

Supplemental Information

The impact of package selection and versioning on single-cell RNA-seq analysis

Joseph M Rich^{1,2}, Lambda Moses¹, Pétur Helgi Einarsson³, Kayla Jackson^{1,2}, Laura Luebbert¹, A. Sina Boeshaghi¹, Sindri Antonsson³, Delaney K. Sullivan^{1,4}, Nicolas Bray⁵, Páll Melsted³, and Lior Pachter^{*1,6,7}

¹*Biology and Biological Engineering, California Institute of Technology,
Pasadena, CA, 91125, USA*

²*USC-Caltech MD/PhD Program, Keck School of Medicine, Los Angeles,
CA, 90033, USA*

³*Faculty of Industrial Engineering, Mechanical Engineering and Computer
Science, Reykjavík, Iceland*

⁴*UCLA-Caltech Medical Scientist Training Program, David Geffen School of
Medicine, University of California, Los Angeles, Los Angeles, CA, 90095,
USA*

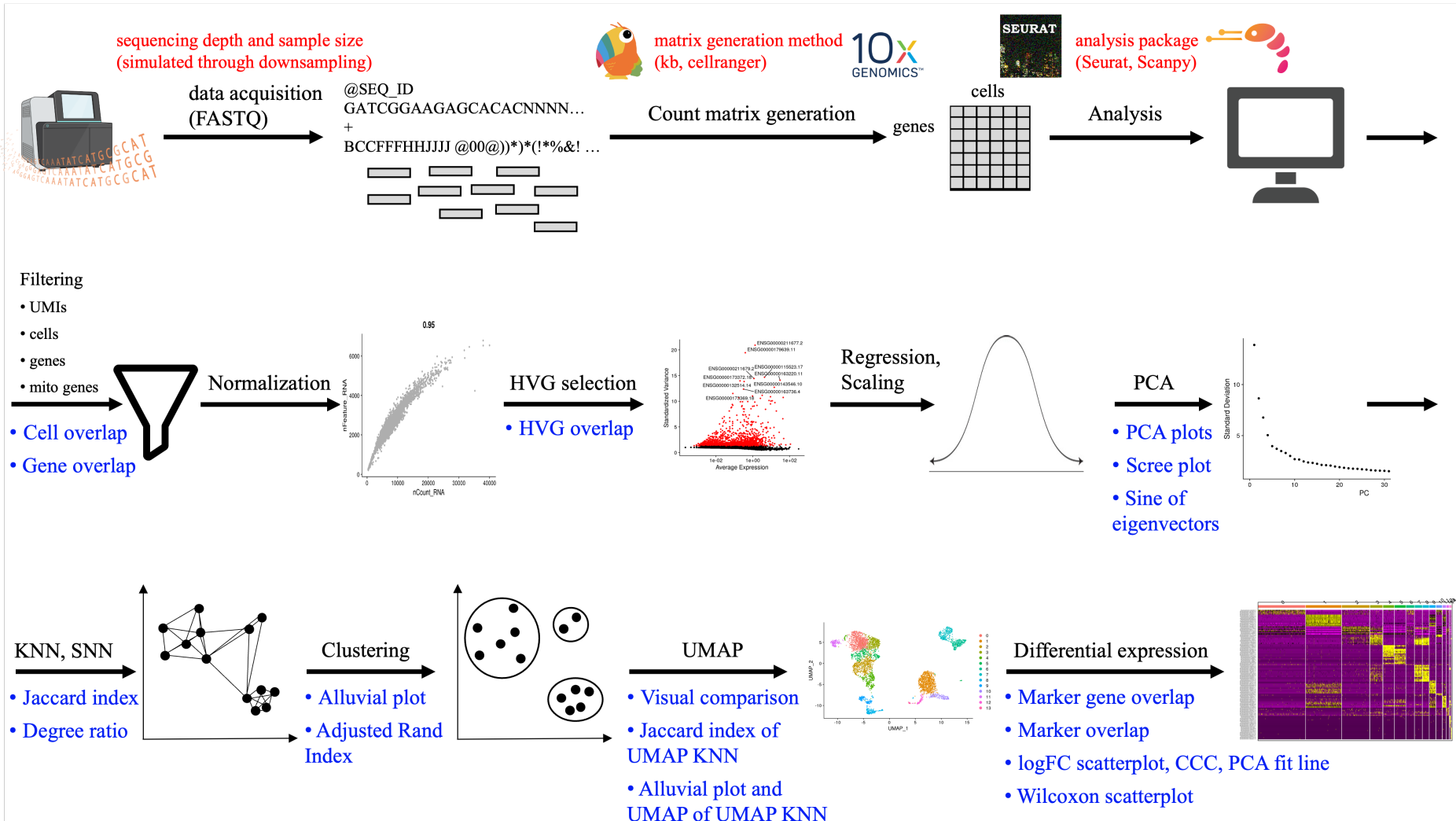
⁵*Boston, MA*

⁶*Computing and Mathematical Sciences, California Institute of Technology,
Pasadena, CA, 91125, USA*

