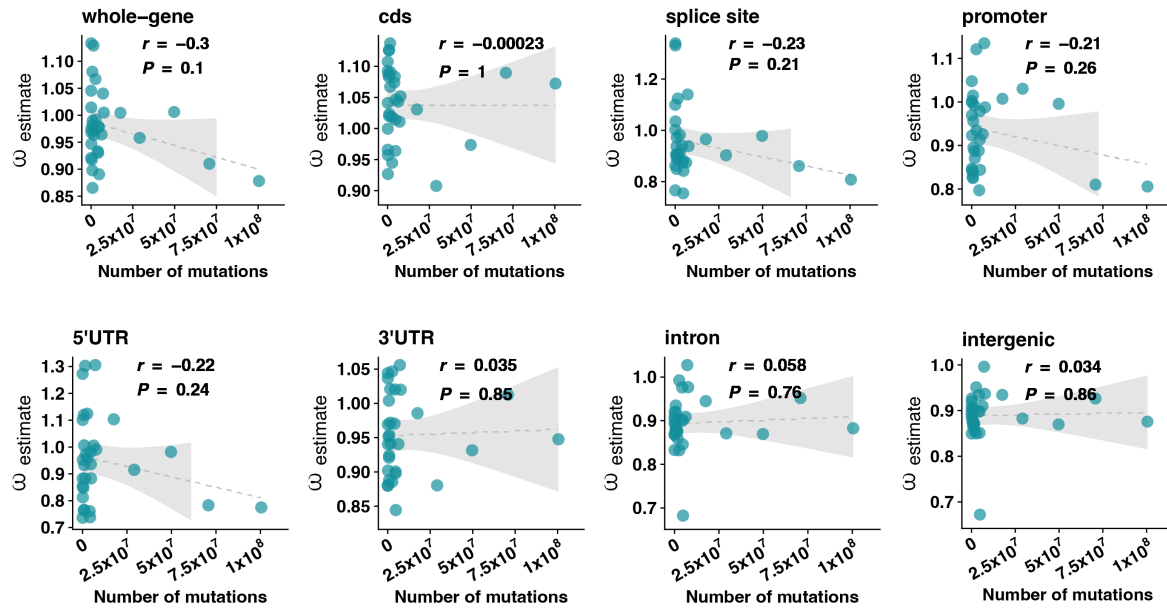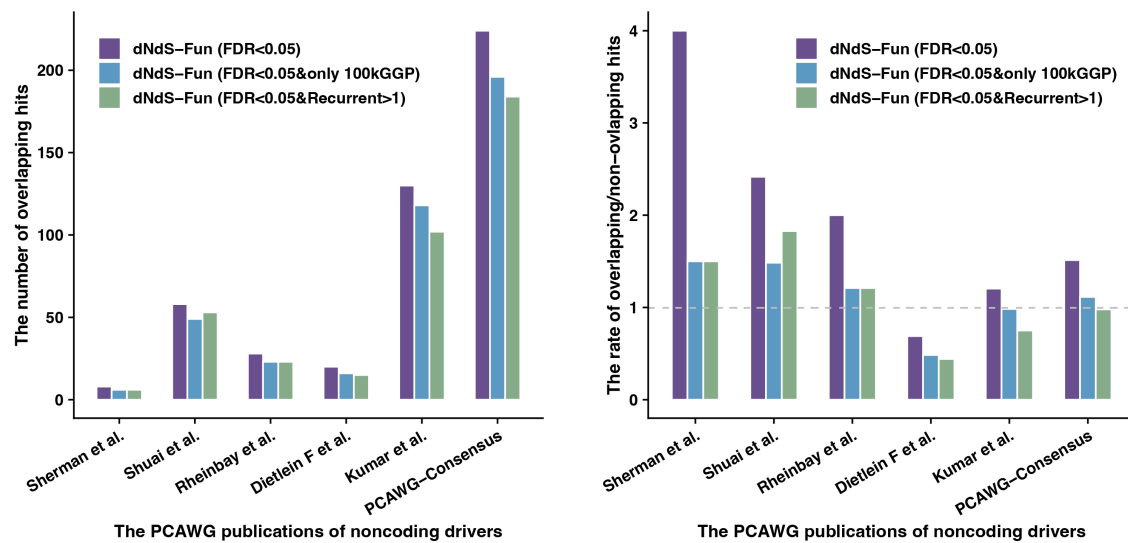
**Supplementary Figure 6. Comparison of the performance of dNdS-Fun with DIGDriver, DriverPower, and LARVA using real data.** a) Precision, recall, and F1-score for genes identified across all tissues and in eight tissues that overlap between 100kGP and PCAWG (i.e., melanoma, breast, bladder, colorectal, prostate, renal, lung, and ovarian cancers), hereafter referred to as overlapping tissues. b) Precision, recall, and F1-score for genes identified in each individual tissue, with *P*-values computed using Student's *t*-test to compare dNdS-Fun with the other methods. c) Precision, recall, and F1-score for genes identified across all tissues, and in overlapping tissues, excluding hypermutator cancers (i.e., colorectal, melanoma, and endometrial cancers). d) Precision of genes identified in all tissues and in overlapping tissues from 100kGP and PCAWG, before and after excluding hypermutator cancers. e) Jaccard similarity of genes identified between 100kGP and PCAWG, before and after excluding hypermutator cancers.

