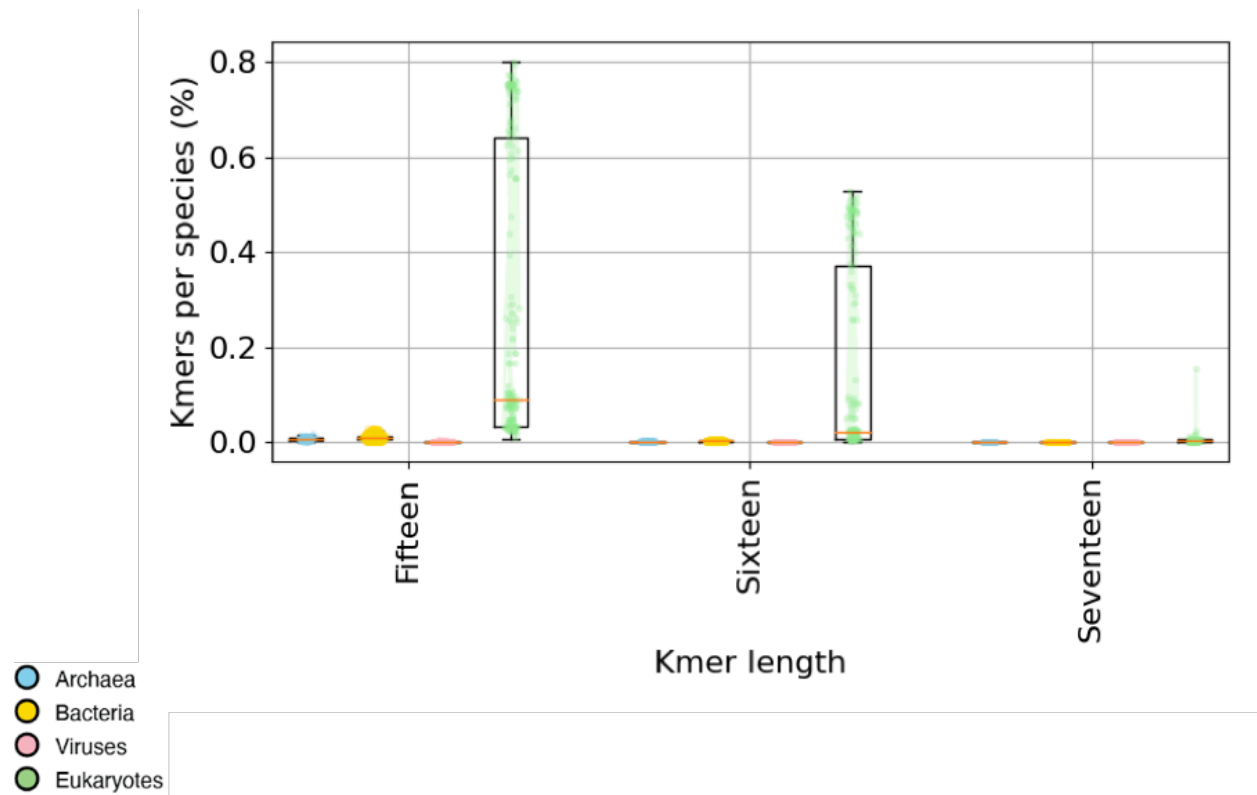
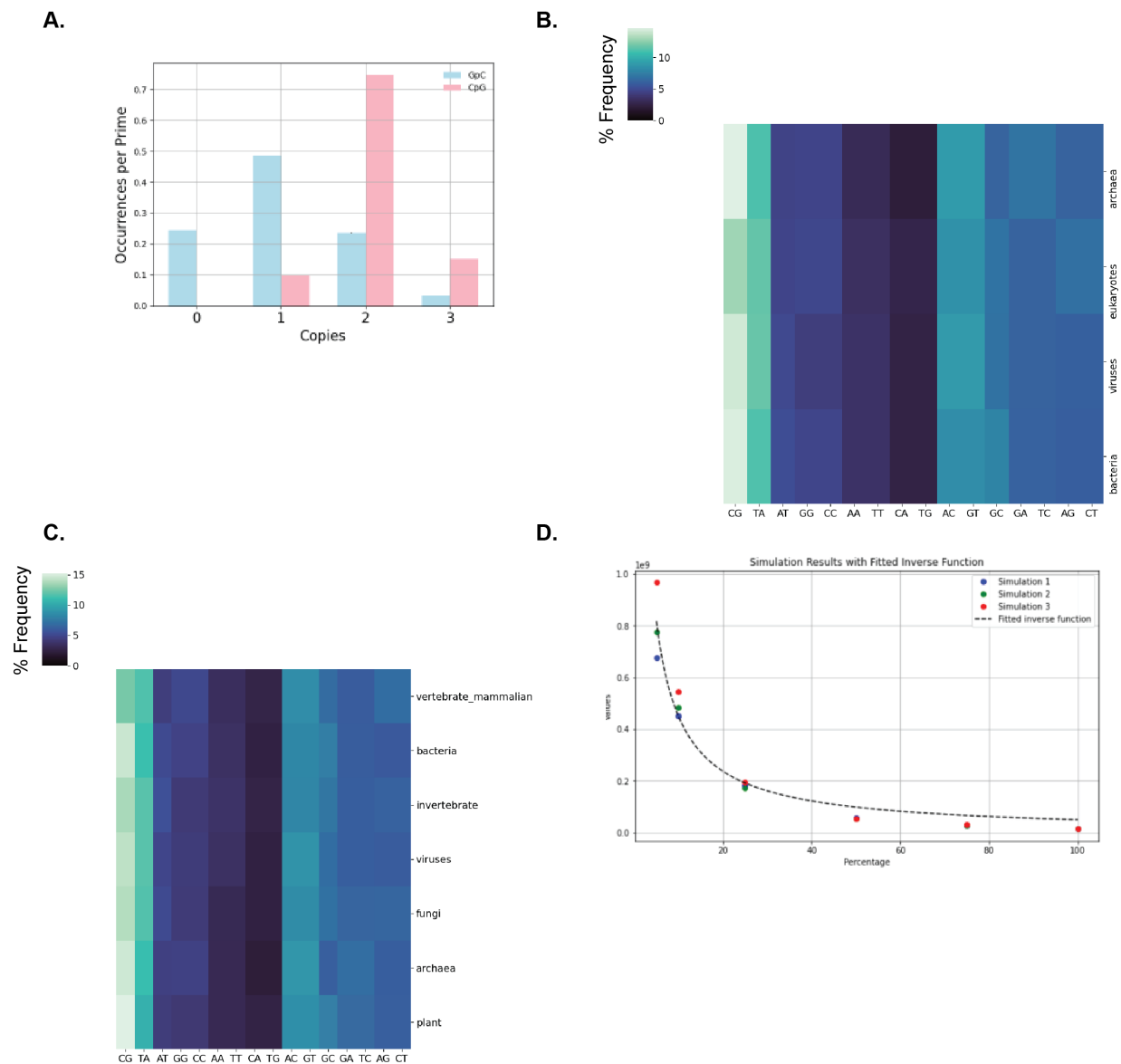