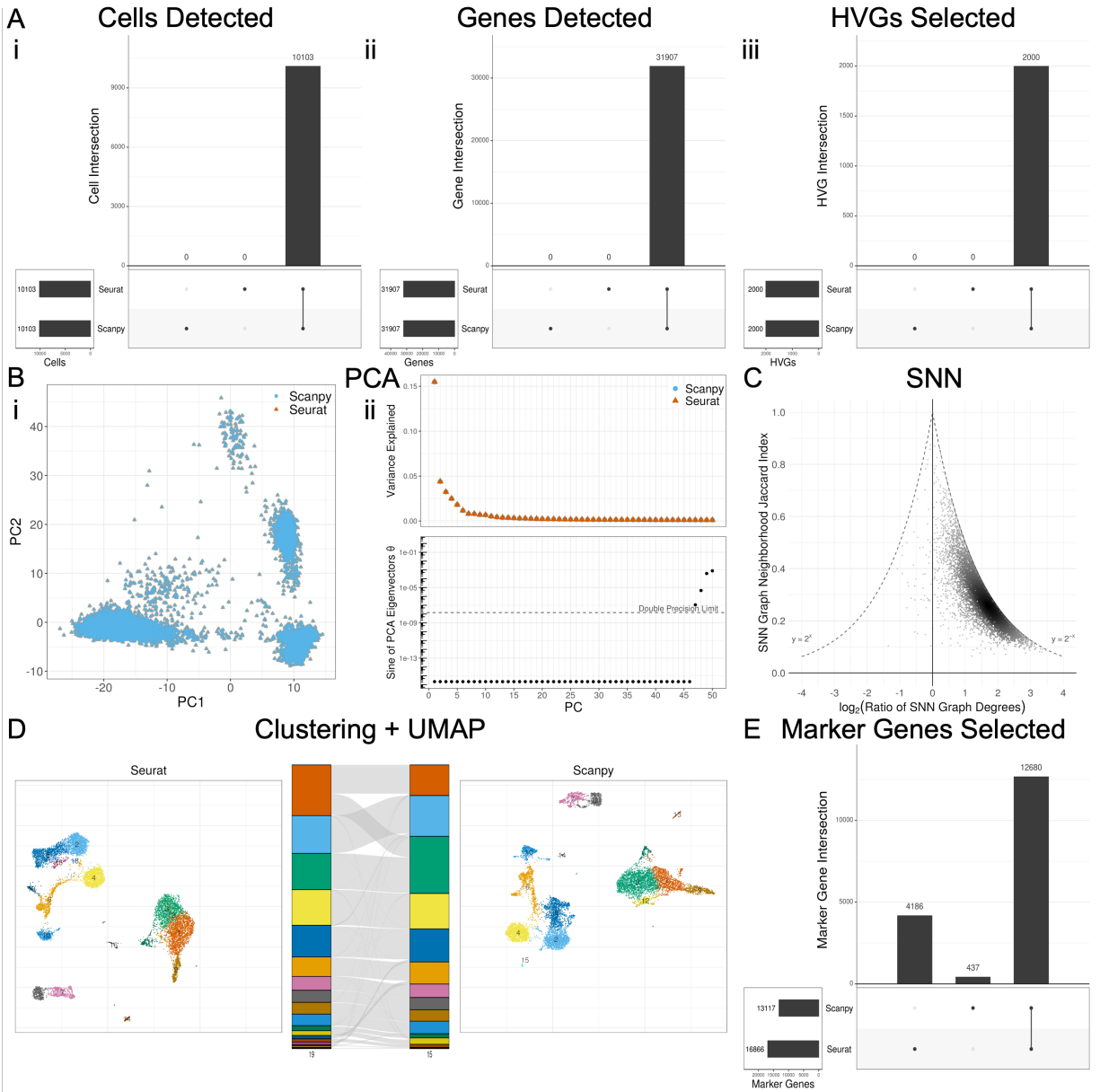
⁷*Lead Contact*

April 9, 2024

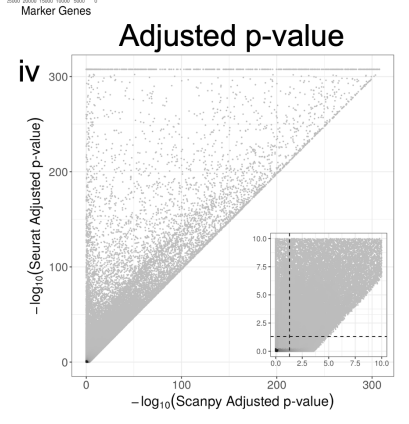
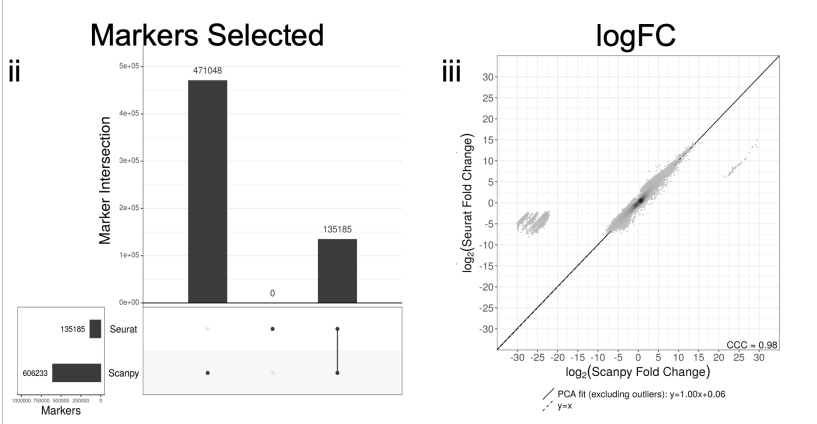
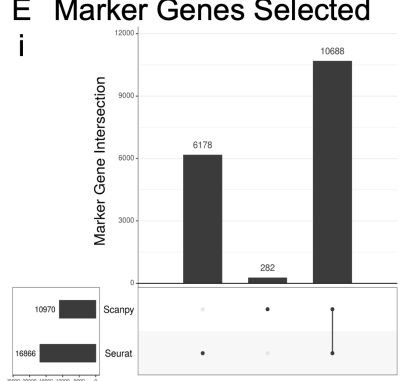
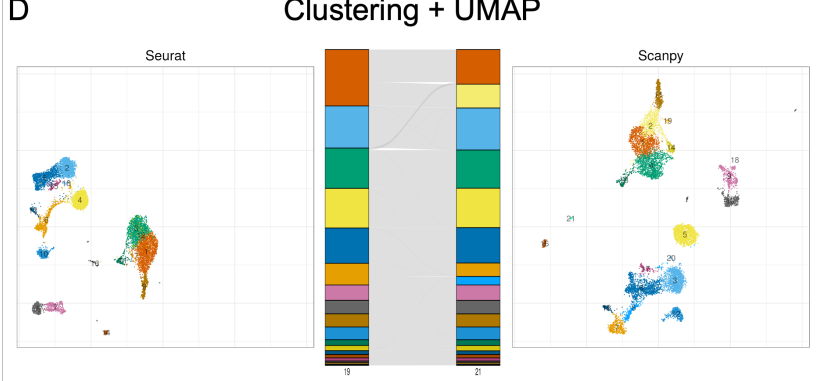
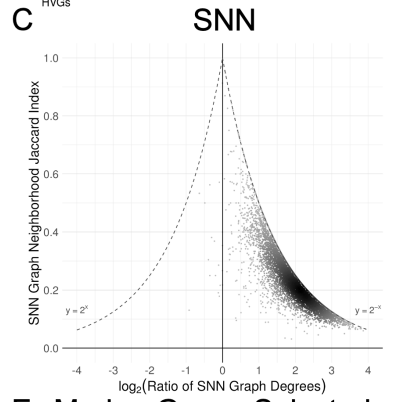
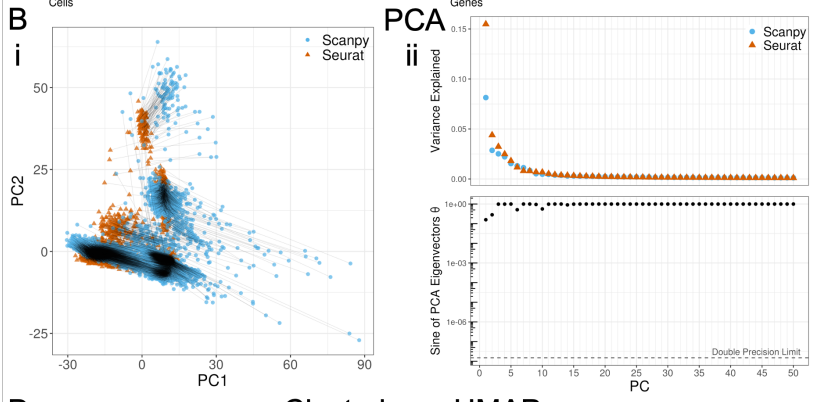
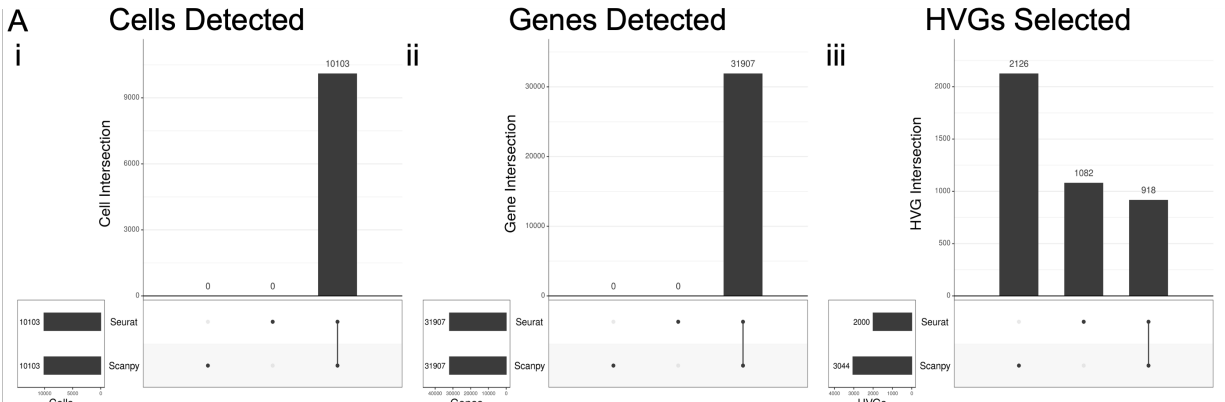
*Correspondence: lpachter@caltech.edu.



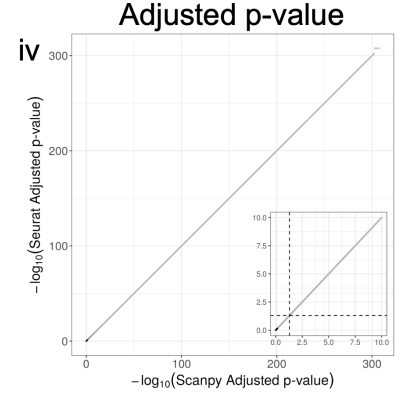
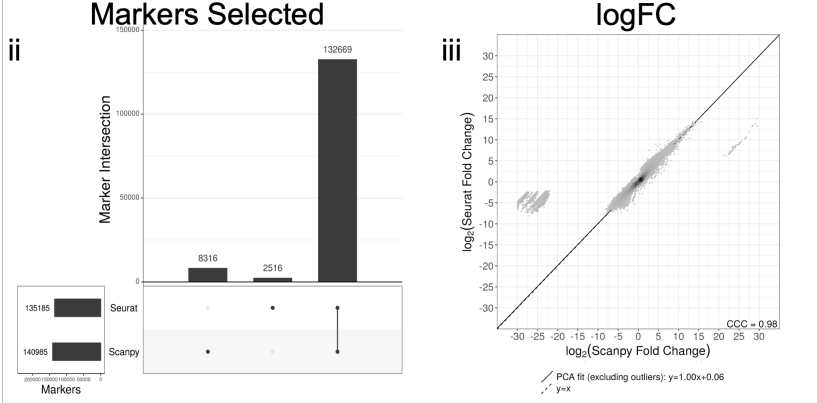
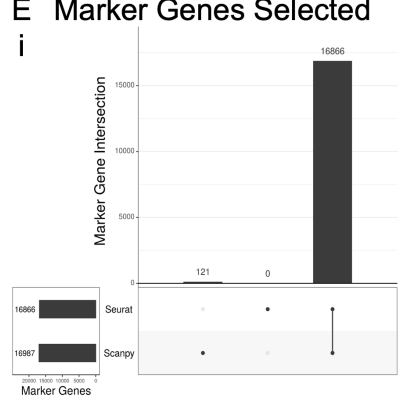
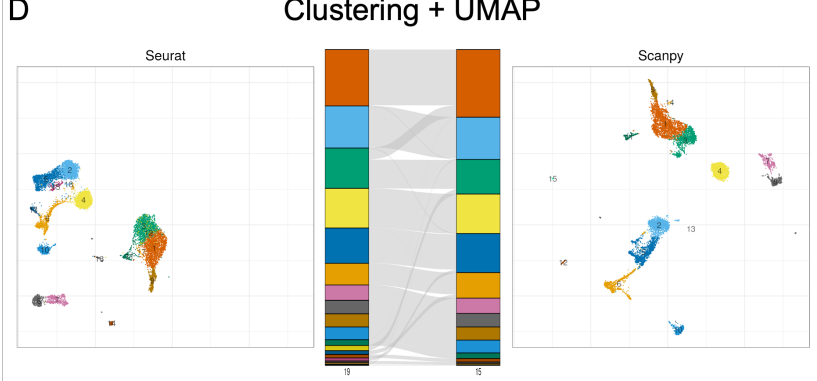
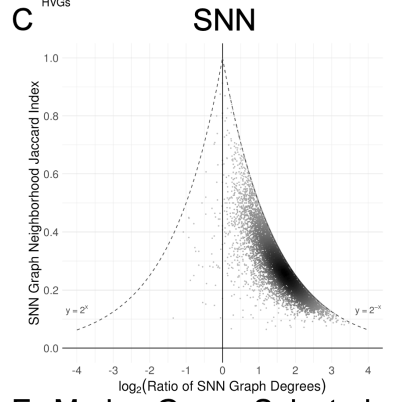
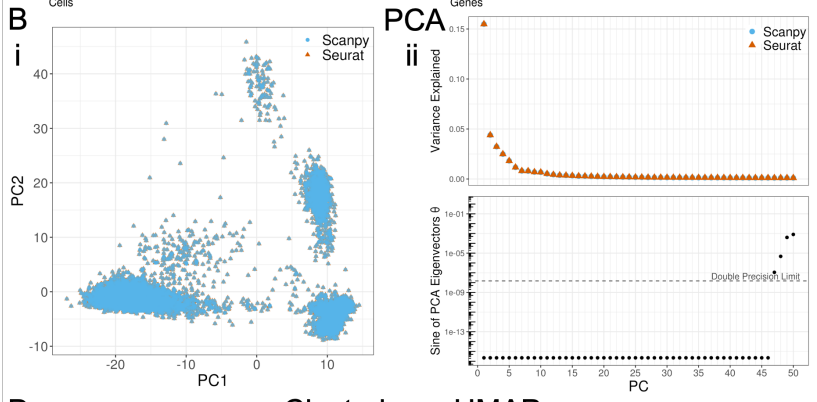
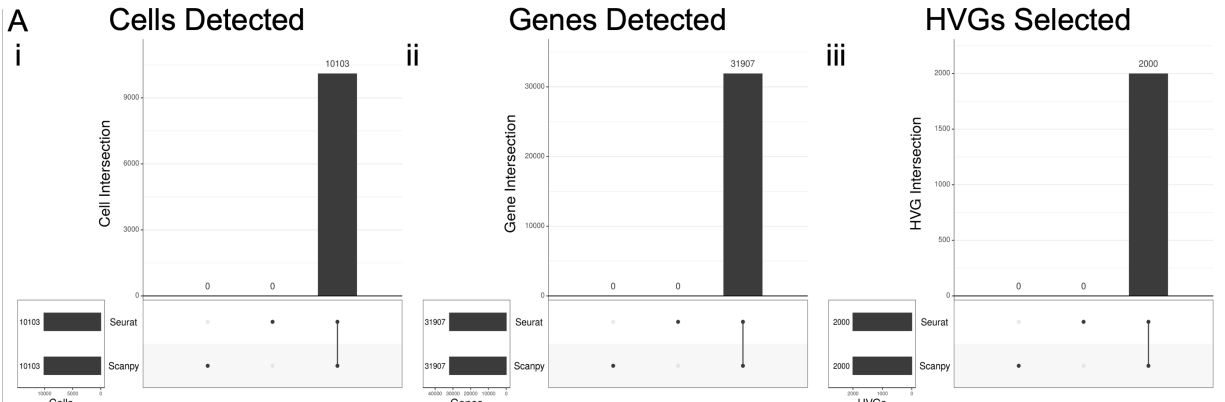
Supplemental Figure 1: scRNA-seq matrix generation and analysis overview. Red = sources of variability explored in this study; Blue = plots and metrics used to assess function similarity