**Supplementary Figure 7. Comparison of ω estimates between different cohorts and sequencing platforms.** a) Comparisons of ω estimates between 100kGP and PCAWG across different functional categories. b) The top panel displays selection signatures in coding regions across the TCGA, 100kGP, and PCAWG datasets. The bottom panel compares the ω estimates in coding regions between different sequencing platforms: WGS (100kGP and PCAWG) vs. WES (TCGA). The Pearson correlation coefficient (r) and *P*-value are shown in the top left corner of each plot. The grey shaded area represents the 95% CI.
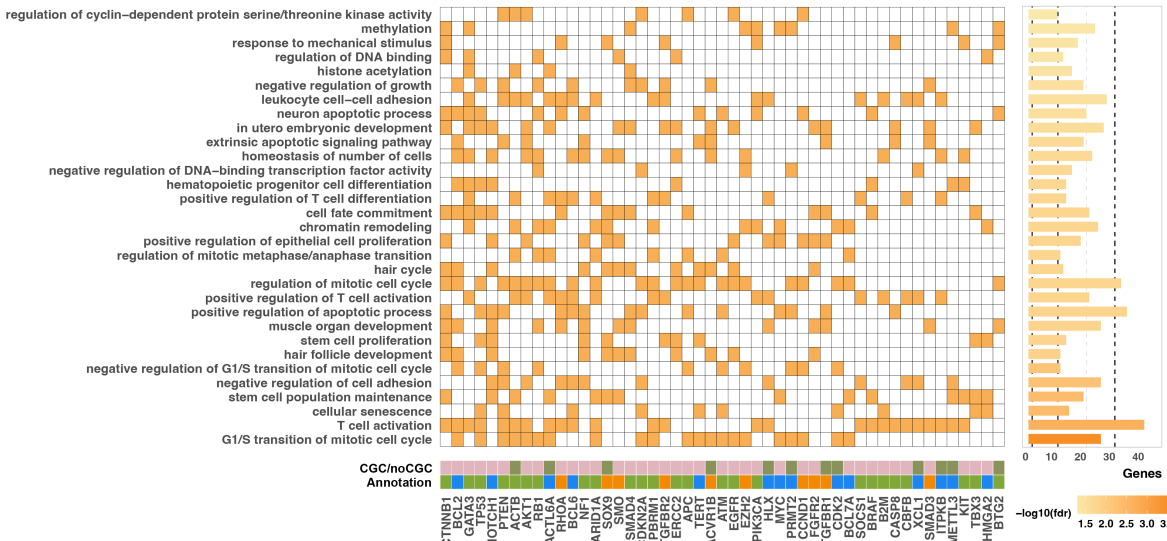
**Supplementary Figure 8. Correlation between number of mutations and $\omega$ estimate across cancer types in each of the eight functional categories**. Each dot represents the $\omega$ estimate from a cancer type. The Pearson correlation coefficient (r) and *P*-value are shown in the top right corner of each plot. The grey shaded area represents the 95% CI.

**Supplementary Figure 9. Comparison of genes under selection identified in the 100kGP and PCAWG datasets.** a) The number of overlapping positively and negatively selected genes identified in both datasets. b) Replication in the 100kGP dataset for genes under selection identified in PCAWG, categorized by their inclusion in the CGC. c) Evidence of selection for genes identified in both the 100kGP and PCAWG datasets.