## Supplemental Figures Table of Contents

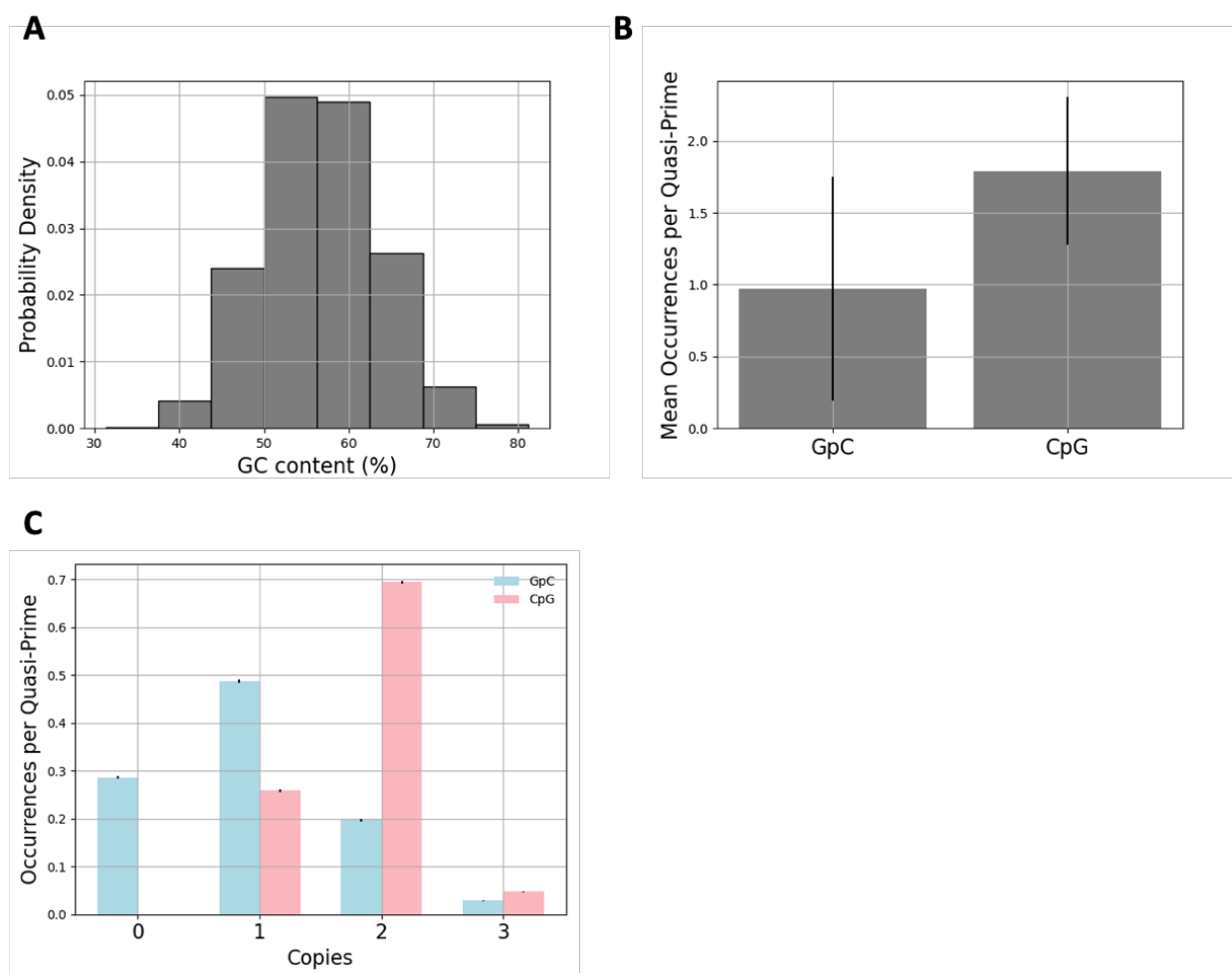
Supplementary Figure 1: Proportion of k-mers found in each reference genome for kmer lengths of 15bp, 16bp, and 17bp across the taxonomic subdivisions. ....	2
Supplementary Figure 2: Features of nucleic quasi-primes.....	3
Supplementary Figure 3: GC content distribution of human quasi-prime sequences.....	4
Supplementary Figure 4: Quasi-prime sequences in genomic sub-compartments. ....	5
Supplementary Figure 5:.....	7
Supplementary Figure 6: Disease association for quasi-prime genes with DisGeNET .....	9
Supplementary Figure 7: Disease association for quasi-prime genes.....	10
Supplementary Figure 8: Single Cell Analysis Metrics .....	13
Supplementary Figure 9: Single Cell Transformation and Variance Analysis.....	14
Supplementary Figure 10: Differential Expression based on quasi-prime genes among cell types. ....	15
Supplementary Figure 11: Gene set enrichment analysis of genes associated with quasi-prime sequences among cell types found in the human single-cell primary motor cortex brain atlas ..	16
Supplementary Figure 12: Gene set enrichment analysis of differentially expressed genes associated with quasi-prime sequences display several disease associations among cell types found in the human single-cell primary motor cortex brain atlas.....	17
Supplementary Figure 13: Variant Characterization of Human Quasi-primes. ....	18