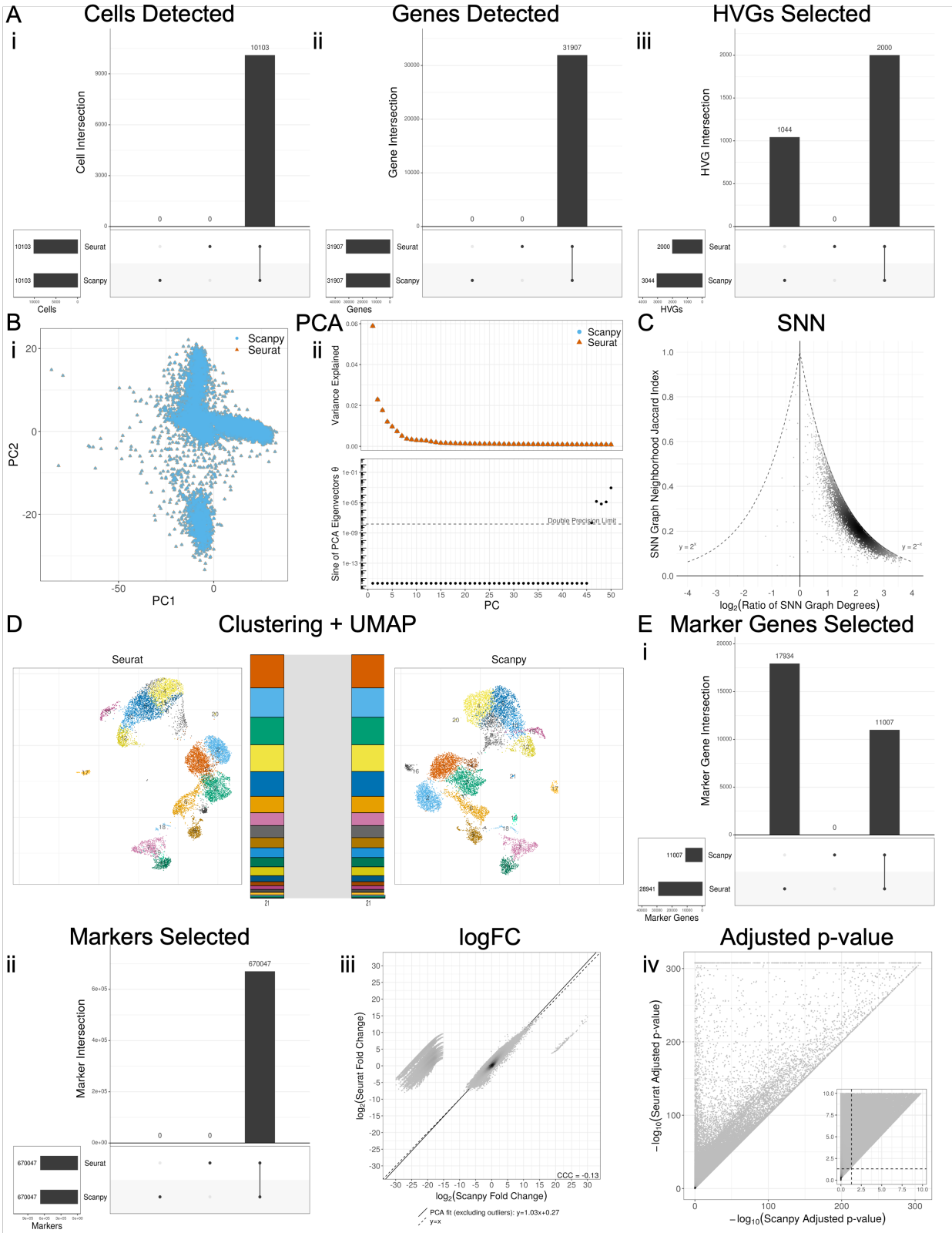
Supplemental Figure 2: Seurat vs. Scanpy, aligned function arguments (Seurat-like), sequential analysis. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (Seurat/Scanpy) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis UpSet plot through overlap of all significant ($p < 0.05$) marker genes across all clusters.



Supplemental Figure 3: Seurat vs. Scanpy, default function arguments, controlled analysis. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (Seurat/Scanpy) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters. (F) Differential expression analysis with aligned clusters of overlap of all markers. (G) Differential expression analysis with aligned clusters of logFC similarity. (H) Differential expression analysis with aligned clusters of adjusted p-value similarity.

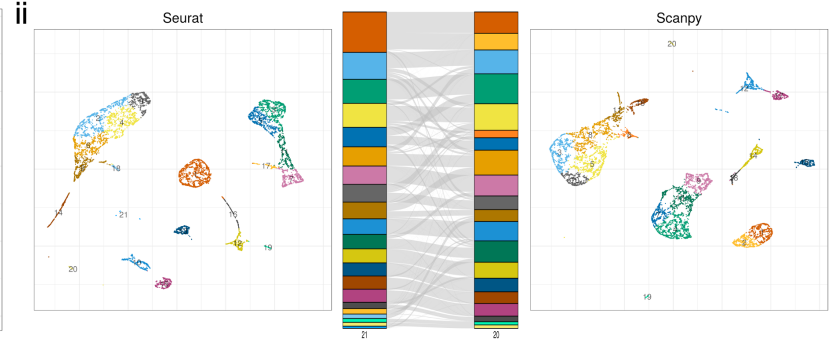
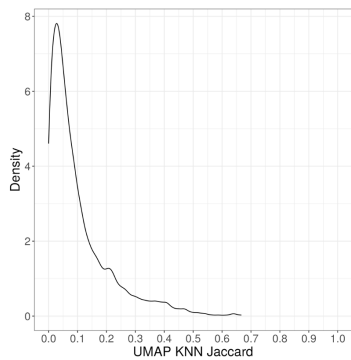


Supplemental Figure 4: Seurat vs. Scanpy, aligned function arguments (Seurat-like), controlled analysis. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (Seurat/Scanpy) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters. (F) Differential expression analysis with aligned clusters of overlap of all markers. (G) Differential expression analysis with aligned clusters of logFC similarity. (H) Differential expression analysis with aligned clusters of adjusted p-value similarity.

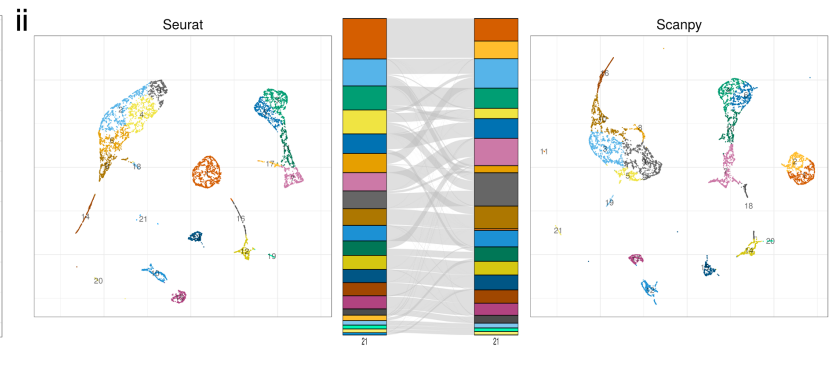
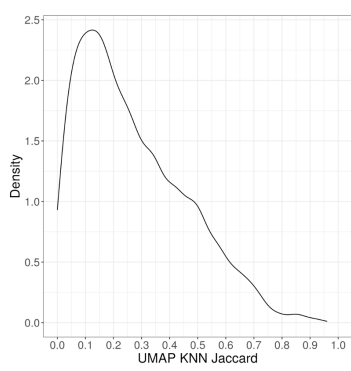


Supplemental Figure 5: Seurat vs. Scanpy, aligned function arguments (Scanpy-like), controlled analysis. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (Seurat/Scanpy) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters. (F) Differential expression analysis with aligned clusters of overlap of all markers. (G) Differential expression analysis with aligned clusters of logFC similarity. (H) Differential expression analysis with aligned clusters of adjusted p-value similarity.

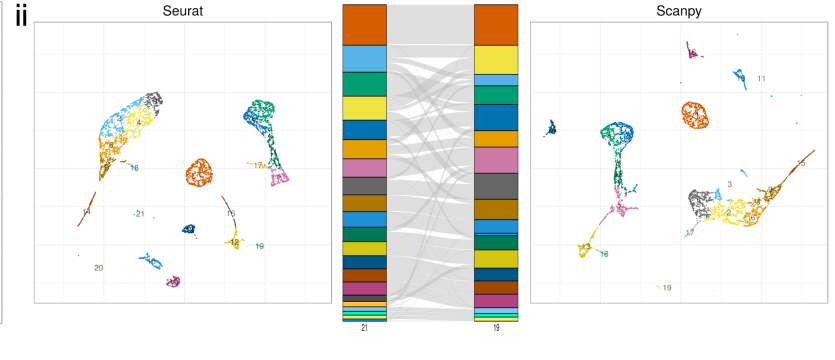
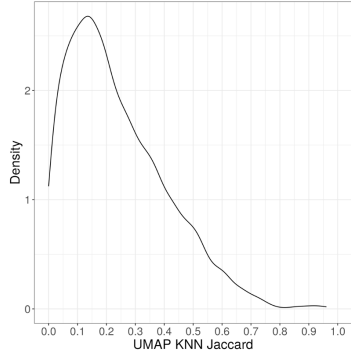
A
i
Defaults



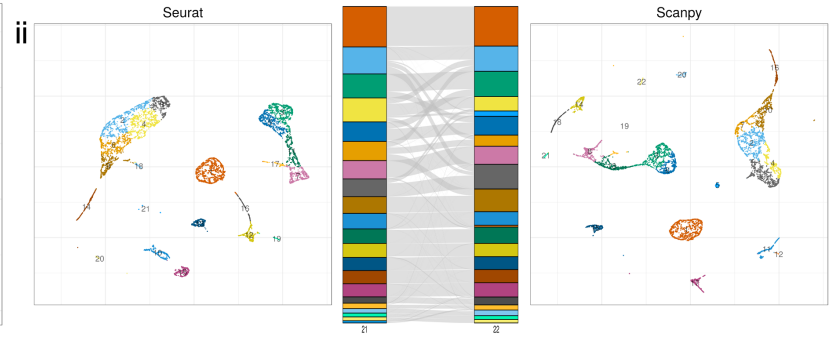
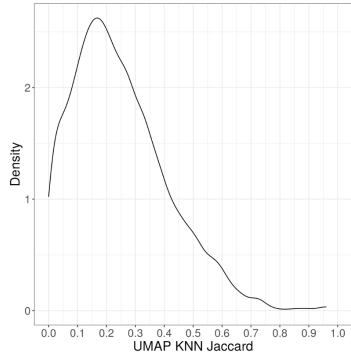
B
i
Aligned methods



