

1 *Building the Mega Single Cell Transcriptome Ocular Meta-Atlas*

2 Vinay S Swamy¹, Temesgen D Fufa², Robert B Hufnagel², and David M McGaughey^{1,□}

3 March 26, 2021

4 The development of highly scalable single cell transcriptome technology has resulted in the creation of thousands of
5 datasets, over 30 in the retina alone. Analyzing the transcriptomes between different projects is highly desirable as
6 this would allow for better assessment of which biological effects are consistent across independent studies.
7 However it is difficult to compare and contrast data across different projects as there are substantial batch effects
8 from computational processing, single cell technology utilized, and the natural biological variation. While many
9 single cell transcriptome specific batch correction methods purport to remove the technical noise it is difficult to
10 ascertain which method functions works best. We developed a lightweight R package (scPOP) that brings in batch
11 integration methods and uses a simple heuristic to balance batch merging and celltype/cluster purity. We use this
12 package along with a Snakefile based workflow system to demonstrate how to optimally merge 766,615 cells from
13 34 retina datasets and three species to create a massive ocular single cell transcriptome meta-atlas. This provides a
14 model how to efficiently create meta-atlases for tissues and cells of interest.

15 ¹ Bioinformatics Group, Ophthalmic Genetics & Visual Function Branch, National Eye Institute, National

16 Institutes of Health

17 ² Medical Genetics and Ophthalmic Genomics Unit, National Eye Institute, National Institutes of Health

18 □ Correspondence: [David M McGaughey <mcgaugheyd@mail.nih.gov>](mailto:mcgaugheyd@mail.nih.gov)

19 **Introduction**

20 **A plethora of single-cell transcriptome studies in the retina**

21 The retina contains a multitude of cell types that, in total, are responsible for turning light information into
22 signal for the brain to interpret as vision. Very briefly, the photoreceptors (rods and cones) are responsible for
23 capturing the photons. The retinal pigmented epithelium (RPE) behind the photoreceptors physically support the
24 rods and cones by processing byproducts of the visual cycle. Müller glia serve as support cells for the neurons. The
25 retinal bipolar cells transmit the electrical signal from the photoreceptors to the retinal ganglion cells. Horizontal and
26 amacrine regulate and help interpret signals from the photoreceptors. The signal is relayed via the retinal ganglion
27 projections through the optic stalk to the brain (for review see¹). Since 2000 many groups have investigated gene
28 expression in small numbers of individual cells of the retina.²⁻⁵

29 The recent introduction of lower cost and high throughput single cell sequencing technology has led to an
30 explosion of research across many fields. As of early 2021, over 40 million cells have been sequenced across over
31 1,200 studies and the average size of each study starting in 2020 is over 100,000 cells.⁶ The retina was used as the
32 source tissue in one of the earliest works in the high throughput single cell transcriptomics field.⁷ As of late 2020,
33 over a twenty published studies, cumulatively containing over a million cells, have used single cell technology to
34 profile cell type specific gene expression patterns, cell fate trajectory, tissue and cell differentiation, and disease
35 perturbation across multiple mammalian species.⁷⁻²⁹

36 While the gene - cell count tables are generally made available in repositories like the Gene Expression
37 Omnibus (GEO), there are no requirements to uniformly process the data. This means the count tables cannot be
38 used in cross-study comparisons as even small differences in the computational pipeline (aligner, transcriptome
39 reference, etc.) create study-specific effects. This issue can be addressed only by re-quantifying the data in a uniform
40 computational environment. Fortunately, due to the continued development of computationally light-weight gene
41 quantification tools in the single-cell space (e.g kallisto bustools, alevin-fry), re-quantification does not require
42 massive compute and time resources.^{30,31}

43 Still, even after re-quantification under identical computational conditions there remain study specific batch
44 effects due to the diversity in single cell technologies used and variation in tissue handling and processing across
45 each scientific group. The single cell community has recognized that removal of these technical (also referred to as
46 batch) effects is a critical issue and have independently developed many tools, though it remains unclear which tools
47 and parameters are optimal for a particular dataset.³²⁻⁴³