**Supplementary Figure 1: Proportion of k-mers found in each reference genome for kmer lengths of 15bp, 16bp, and 17bp across the taxonomic subdivisions.** Each dot represents the proportion of k-mers observed in a reference genome. The majority of k-mers are found in a minority of the species studied across taxonomies. Error bars represent standard deviation.

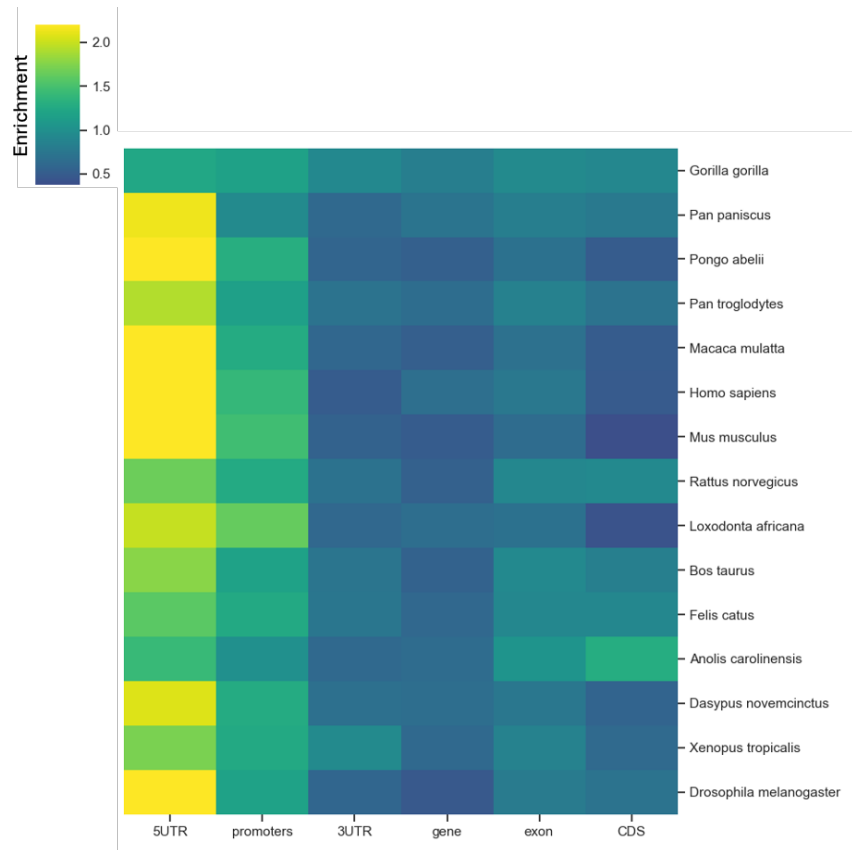


**Supplementary Figure 2: Features of nucleic quasi-primes.** **A.** Number of occurrences of GpCs and CpGs for different copy numbers, per nucleic prime sequence. Error bars are derived from bootstrapping with replacement ( $n=1,000$ ) and represent standard deviation. **B.** Nucleotide composition across three domains of life and virus and **C.** in eukaryotes. **D.** Simulation experiments examining how the number of quasi-primes identified changes as a function of the number of genomes available.