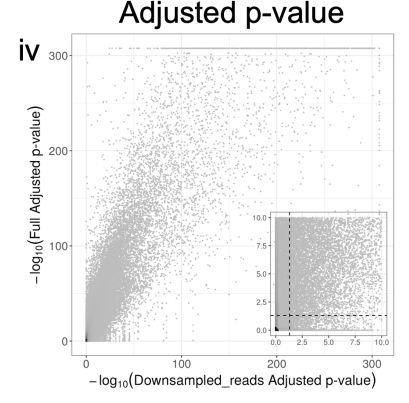
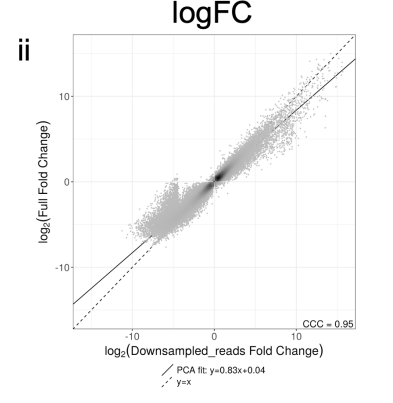
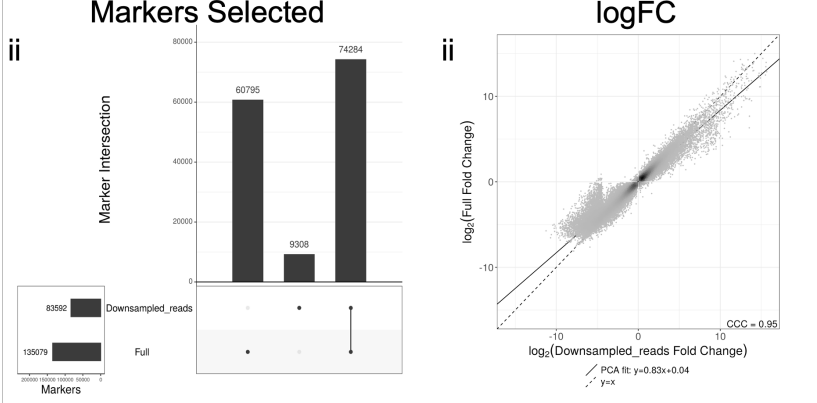
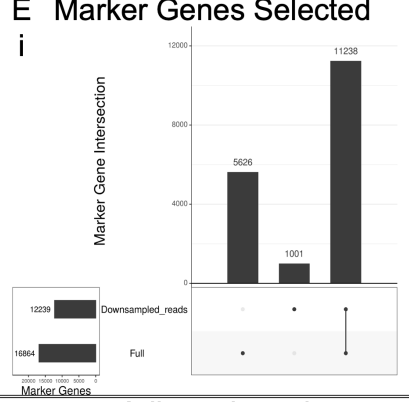
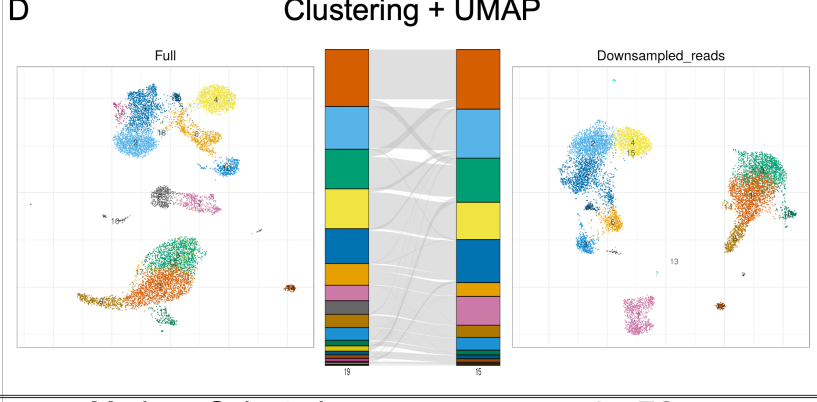
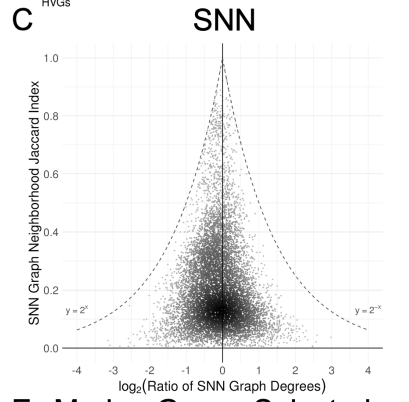
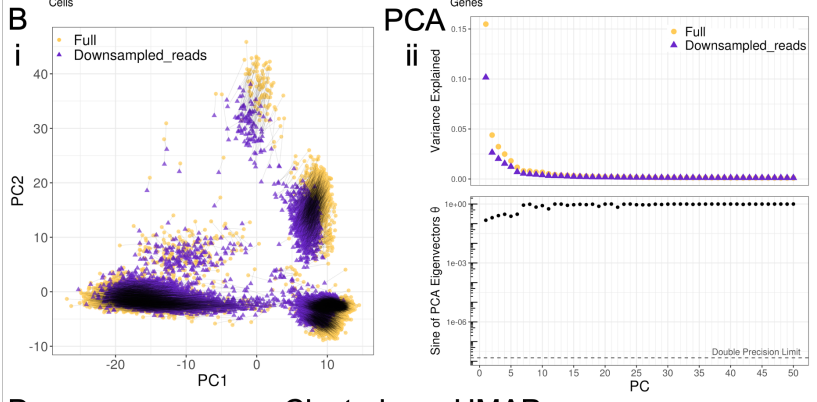
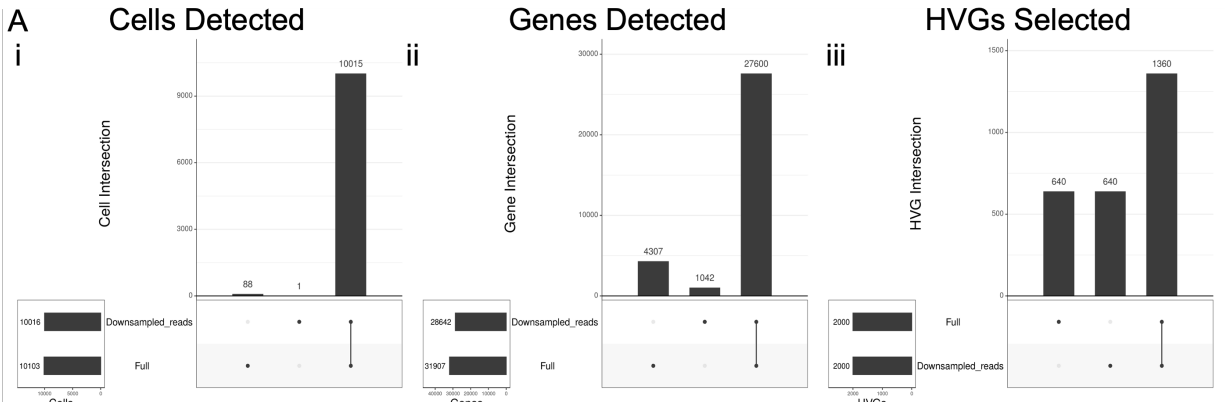
C
i
Controlled input



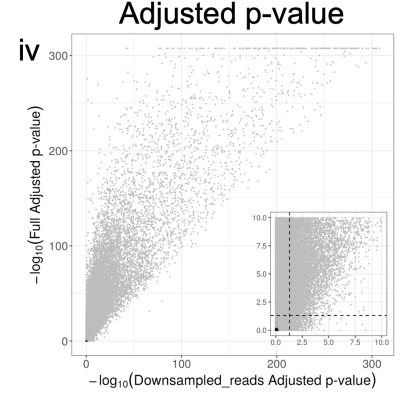
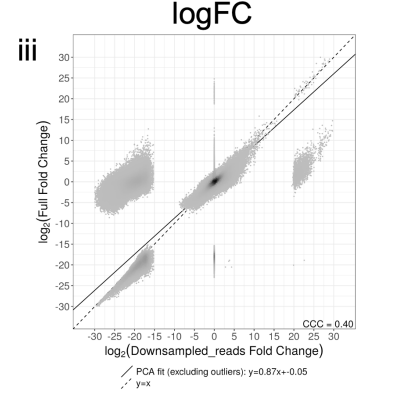
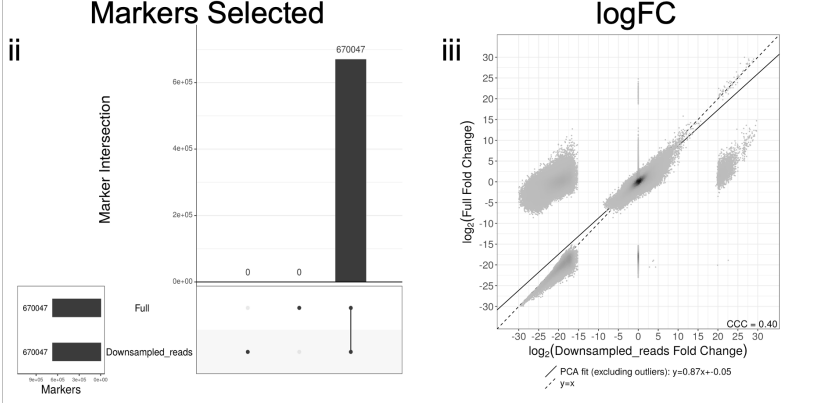
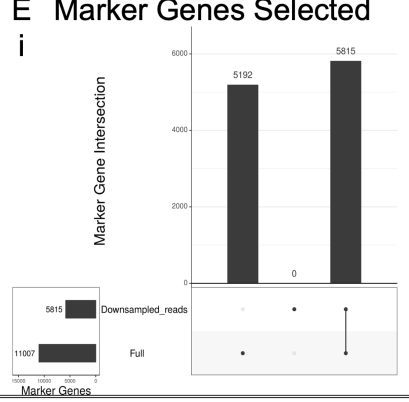
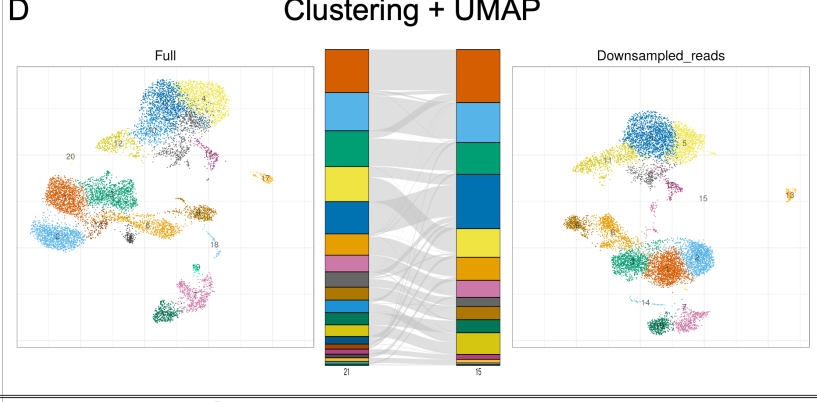
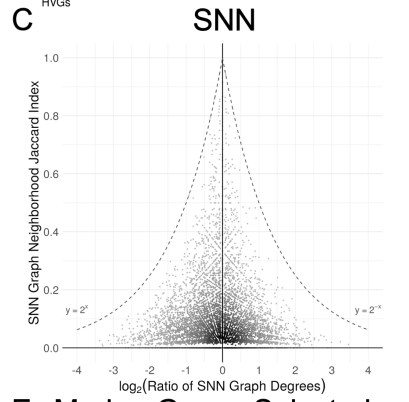
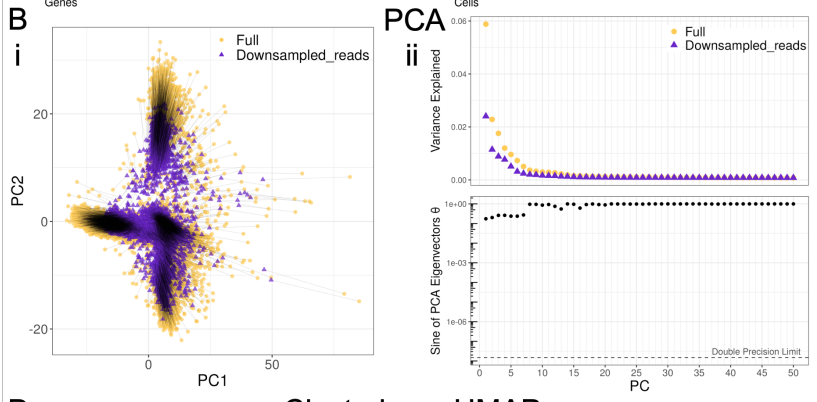
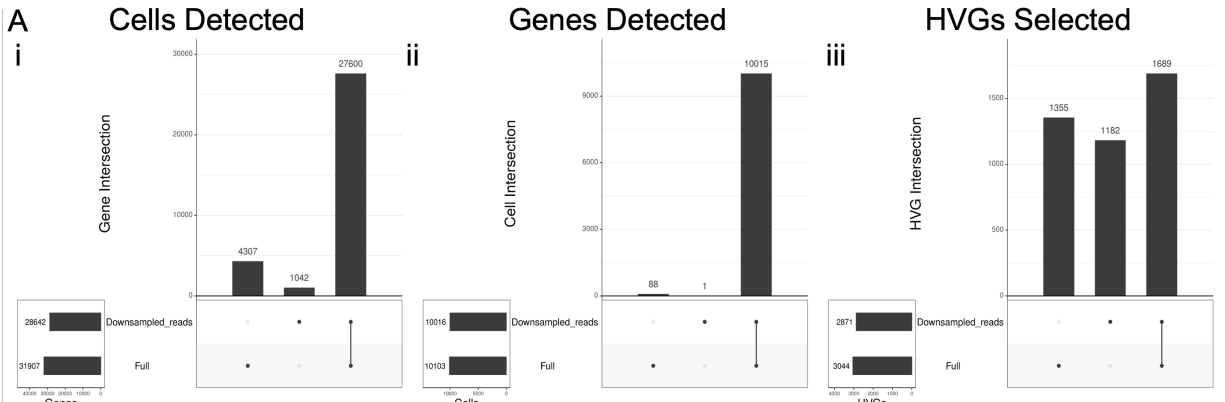
D
i
Aligned, controlled