48 **The projectable meta-atlas**

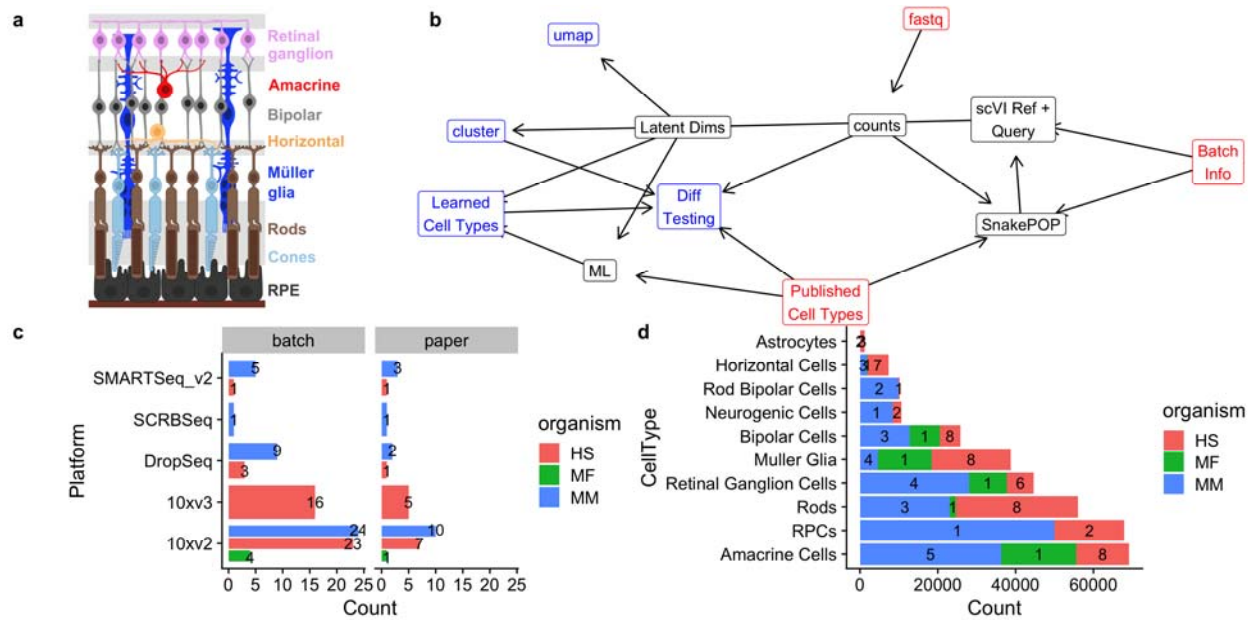
49 We propose that by re-processing publicly available raw single cell transcriptome data in a consistent
50 bioinformatic framework and optimally using batch correction tools we can create a meta-atlas of retina single cell
51 transcriptomes. As there are thousands of possible permutations of single cell tools, references, and parameter
52 choices, we create our meta-atlas (which we refer to as the single cell Eye in a Disk or scEiaD) by benchmarking
53 integration outcomes across multiple important single cell RNA-seq processing parameters (batch removal method,
54 number of hyper-variable genes (HVGs), clustering resolution, etc.). The benchmarking system we developed uses a
55 wide variety of metrics that combine in the R package scPOP (single-cell Pick Optimal Parameters). The scEiaD

56 will be of utility to two communities. First, the ocular community who can both search scEiaD for gene expression
57 across many dimensions (e.g. cluster, cell type, study) and project their own single cell data onto scEiaD for
58 comparison and rich automatic cell labelling. Second, the computational community can use this very large, well-
59 curated dataset to test algorithms for compute efficiency and performance in a diverse environment. As we believe
60 data re-use is a powerful and efficient approach to facilitate discovery, we provide our meta-atlas code-base, the
61 meta-atlas in several data formats, and propose general guidelines to optimally create custom meta-atlases.