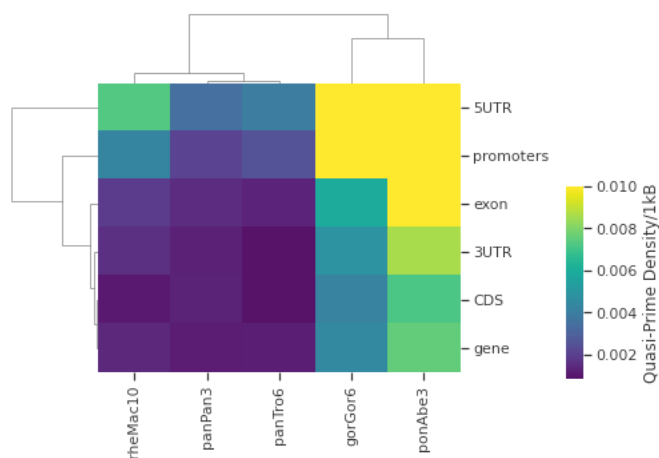


**Supplementary Figure 3: GC content distribution of human quasi-prime sequences. A.** GC content percentage of human quasi-prime sequences. **B.** Average number of GpC and CpG occurrences per human quasi-prime. Error bars show standard deviation. **C.** Number of occurrences of GpCs and CpGs for different copy numbers, per nucleic quasi-prime sequence. Error bars are derived from bootstrapping with replacement ( $n=1,000$ ) and represent standard deviation.