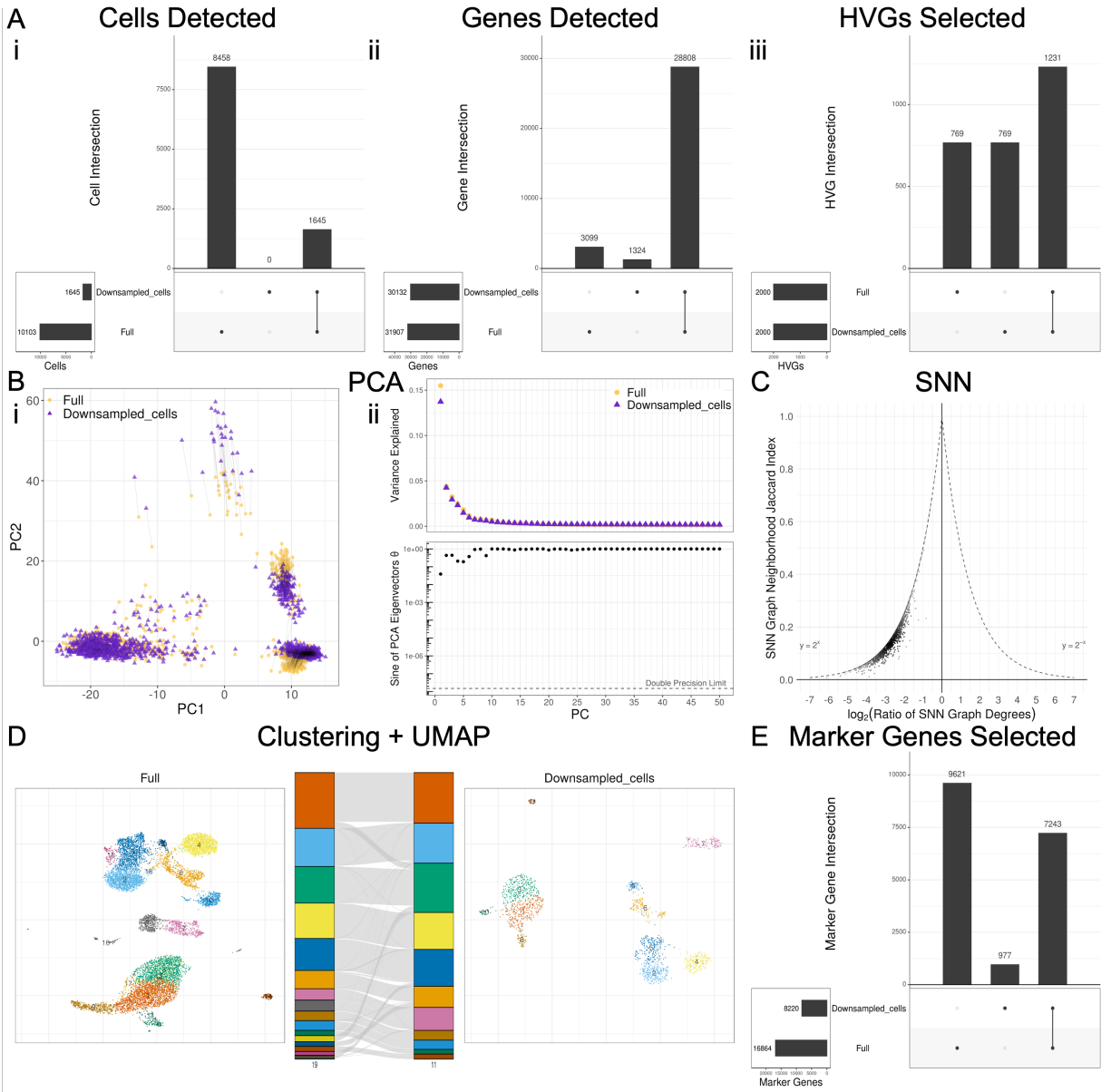
Supplemental Figure 6: Extended UMAP analysis for Seurat vs. Scanpy showing density plot of UMAP-derived KNN graph neighborhood Jaccard indices, alluvial plot of leiden clustering performed on UMAP-derived KNN graph, and UMAP projection of UMAP-derived KNN graph colored with leiden clustering results. (A) Default methods, sequential analysis. (B) Aligned methods, sequential analysis. (C) Default methods, controlled analysis. (D) Aligned methods, controlled analysis.



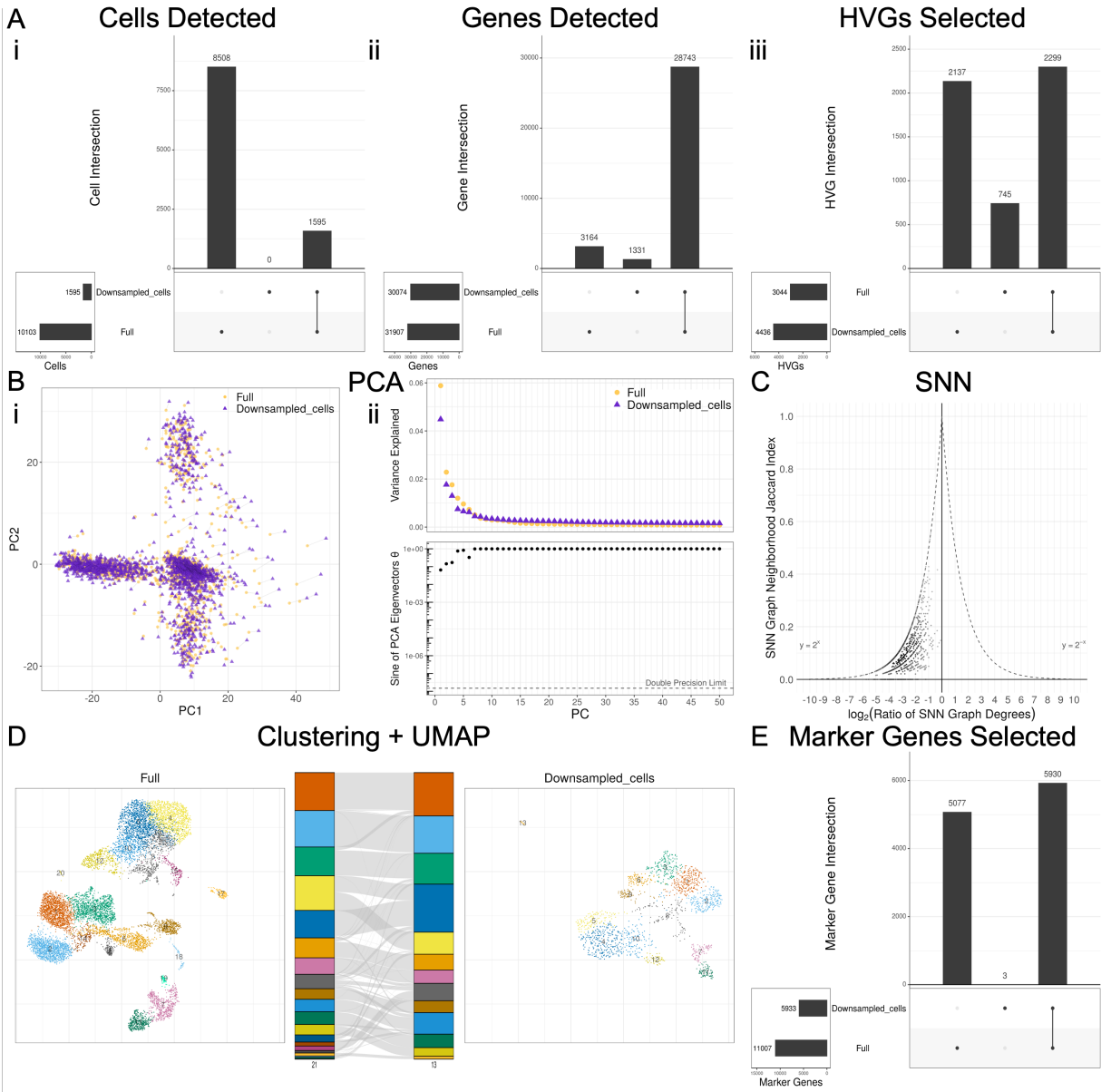
Supplemental Figure 7: Read downsampling to 4% of original size vs. full dataset in Seurat, with a fraction that displays variation comparable to Seurat vs. Scanpy. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (full/downsampled) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters. (F) Differential expression analysis with aligned clusters of overlap of all markers. (G) Differential expression analysis with aligned clusters of logFC similarity. (H) Differential expression analysis with aligned clusters of adjusted p-value similarity.



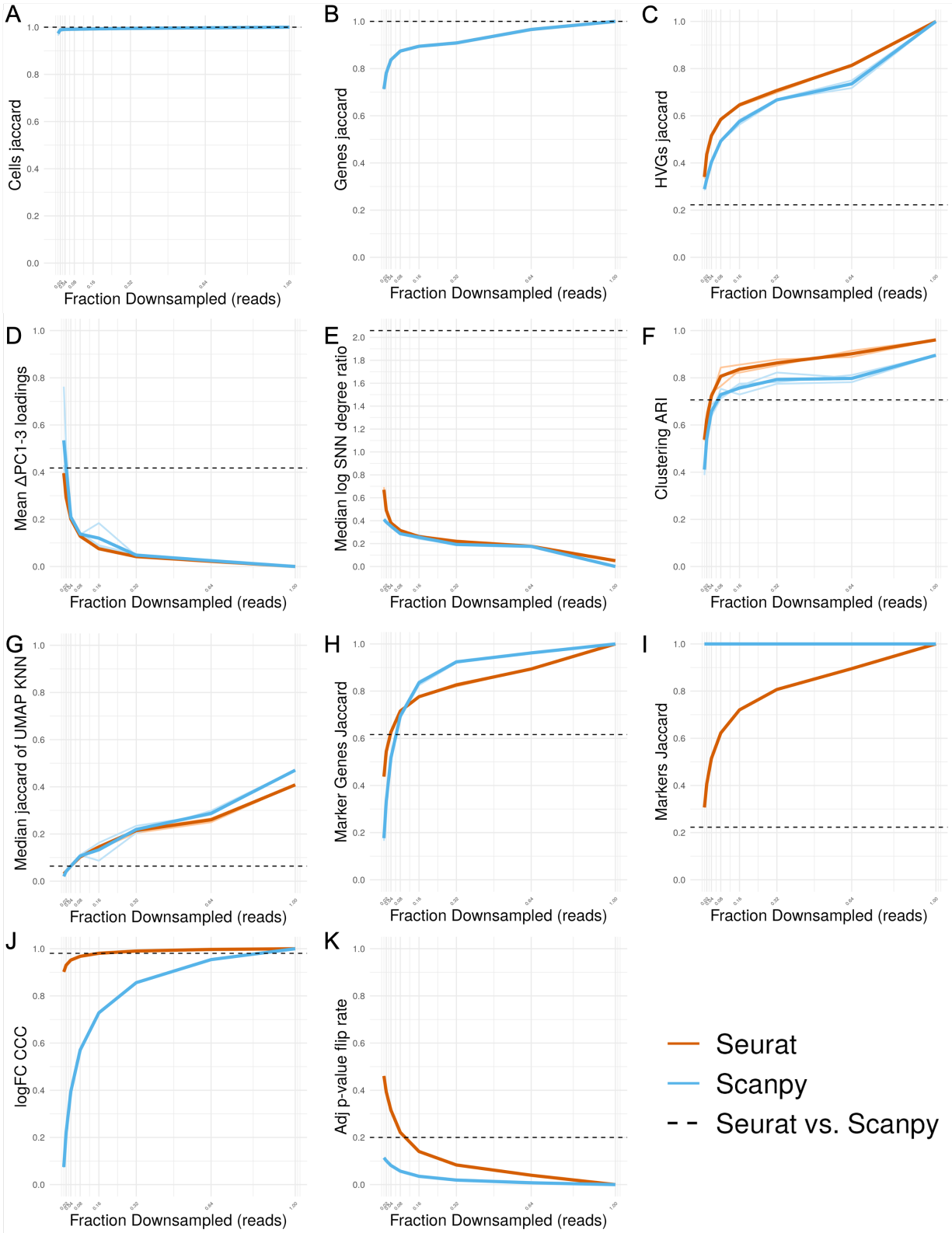
Supplemental Figure 8: Read downsampling to 4% of original size vs. full dataset in Scanpy, with a fraction that displays variation comparable to Seurat vs. Scanpy. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (full/-downsampled) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters. (F) Differential expression analysis with aligned clusters of overlap of all markers. (G) Differential expression analysis with aligned clusters of logFC similarity. (H) Differential expression analysis with aligned clusters of adjusted p-value similarity.