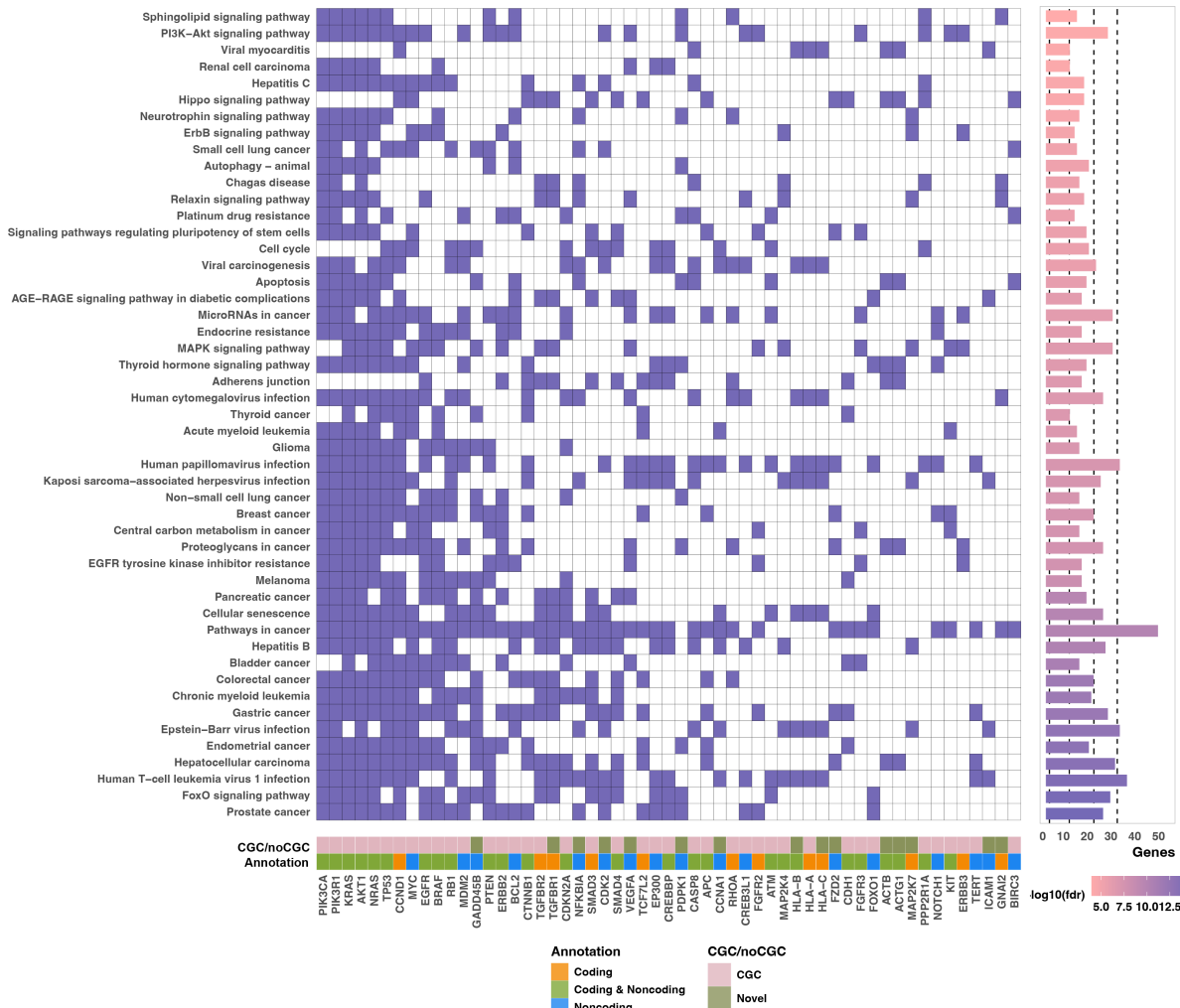
**Supplementary Figure 10. Replication of known noncoding drivers in positively selected genes as identified by dNdS-Fun.** The known noncoding drivers were previously reported in papers from PCAWG and are detailed in **Supplementary Table 9**. This analysis includes genes with mutations in key regulatory regions: promoters, splice sites, 5'UTRs and 3'UTRs. Each hit represents a unique genomic element-tumor pair where replication of the known noncoding driver was observed.