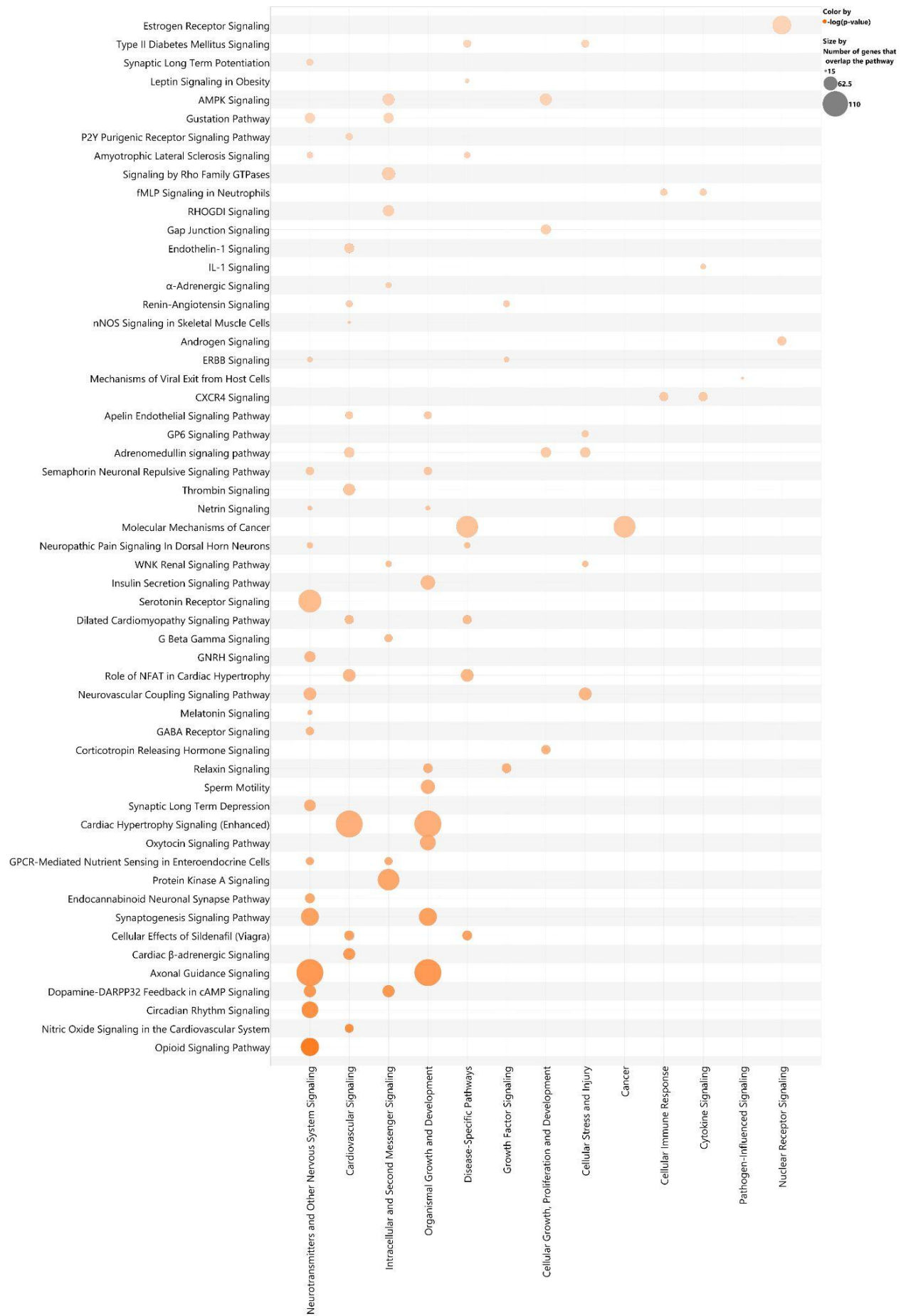
**A.**



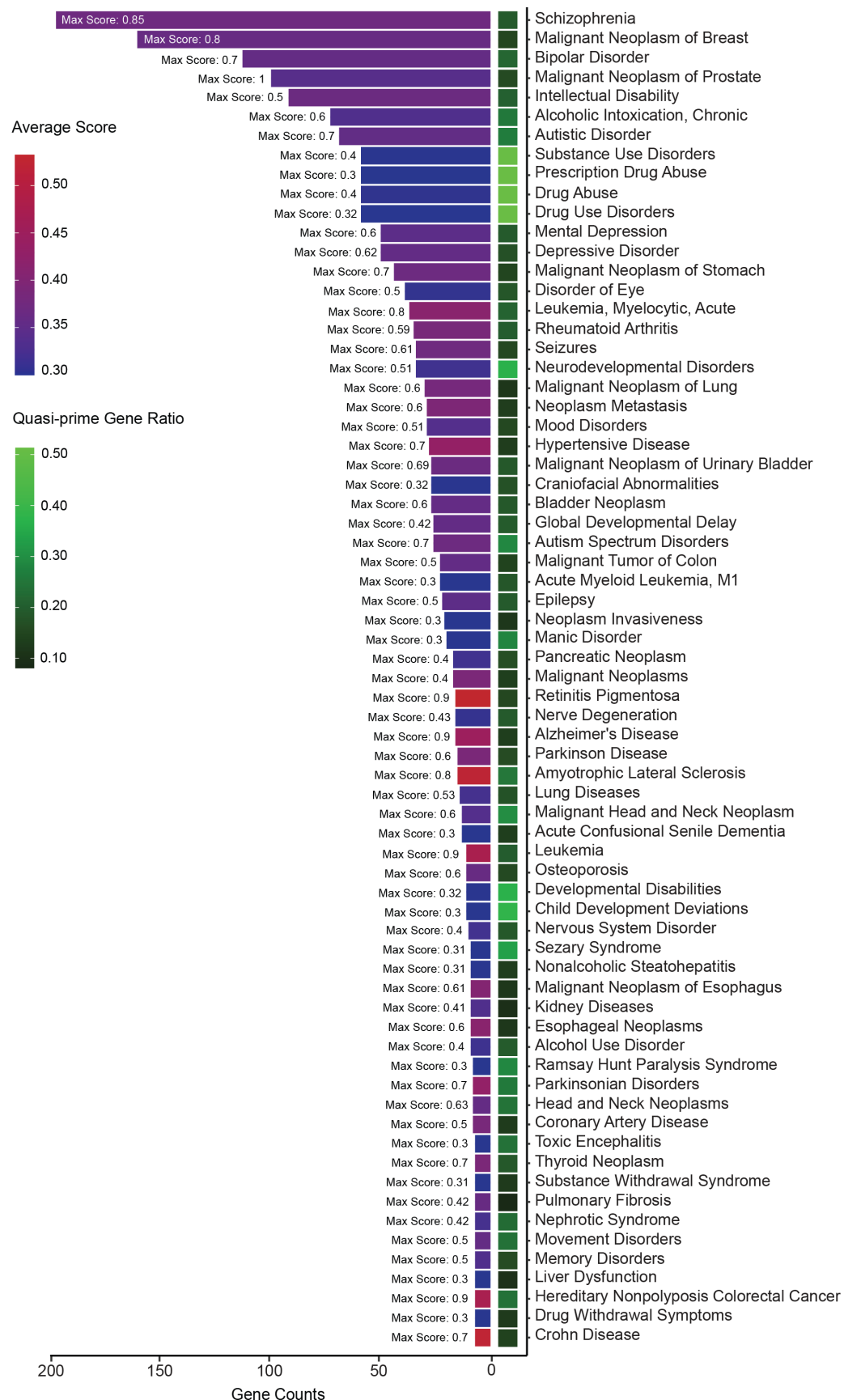
**B.**



**Supplementary Figure 4: Quasi-prime sequences in genomic sub-compartments. A.** Density of quasi-primes across genic regions of multiple primate genomes, rodents, and other mammals and invertebrates and **B.** of five non-human primate genomes. NCBI RefSeq annotation was used for all non-human primates.



**Supplementary Figure 5:** Bubble plot showing the pathway categories (Y axis) of the enriched canonical pathways (X axis) enriched in the quasi-prime gene set. The color represents the adjusted  $p$ -value of the enrichment (Increasing orange hue, the lower the  $p$ -value. The cutoff of  $p$ -value was set at  $<0.001$ ). The size of each bubble represents the number of overlapping genes in the pathway.





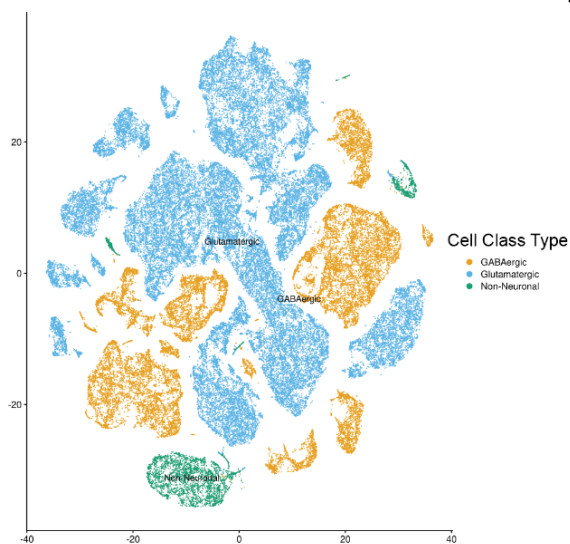
### **Supplementary Figure 6: Disease association for quasi-prime genes with DisGeNET**

Full gene disease bar plot presenting the count of associated genes across different disease conditions for all diseases with more than 6 associated genes. The average disease association score for all quasi-prime genes for a given disease within “ALL” DisGeNET databases is visualized in the color of the bar. The max score for all genes within a disease is displayed in text on figure. The number of quasi-prime genes out of the total genes annotated in the database is represented as a ratio shown in the heatmap tile.