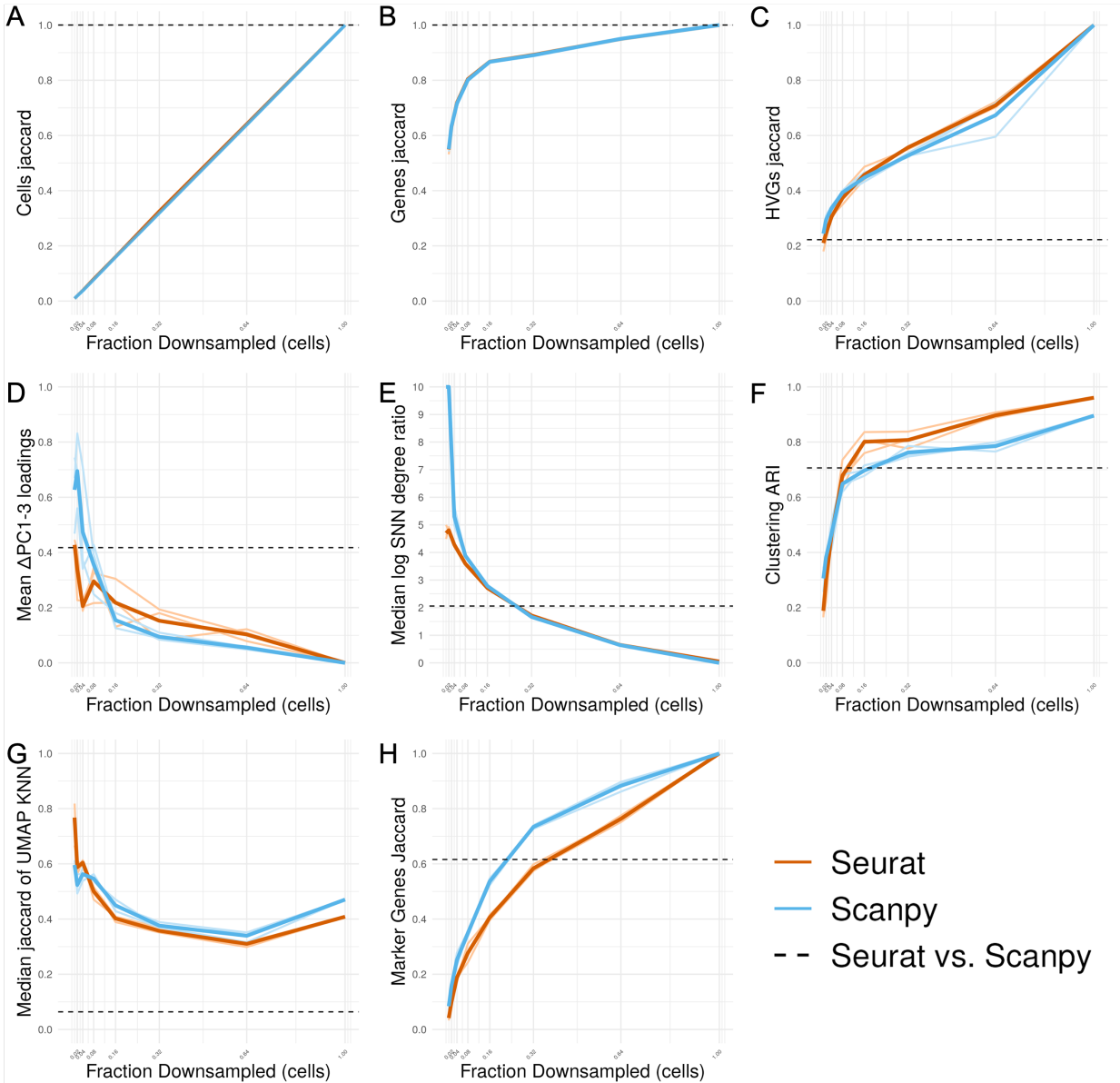
Supplemental Figure 9: Cell downsampling to 16% of original size vs. full dataset in Seurat, with a fraction that displays variation comparable to Seurat vs. Scanpy. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (full/downsampled) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis UpSet plot through overlap of all significant ($p < 0.05$) marker genes across all clusters.



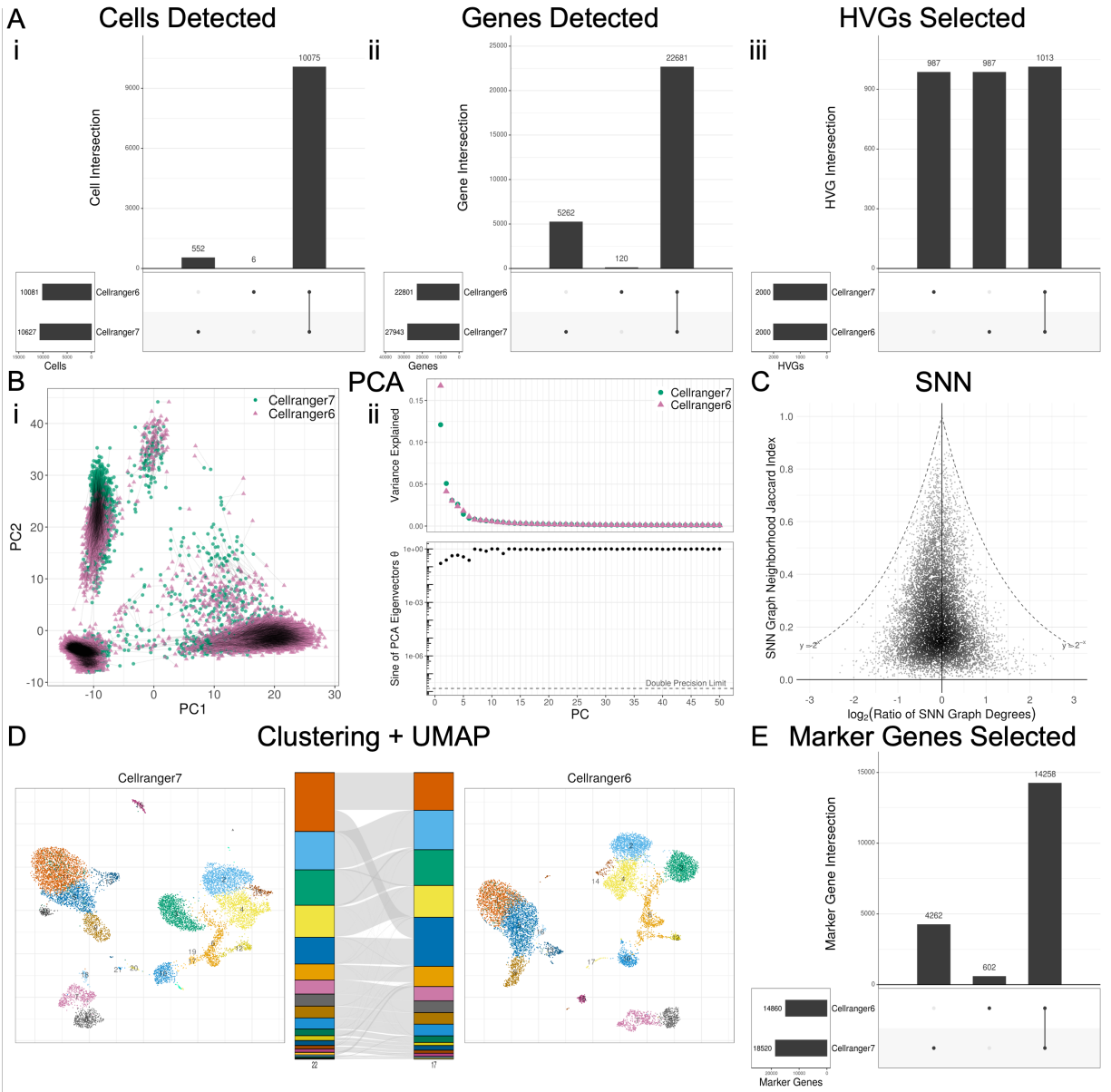
Supplemental Figure 10: Cell downsampling to 16% of original size vs. full dataset in Scanpy, with a fraction that displays variation comparable to Seurat vs. Scanpy. (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. Black lines = point mapping between conditions. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (full/downsampled) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis UpSet plot through overlap of all significant ($p < 0.05$) marker genes across all clusters.



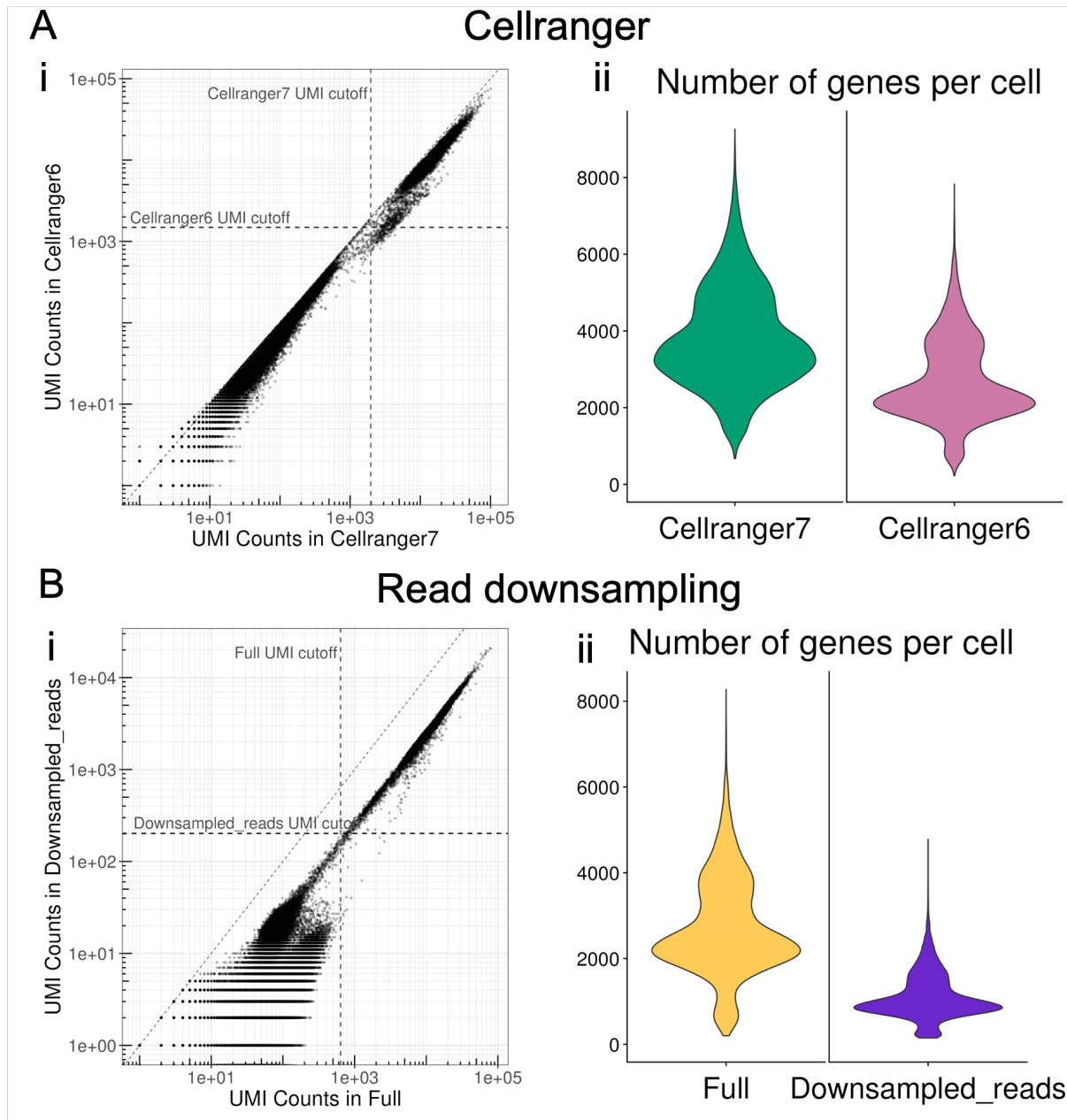
Supplemental Figure 11: Metrics vs. fraction of reads downsampled. (A) Jaccard index of cell overlap. (B) Jaccard index of gene overlap. (C) Jaccard index of HVG overlap. (D) Mean of difference in corresponding PC loadings for PCs 1-3. (E) Median log(SNN degree ratio). (F) Adjusted Rand Index of clusters. (G) Median Jaccard index of UMAP-derived KNN graphs for each cell. (H) Jaccard index of sets of significant ($p < 0.05$) marker genes across all clusters. (I) Jaccard index of all markers. (J) CCC of logFC scatterploz. (K) Fraction of adjusted p-value that flipped across $p=0.05$ threshold. Orange = Seurat; Blue = Scanpy; lighter lines = replicates of downsampling across random seeds; dashed line = value recorded by Seurat vs. Scanpy with default settings.