62 **Results**

63 **We identify 33 ocular scRNA datasets across 3 species**

64 The first step in building a meta-atlas is identifying studies to draw the data from. We identified ocular
65 single cell RNA sequencing (scRNA) studies by querying PubMed, the Sequence Read Archive (SRA), and the
66 European Nucleotide Archive (ENA) for the inclusive terms “retina”, “single cell”, “scRNA”, “ocular”, “eye”,
67 “transcriptome.” We then hand filtered the results to only keep ocular and normal (non-perturbed or mutagenized)
68 data from single cell RNA-seq technology. On December 2020 we identified 33 deposited datasets that have been
69 published in 27 publications (Figure 1). To provide a non-ocular reference we also downloaded the raw sequence
70 data from the Tabula Muris project for re-processing. In cases where the fastq file from the SRA was not processed
71 properly, we acquired the bam files (SRA or personal correspondence) and re-extracted the fastq. After downloading
72 all the data we had 11 TB across 2470 fastq file sets.



73

74 *Figure 1: a. Schematic of the retina with major cell types delineated b. Simplified directed workflow of major steps*
 75 *in scEiaD creation from raw counts to gene counts, benchmarking optimal integration methods (SnakePOP) to*
 76 *produce batch corrected latent dimensions (Latent Dims), then downstream analysis outputs like clustering,*
 77 *differential gene testing (Diff Testing), and 2D UMAP visualization. c. Counts of published papers and batches*
 78 *(unique biological samples) for each scRNA technology, split by organism d. Cell type counts extracted from*
 79 *published studies for the more common retina cell types, split by species. Count of study accessions for each species*
 80 *overlaid on bar plot.*

81 Transcriptome quantification across multiple technologies

82 Droplet and well based scRNA-seq technologies require different quantification approaches as the former
 83 have UMI and multiple cells are quantified within a single file. We wrote a Snakemake based pipeline
 84 (SnakeQUANT, see methods) to quantify and merge both droplet and well based technologies into a single matrix
 85 for downstream processing. For well based data we perform both gene and isoform level quantification; for droplet
 86 base technologies we quantify both exonic and intronic gene-level expression to facilitate calculation of RNA
 87 velocity. In total we quantify 6.7e+10 molecules, finding 1.1e+10 unique molecules with a mean pseudoalignment
 88 rate of 66.5%. Across the 766,615 cells (post QC) we have an average of 3,410 RNA counts across 978 unique genes
 89 (see Supplemental File kallisto_stats.tsv and splicing_stats.tsv for more details).

90 1,204,269 cells before quality control

91 Gene-level counts were quantified with the kallisto bustools pseudo-aligner for both the droplet and well
 92 based samples. After empty droplet removal, we had 1,204,269 cells. We then removed cells which had more than

93 10% mitochondrial reads across all gene counts or fewer than 200 unique genes quantified. After these standard
94 quality control steps we were left with 790,072 cells (Supplemental Figure 1).

95 A core objective of many scRNA based studies is labeling the cell types. As this information is crucial to
96 assess dataset integration and provide an accurate reference for user querying, we extracted individual cell labels
97 with a combination of inspecting the GEO web site, supplemental information from the publication, web resources
98 (e.g. a web app was created for the paper), and personal correspondence. After normalizing cell type name
99 nomenclature, we obtained labels for 375,966 cells across 33 cell types (Supplemental Table 1).

100 **Running 11 tools in a Snakemake-based system**

101 Disentangling the technical and biological effects when integrating multiple datasets is crucial. We define
102 batch as each unique biological sample and assume each study is at least one unique sample. We studied the
103 metadata and methods of each study to identify the unique biological samples. Within the current scEiaD data set we
104 identified 86 batches across 34 deposited datasets in 26 published papers.

105 A wide variety of methods have been written for scRNA-seq integration. As we were uncertain which
106 would perform the best, we ran 11 tools with a commonly used set of key parameters like number of hyper variable
107 genes (HVG), size of latent dimensional space, and the number of nearest neighbors for the louvain clustering
108 algorithm. The Snakemake system was used to automate the running of the wide variety of tools. In total 5,591 jobs
109 were run to quantify gene expression and build the unified Seurat objects (SnakeQUANT) and 2,446 jobs were run
110 to assess integration performance (SnakePOP).