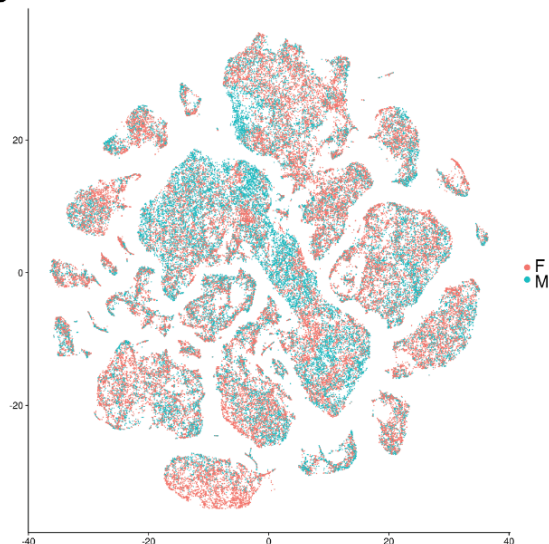


class. **B-C.** gnomAD pLOF constraint gene analysis. **B.** The odds ratio enrichment of highly constrained human quasi-prime genes (with pLI  $\geq 0.9$  and LOEUF  $< 0.35$ ) is displayed as a bar plot. The enriched gene set was tested for statistical significance with a hypergeometric test. Significance levels are indicated as follows: \*  $p < 1E-25$ , \*\*  $p < 1E-35$ , \*\*\*  $p < 1E-45$ , and \*\*\*\*  $p < 1E-55$ . **C.** pLOF variant type enrichment across different LOEUF Decile bins for human quasi-prime genes.