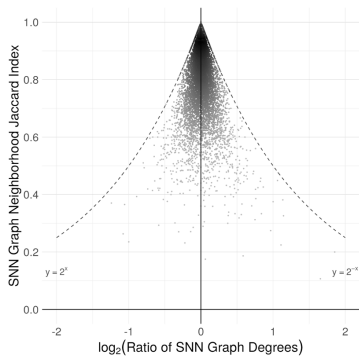
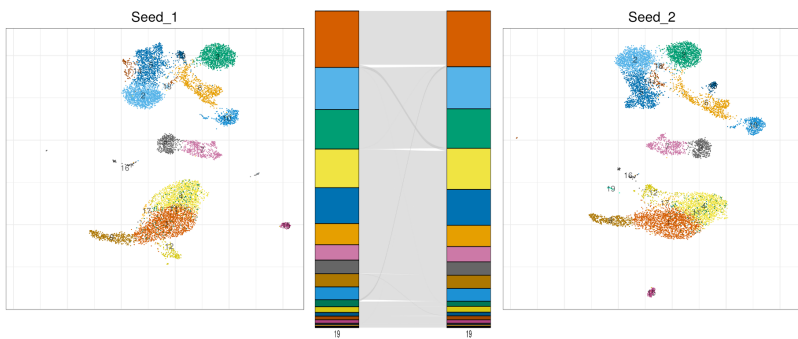
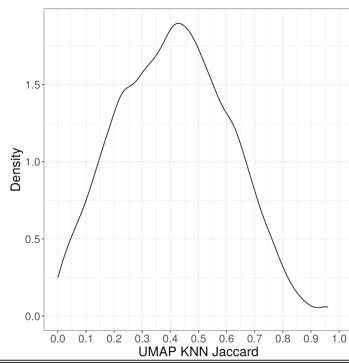
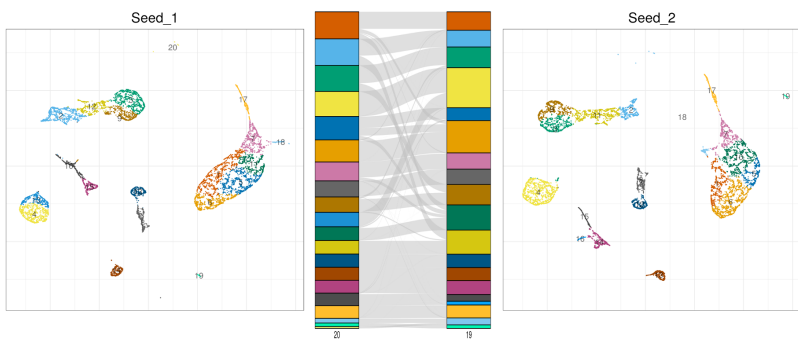
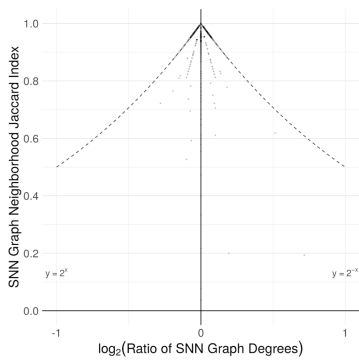
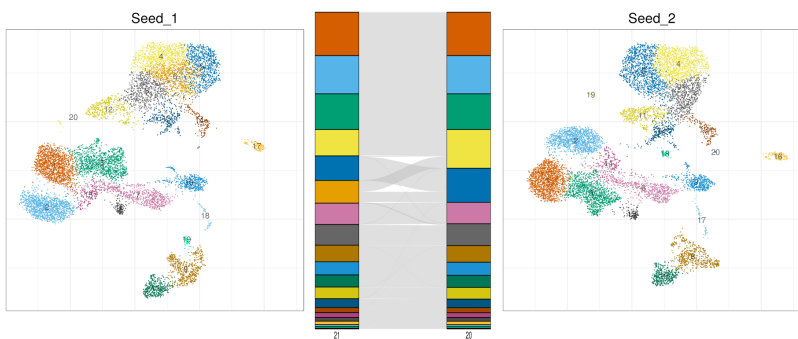
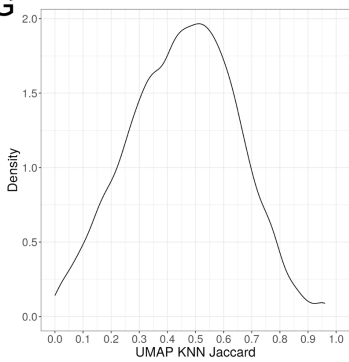
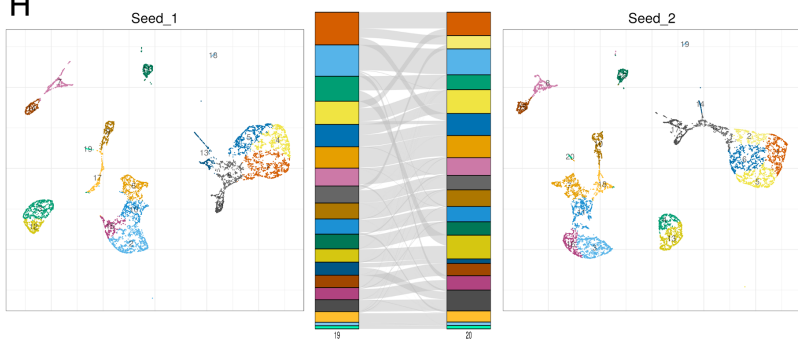
Supplemental Figure 12: Metrics vs. fraction of cells downsampled. (A) Jaccard index of cell overlap. (B) Jaccard index of gene overlap. (C) Jaccard index of HVG overlap. (D) Mean of difference in corresponding PC loadings for PCs 1-3. (E) Median log(SNN degree ratio). (F) Adjusted Rand Index of clusters. (G) Median Jaccard index of UMAP-derived KNN graphs for each cell. (H) Jaccard index of sets of significant ($p < 0.05$) marker genes across all clusters. (I) Jaccard index of all markers. (J) CCC of logFC scatterplot. (K) Fraction of adjusted p-value that flipped across $p=0.05$ threshold. Orange = Seurat; Blue = Scanpy; lighter lines = replicates of downsampling across random seeds; dashed line = value recorded by Seurat vs. Scanpy with default settings.



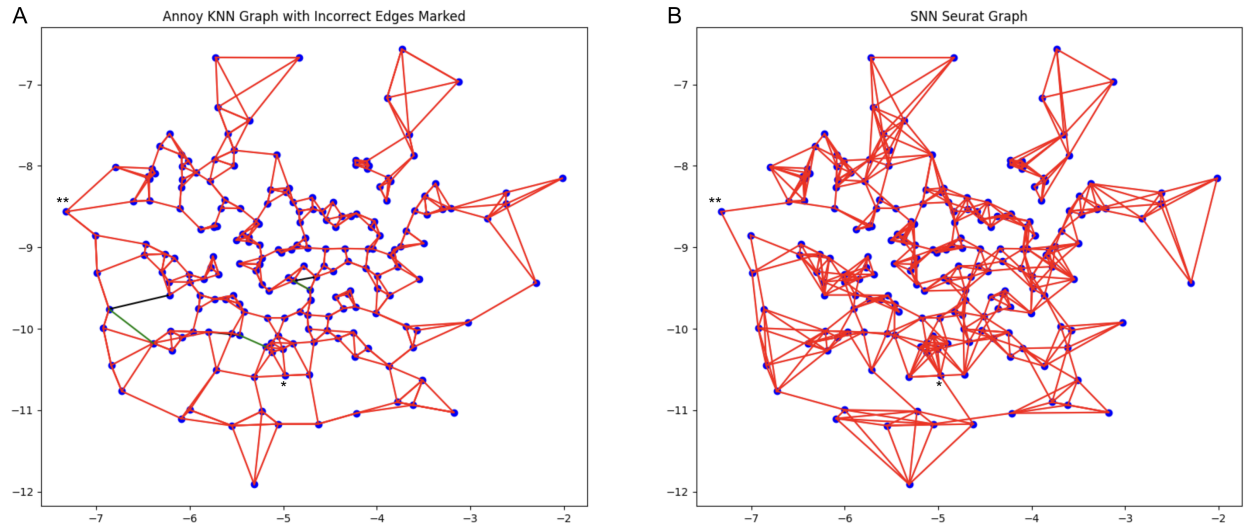
Supplemental Figure 13: Cell Ranger version control (v7 vs. v6). (A) Filtering and HVG selection analysis UpSet plots consisting of overlap of sets of cells, genes, and HVGs. (B) PCA analysis through projection onto first 2 PCs, Scree plot eigenvalue comparison, and sine of eigenvectors. (C) KNN/SNN analysis through SNN neighborhood Jaccard index and degree ratio (v7/v6) per cell. (D) Clustering and UMAP analysis through UMAP plots of each condition, with alluvial plot showing cluster assignment mapping and degree of agreement. (E) Differential expression analysis through overlap of all significant ($p < 0.05$) marker genes across all clusters.



Supplemental Figure 14: UMI counts pre-filtering and number of features under different conditions. (A) Cell Ranger v7 vs. v6. (B) Full size dataset vs. downsampled reads to 4%.

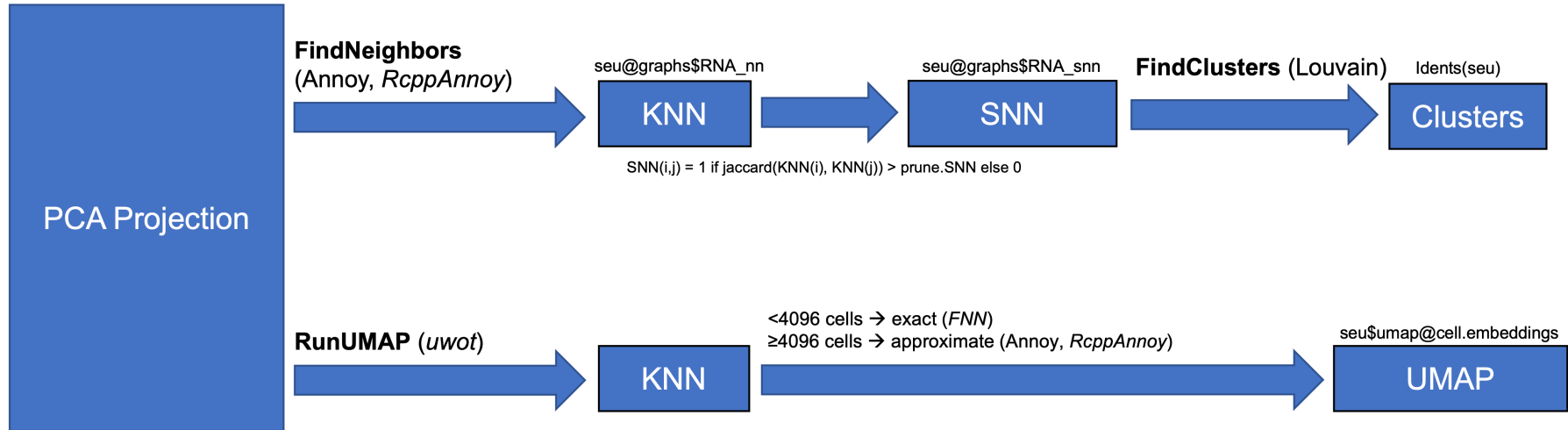
A**B****Seurat****C****D****E****F****Scanpy****G****H**

Supplemental Figure 15: Variability introduced by differences in random seeds for functions involving randomness. (A) Seurat SNN graph generation (Annoy), with different seeds for KNN graph generation. (B) Seurat clustering (Louvain), with different seeds for Louvain; and UMAP-projection (uwot), with different seeds for UMAP projection. (C) Seurat KNN graph generation from UMAP space, with different seeds for UMAP projection. (D) Seurat clustering (Leiden) and UMAP projection from UMAP space, with different seeds for UMAP projection. (E) Scanpy SNN graph generation (NNDescent), with different seeds for KNN graph generation. (F) Scanpy clustering (Leiden), with different seeds for Leiden; and UMAP-projection (UMAP-learn), with different seeds for UMAP projection. (G) Scanpy KNN graph generation from UMAP space, with different seeds for UMAP projection. (H) Scanpy clustering (Leiden) and UMAP projection from UMAP space, with different seeds for UMAP projection.