111 The two key metrics which have to be balanced in order to optimize integration performance are cell type
112 or cluster purity (where different cell types or clusters should be homogenous) and batch mixing (the same cell types
113 should be similar across independent studies). While these can be visually assessed by looking at marker gene
114 expression across the 2D UMAP projection, it is more rigorous and scalable to quantify these diametrically opposed
115 characteristics.

116 **scPOP wraps several different methods for measuring integration performance.**

117 Multiple methods have been proposed to quantitatively evaluate batch correction. Some of these metrics
118 evaluate the concordance between sets of labels, while other compute distances between the individual data points of

119 a given set of labels. While any one of these methods can be useful, we propose that calculating and evaluating them
120 in tandem provides greater accuracy for dataset integration. We developed the R package scPOP, a lightweight, low
121 dependency R package which brings together the Local Simpson Index (LISI), Average Silhouette Width (ASW),
122 Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI) metrics from the R packages Harmony,
123 kBET and aricode, respectively. The LISI and ASW were used to measure batch mixing (where lower is better), cell
124 type mixing (higher is better), and cluster mixing (higher is better). NMI and ARI were used to assess the
125 consistency of cell type to cluster assignment (where 1 is perfect correspondence between cluster and cell type).

126 To visualize the interplay between batch mixing and cell type distinction we plot the batch mixing LISI
127 score (which has been multiplied by -1) on the y-axis (higher is better) against the cluster LISI on the x-axis (higher
128 is better). The best performer on both metrics will be in the top right corner (Supplemental Figure 2a). In the same
129 manner we plot the silhouette metric (Supplemental Figure 2 b). To merge the different scores we define
130 $sumZScale = \sum scale(m)$ where m is a metric (LISI by batch, LISI by cluster, LISI by cell type, silhouette by
131 batch, silhouette by cluster, silhouette by celltype, NMI, and ARI).

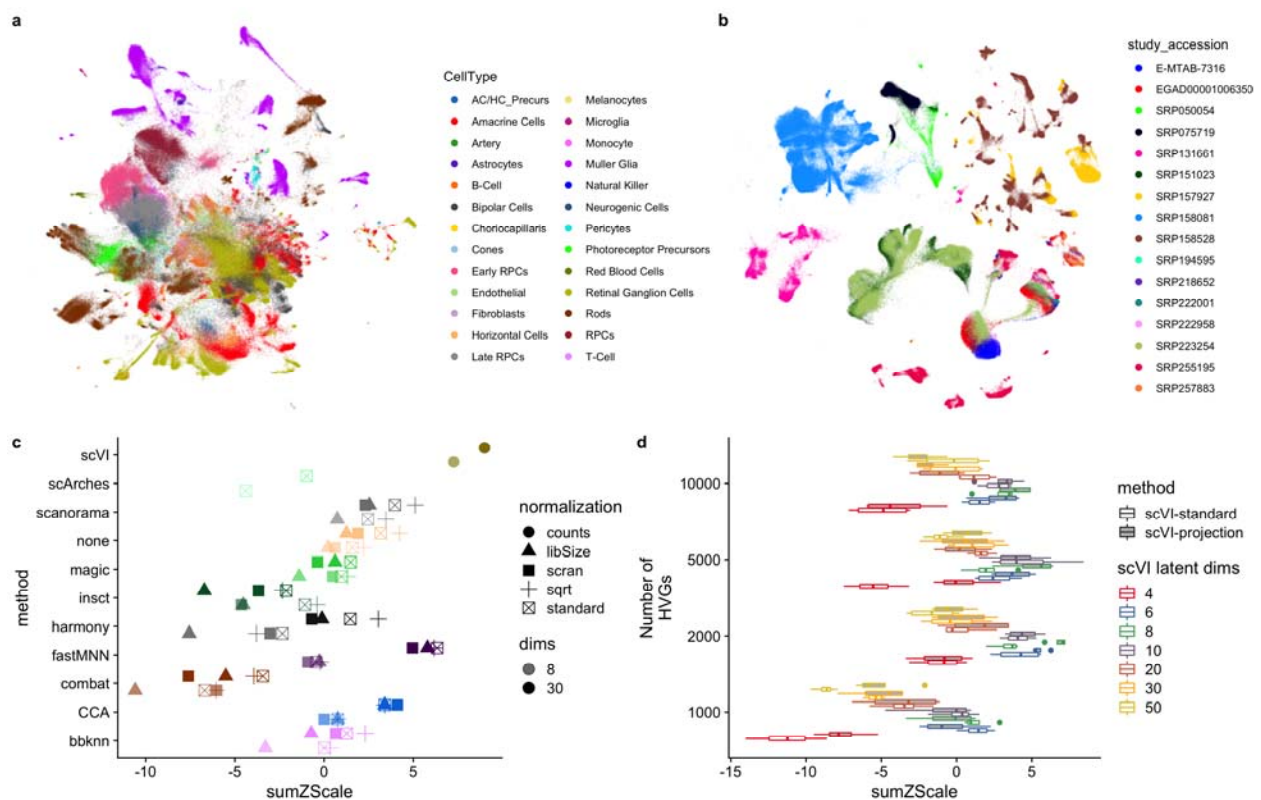
132 On one extreme we have ComBat, which merges together different batches very well, but also mixes
133 together the distinct cell types (Figure 2a). The other extreme is not using any batch integration method, where you
134 see very distinct groups of cells, but also each nearly study is has a distinct region in the UMAP (Figure 2b). In our
135 scEiaD dataset we see that like ComBat, Harmony and CCA are weighed more towards batch mixing then cluster
136 and cell type purity. Scanorama, fastMNN, and bbknn prioritize cleanly separating the clusters. With our scEiaD
137 meta-atlas insct, trVAE, and magic do not perform particularly well in batch mixing or cluster purity. Across the 11
138 integration methods we tested, scVI achieved the highest sumZScale for out scEiaD meta-atlas.