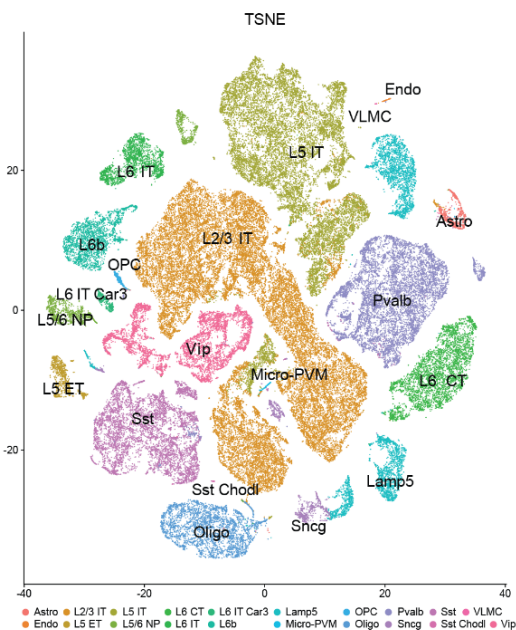
A



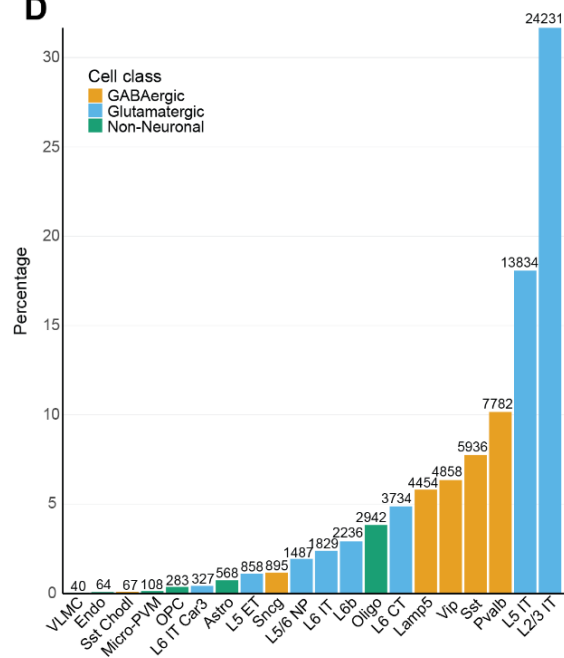
B



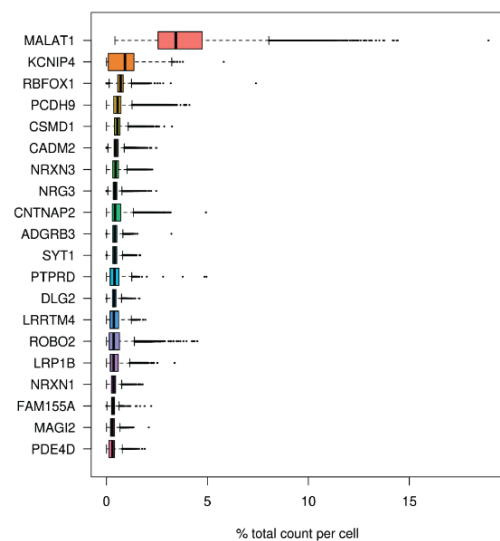
C



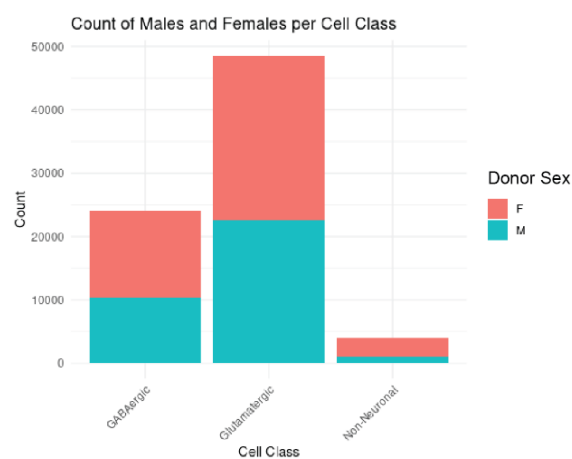
D



E

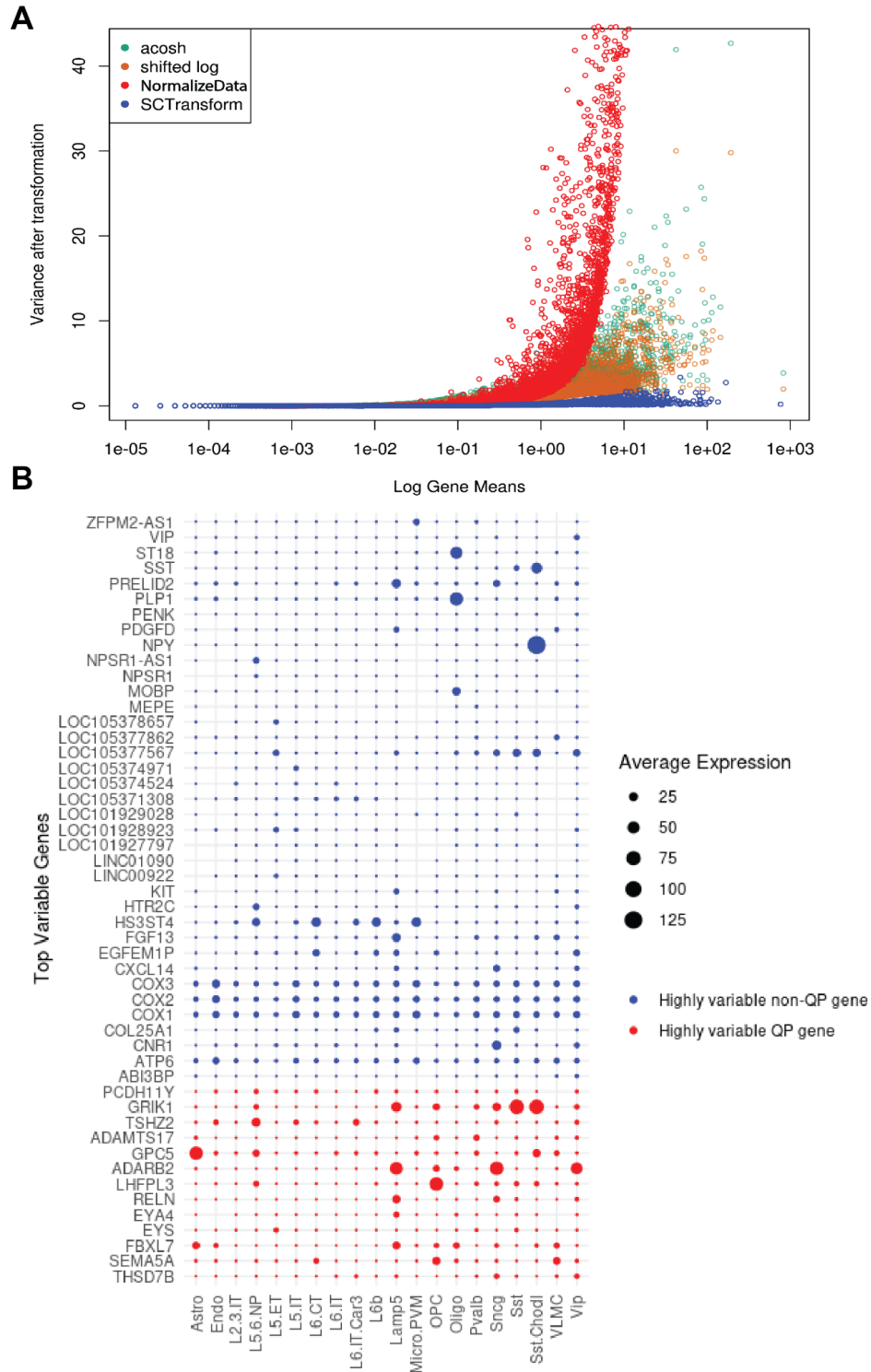


F



### **Supplementary Figure 8: Single Cell Analysis Metrics**

Single cell analysis metrics of M1 Primary Motor Cortex from Single Cell Brain Atlas. **A.** t-SNE plot showing GABAergic, Glutamatergic and non-neuronal cell labels. **B.** t-SNE plot showing cells colored by gender. **C.** t-SNE plot showing cell types found in the human brain atlas (Bakken et al. 2021). **D.** Bar plot shows the percentage of total cells and the count above each bar for each cell type. Bars are colored with cell class. **E.** Box plot sorted by the percentage of genes count expression per cell. **F.** Bar plot representing the differences in counts of cells by donor sex per cell class.