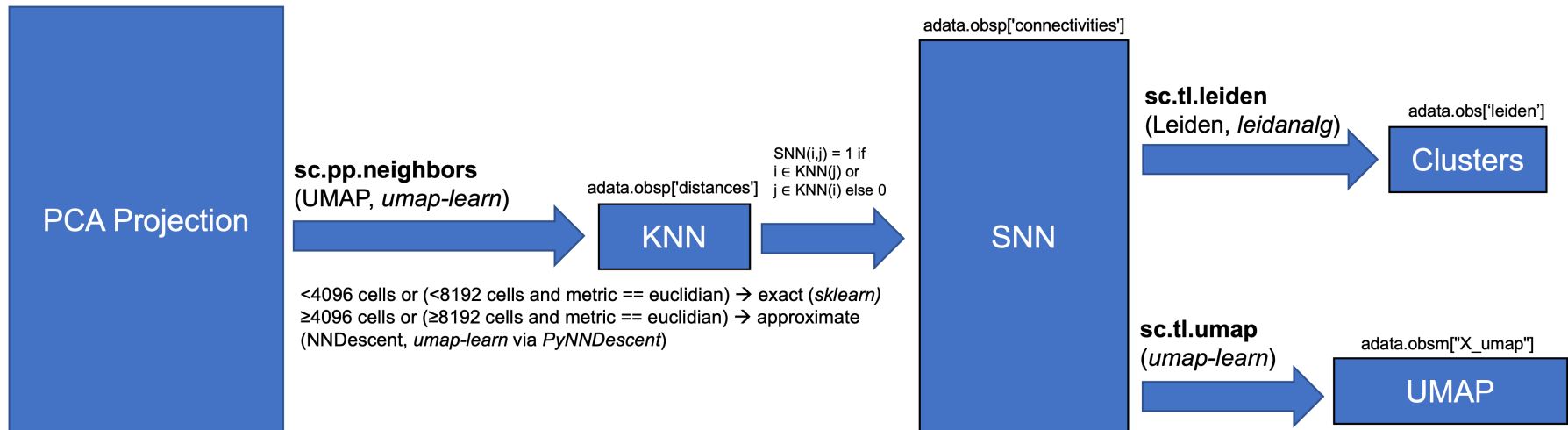


Supplemental Figure 16: Comparison of KNN and SNN graphs. (A) KNN graph ($k=3$) constructed with the approximate Annoy algorithm (directedness not shown). Black = incorrectly added edges; Green = incorrectly omitted edges. (B) SNN graph with Seurat's pruning method. An example hub node is marked with *. An example peripheral node is marked with **. Note: the SNN graph with Scanpy's method is simply the undirected KNN graph.

Seurat v5.0 (default)



Scanpy v1.9 (default)



Supplemental Figure 17: Schematic of protocol describing generation and downstream use of KNN and SNN graphs. Top is for Seurat, bottom is for Scanpy. **KNN(i)** refers to the full set of k-nearest neighbors for cell i. **Bold** = function/method name; *Italics* = outside package utilized.

Seurat v5.0.2	Filtering	Normalization	HVG selection	Regression, Scaling	PCA	SNN	Clustering	UMAP	DE
Function nam	CreateSeuratObject	NormalizeData	FindVariableFeatures	ScaleData	RunPCA	FindNeighbors	FindClusters	RunUMAP	FindAllMarkers
Defaults	<p>min.cells = 0</p> <p>min.features = 0</p>	<p><i>scale.factor = 10000</i></p> <p><i>normalization.method = "LogNormalize"</i></p>	<p>selection.method = "vst"</p> <p>loess.span = 0.3</p> <p><i>clip.max = "auto"</i></p> <p><i>mean.function = FastExpMean</i></p> <p><i>dispersion.function = FastLogVMR</i></p> <p>num.bin = 20</p> <p><i>binning.method = "equal_width"</i></p> <p>nfeatures = 2000</p> <p>mean.cutoff = c(0.1, 8)</p> <p>dispersion.cutoff = c(1, Inf)</p>	<p>vars.to.regress = NULL</p> <p><i>latent.data = NULL</i></p> <p><i>split.by = NULL</i></p> <p><i>model.use = "linear"</i></p> <p><i>use.umi = FALSE</i></p> <p><i>do.scale = TRUE</i></p> <p>do.center = TRUE</p> <p>scale.max = 10</p>	<p>npcs = 50</p> <p><i>rev.pca = FALSE</i></p> <p><i>weight.by.var = TRUE</i></p> <p><i>approx = TRUE</i></p> <p>features = NULL (HVGs)</p>	<p><i>reduction = "pca"</i></p> <p>dims = 1:10</p> <p>k.param = 20</p> <p><i>prune.SNN = 1/15</i></p> <p><i>nn.method = "annoy"</i></p> <p><i>n.trees = 50</i></p> <p>annoy.metric = "euclidean"</p> <p>nn.eps = 0</p> <p><i>l2.norm = FALSE</i></p>	<p><i>modularity.fxn = 1</i></p> <p><i>initial.membership = NULL</i></p> <p><i>node.sizes = NULL</i></p> <p>resolution = 0.8</p> <p>method = "matrix"</p> <p>algorithm = 1 (original Louvain)</p> <p><i>n.start = 10</i></p> <p><i>n.iter = 10</i></p> <p><i>group.singletons = TRUE</i></p>	<p>umap.method = "uwot"</p> <p><i>n.neighbors = 30L</i></p> <p>n.components = 2L</p> <p><i>metric = "cosine"</i></p> <p>n.epochs = NULL</p> <p>learning.rate = 1</p> <p>min.dist = 0.3</p> <p>spread = 1</p> <p><i>set.op.mix.ratio = 1</i></p> <p><i>local.connectivity = 1L</i></p> <p>repulsion.strength = 1</p> <p>negative.sample.rate = 5</p> <p>a = NULL</p> <p>b = NULL</p> <p>dens.lambda = 2</p> <p>dens.frac = 0.3</p> <p>dens.var.shift = 0.1</p>	<p><i>logfc.threshold = 0.1^a</i></p> <p><i>test.use = "wilcox"</i></p> <p><i>min.pct = 0.01^b</i></p> <p><i>min.diff.pct = -Inf</i></p> <p><i>only.pos = FALSE</i></p> <p><i>max.cells.per.ident = Inf</i></p> <p><i>latent.vars = NULL</i></p> <p><i>min.cells.feature = 3</i></p> <p><i>min.cells.group = 3</i></p> <p><i>mean.fxn = NULL</i></p> <p><i>return.thresh = 0.01</i></p> <p><i>base = 2</i></p>
(non-default) Arguments to match Scanpy defaults			<p>selection.method = "mean.var.plot"</p> <p>mean.cutoff = c(0.0125, 3)</p> <p>dispersion.cutoff = c(0.5, Inf)</p>	<p>vars.to.regress = c("nCount_RNA", "pct_mt")</p> <p>scale.max = Inf</p>		<p>k.param = 15</p> <p>dims = <i>total_pcs_used</i></p>	<p>algorithm = "leiden"</p> <p>resolution = 1.0</p>	<p>umap.method = "umap-learn"</p> <p>min.dist = 0.5</p>	<p>logfc.threshold = 0</p> <p>min.pct = 0</p> <p>return.thresh = 1.0</p> <p>min.cells.group = 1</p>

^adefault changed from Seurat v4 (was logfc.threshold = 0.25)

^bdefault changed from Seurat v4 (was min.pct = 0.1)

Supplemental Table 1: Seurat function details for scRNA-seq pipeline agreement with Scanpy. Bold = parameter in agreement by default; italics = no analogous parameter in Scanpy; green = equivalent by default; yellow = equivalent with matched arguments; orange = partially incompatible, partially equivalent with matched arguments (depending on algorithm choice); red = incompatible

Scanpy v1.9.5	Filtering	Normalization	HVG selection	Regression, Scaling	PCA	SNN	Clustering	UMAP	DE
Function name	sc.pp.filter_cells, sc.pp.filter_genes	sc.pp.normalize_total	sc.pp.highly_variable_genes	sc.pp.regress_out, sc.pp.scale	sc.tl.pca	sc.pp.neighbors	sc.tl.leiden, sc.tl.louvain	sc.tl.umap	sc.tl.rank_genes_groups
Defaults	<p>min_cells = None</p> <p>min_genes = None</p> <p>max_cells = None</p> <p>max_genes = None</p> <p>min_counts = None (genes)</p> <p>min_counts = None (cells)</p> <p>max_counts = None (genes)</p> <p>max_counts = None (cells)</p>	<p>target_sum = None</p> <p>exclude_highly_expressed = False</p> <p>max_fraction = 0.05</p>	<p>n_top_genes = None</p> <p>min_mean = 0.0125</p> <p>max_mean = 3</p> <p>min_disp = 0.5</p> <p>max_disp = inf</p> <p>span = 0.3</p> <p>n_bins = 20</p> <p>flavor = "seurat"</p>	<p>zero_center = True</p> <p>max_value = None</p>	<p>n_comps = 50 or 1 – minimum dimension size</p> <p>zero_center = True</p> <p>svd_solver = "arpack"</p> <p>use_highly_variable = None (HVGs)</p>	<p>n_neighbors = 15</p> <p>n_pcs = None</p> <p>use_rep = None</p> <p>knn = True</p> <p>method = "umap"</p> <p>metric = "euclidian"</p>	<p>resolution = 1</p> <p>restrict_to = None</p> <p>directed = True</p> <p>use_weights = True</p> <p>n_iterations = -1</p> <p>partition_type = None</p> <p>flavor = "vtraag" (louvain)</p>	<p>min_dist = 0.5</p> <p>spread = 1.0</p> <p>n_components = 2</p> <p>maxiter = None</p> <p>alpha = 1.0</p> <p>gamma = 1</p> <p>negative_sample_rate = 5</p> <p>init_pos = "spectral"</p> <p>a = None</p> <p>b = None</p> <p>method = "umap"</p>	<p>use_raw = None</p> <p>reference = "rest"</p> <p>n_genes = None</p> <p>rankby_abs = False</p> <p>pts = False</p> <p>method = None</p> <p>corr_method = 'benjamini-hochberg'</p> <p>tie_correct = False</p>
(non-default) Arguments to match Seurat defaults			<p>n_top_genes = 2000</p> <p>flavor = "seurat_v3"</p>	<p>Omit sc.pp.regress_out</p> <p>max_value = 10</p>	<p>zero_center = False</p>	<p>n_neighbors = 20</p> <p>n_pcs = 10</p> <p>use_rep = 'X_pca'</p>	<p>function sc.tl.louvain instead of sc.tl.leiden</p> <p>resolution = 0.8</p> <p>n_iterations = 10</p>	<p>min_dist = 0.3</p>	<p>use_raw=True</p> <p>pts = True</p> <p>method = "wilcoxon"</p> <p>corr_method = "bonferroni"</p> <p>tie_correct = True</p> <p>manual filtering (abs(log_fc) ≥ 0.1, pts ≥ 0.01, pts_rest ≥ 0.01, p_value < 0.01)</p>

Supplemental Table 2: Scanpy function details for scRNA-seq pipeline agreement with Seurat. Bold = parameter in agreement by default; italics = no analogous parameter in Seurat; green = equivalent by default; yellow = equivalent with matched arguments; orange = partially incompatible, partially equivalent with matched arguments (depending on algorithm choice); red = incompatible