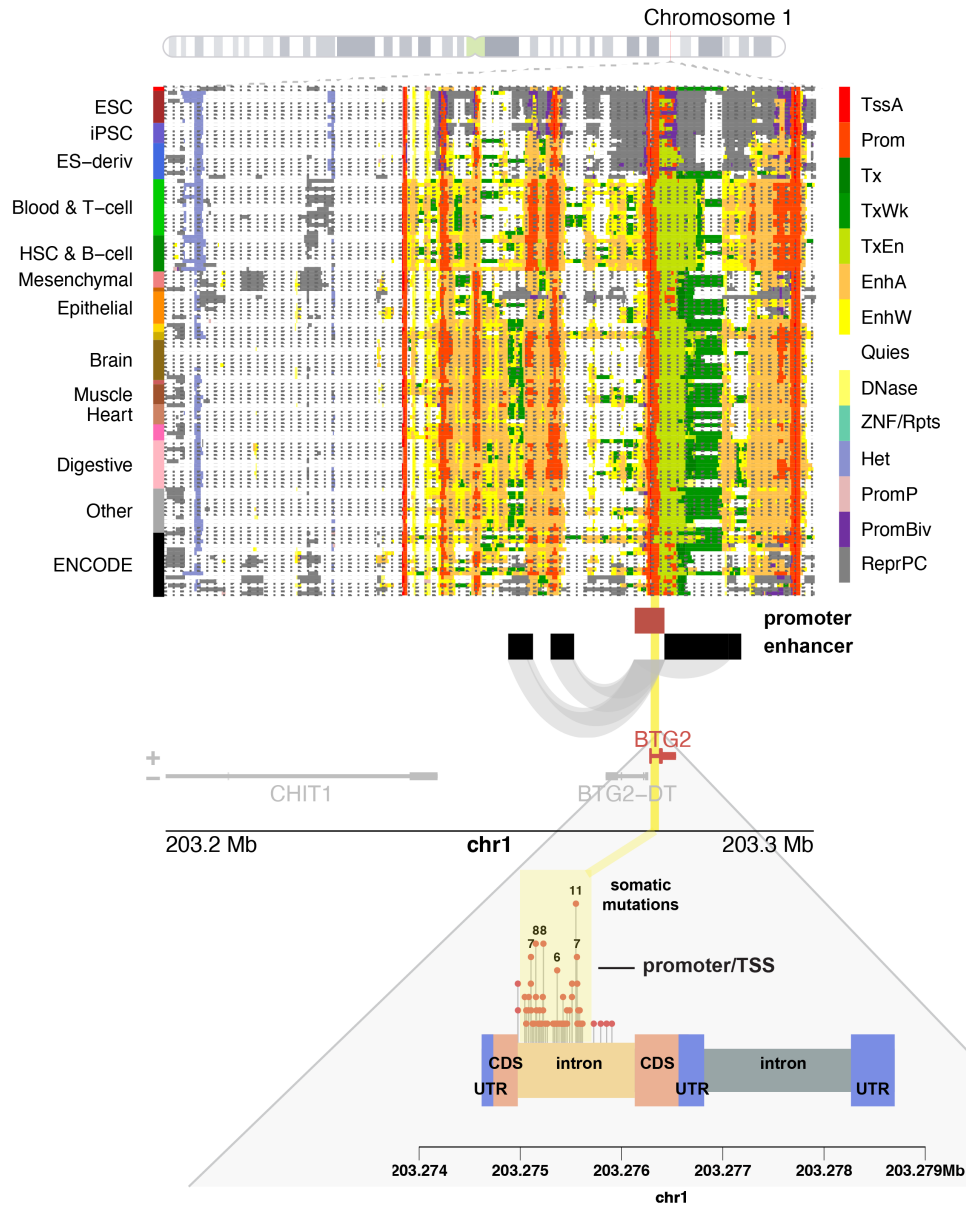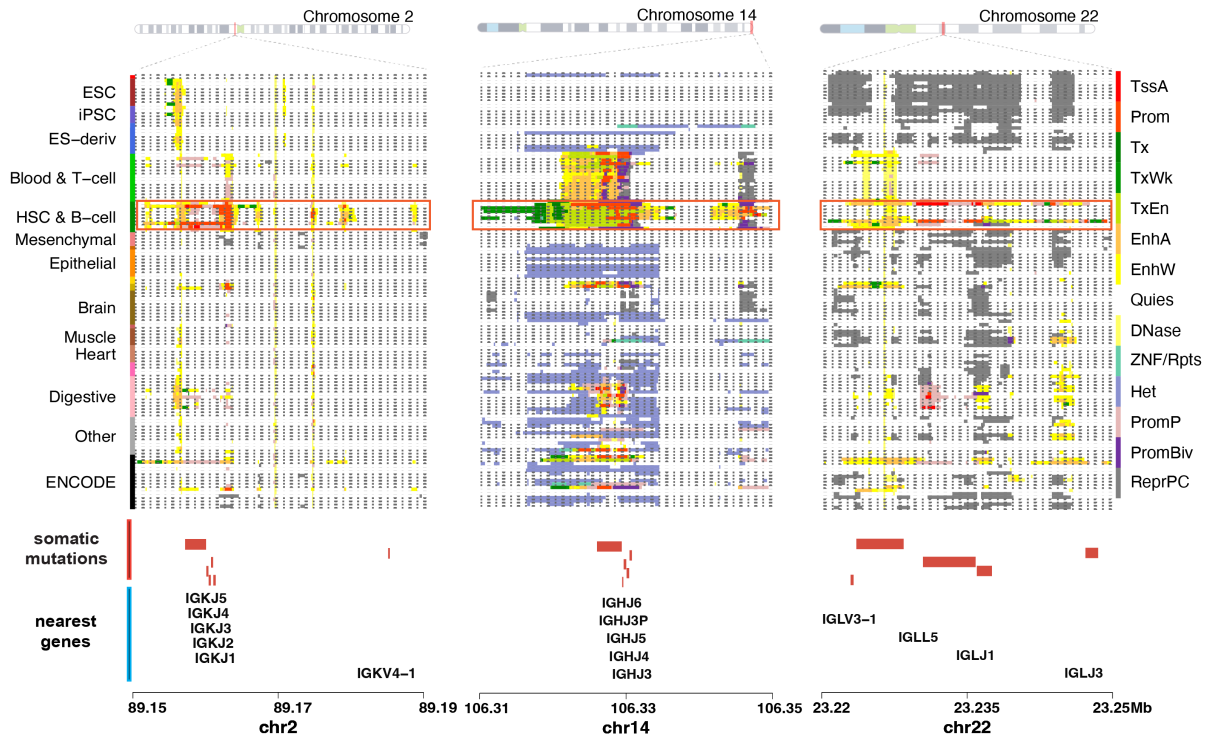
## Gene Ontology (GO) terms

## Kyoto Encyclopedia of Genes and Genomes (KEGG) pathways

**Supplementary Figure 11. Enrichment analysis for positively selected genes identified by dNdS-Fun.** The plot shows the significant GO terms and KEGG pathways. The bar plot on the right shows the number of identified genes associated with each term or pathway. Classification under 'Known/Novel' distinguishes between genes that are listed in the CGC and those that are not. The 'Annotation' classification details the type of mutation evidence observed: coding only, noncoding only, or both coding and noncoding.

**Supplementary Figure 12. Selection signatures in the *BTG2* gene.** The top track shows the 14 chromatin state annotations across 127 samples from the Roadmap Epigenomics Mapping Consortium (REMC), representing various primary cells and tissue types, each row color-coded according to chromatin state. The middle track shows the Hi-C interaction loops, collected from two blood cell samples, indicating regulatory interactions. The bottom part of the figure illustrates recurrent somatic mutations located in the intronic regions of the *BTG2* gene as identified in the PCAWG dataset, with numbers indicating the mutation count at each genomic site. The "promoter/TSS" annotation for this intronic region is inferred based on the chromatin marks shown.

**Supplementary Figure 13. Selection signatures in the Immunoglobulin (Ig) gene loci.** The top track displays the 14 chromatin state annotations from 127 samples provided by the REMC, representing different primary cells and tissue types. Each chromatin state is indicated by a unique color. Notably, items related to HSC and B-cell are highlighted in red. The bottom track shows the regions with evidence of positive selection identified in lymphoma cancers. These selected elements are annotated with the nearest genes located within the Ig gene loci, indicated at their corresponding positions on the genome.