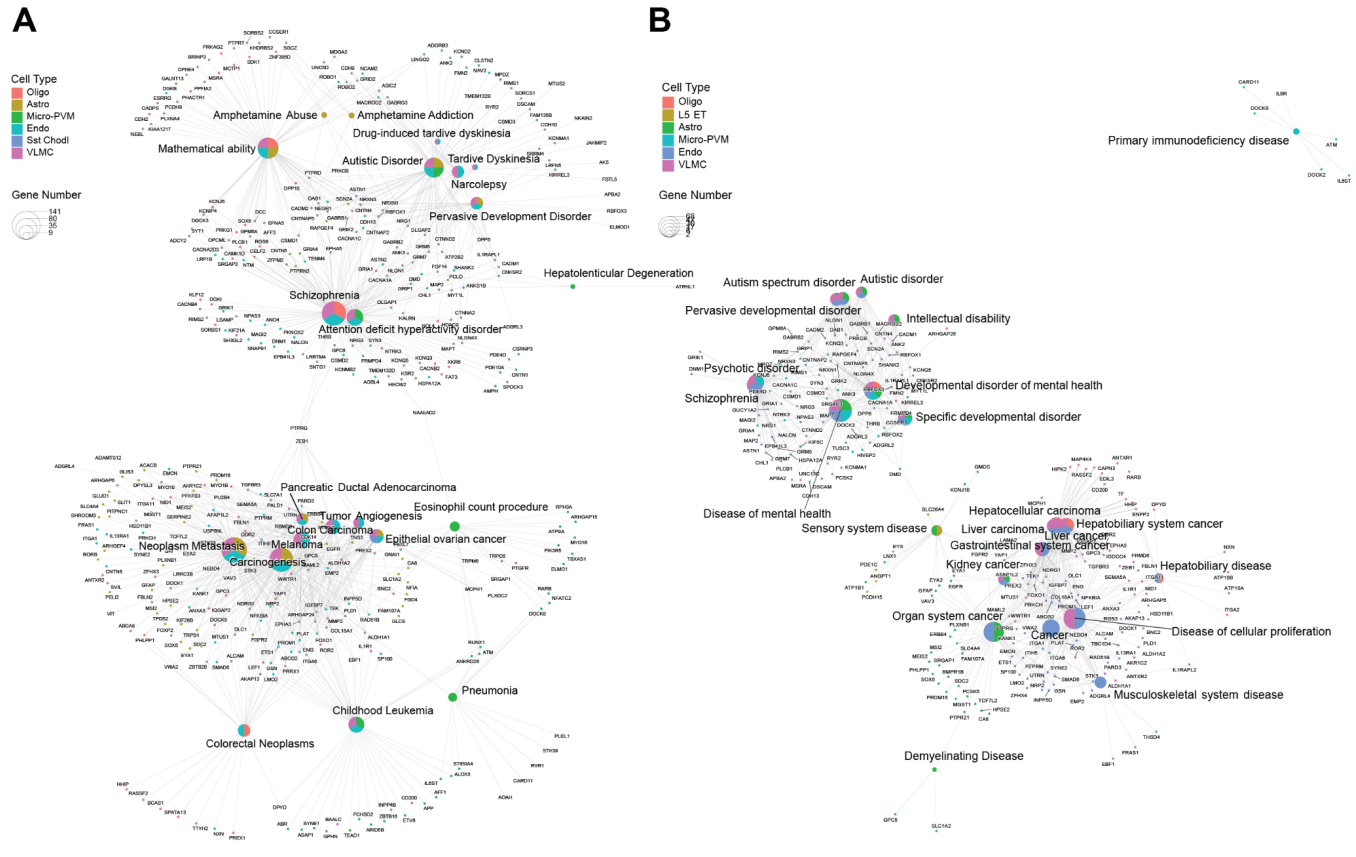


**Supplementary Figure 9: Single Cell Transformation and Variance Analysis. A.** Variance after transformation using four methods, acosh, shifted log, NormalizeData (Seurat), and SCTransform (Seurat). **B.** Dot plot representing the average expression across cell types between quasi-prime and non-quasi prime genes for the top 25 genes in the dataset.









**Supplementary Figure 12: Gene set enrichment analysis of differentially expressed genes associated with quasi-prime sequences display several disease associations among cell types found in the human single-cell primary motor cortex brain atlas. Significant ( $p$  value  $< 0.05$ ) differentially expressed genes with an absolute  $\log_2$  fold change  $> 1$  amongst cell types were used in GSEA analysis as represented by a network graph. **A.** DisGeNET “ALL” database. **B.** Disease ontology database. Different cell types are color coded and pie charts reflect the cell types at which each the diseases shown were associated.**