139 **Different normalization methods alter integration performance**

140 There are several normalization approaches that have been used or published. The “standard” approach that
141 the popular analysis packages Seurat and anndata use by default is to, per cell, divide the counts by the sum counts
142 for the cell, multiply by a scaling factor, then log transform. This helps make the count distribution more normal,
143 which is an assumption that many algorithms require. In contrast, the scran normalization method groups cells into
144 pools and normalizes across the pool summed counts instead of the individual cell counts. We also use the square
145 root (sqrt) normalization which replaces the log transformation with a square root. Library size (libSize)

146 normalization omits the sqrt or log transformation. Finally some methods, like scVI, directly use the raw counts data
 147 for modeling the data.

148 As expected the libSize normalization which omits the log or square root scaling generally performs the
 149 worst (Figure 2c). We see that the remaining normalization techniques alter the batch correction performance,
 150 though the outcome differs across the different methods. We also see that changing the number of latent dimensions
 151 (8 or 30) can occasionally dramatically change performance. These results demonstrate the importance of assessing
 152 performance in a rigorous manner across many parameters.



153
 154 *Figure 2: a. Example of a method (combat) which has a high level of batch blending, but poor separation of cell*
 155 *types (colored by cell type). b. no batch correction cleanly separates cell types but does not mix batches (colored by*
 156 *study). c. sumZScale (higher is better) for each method across a variety of data normalizations. All methods shown*
 157 *here use 2000 HVG, lowain clustering, and 8 latent dimensions. Each color is a different method. d. Boxplot of 4*
 158 *different clustering resolutions for across 1000 to 10000 HVG numbers and 4 to 50 scVI latent dimensions. Open*
 159 *boxes are using scVI-standard and gray boxes are scVI-projection (human reference with the remaining data*
 160 *projected)*

161 Further optimization of scVI with grid search and projection

162 To find the best set of parameters for scVI in our dataset we did a grid search across key parameters: HVG,
 163 latent dimensions, and k nearest neighbors. Furthermore we used a recent advance in scVI capability (\geq version

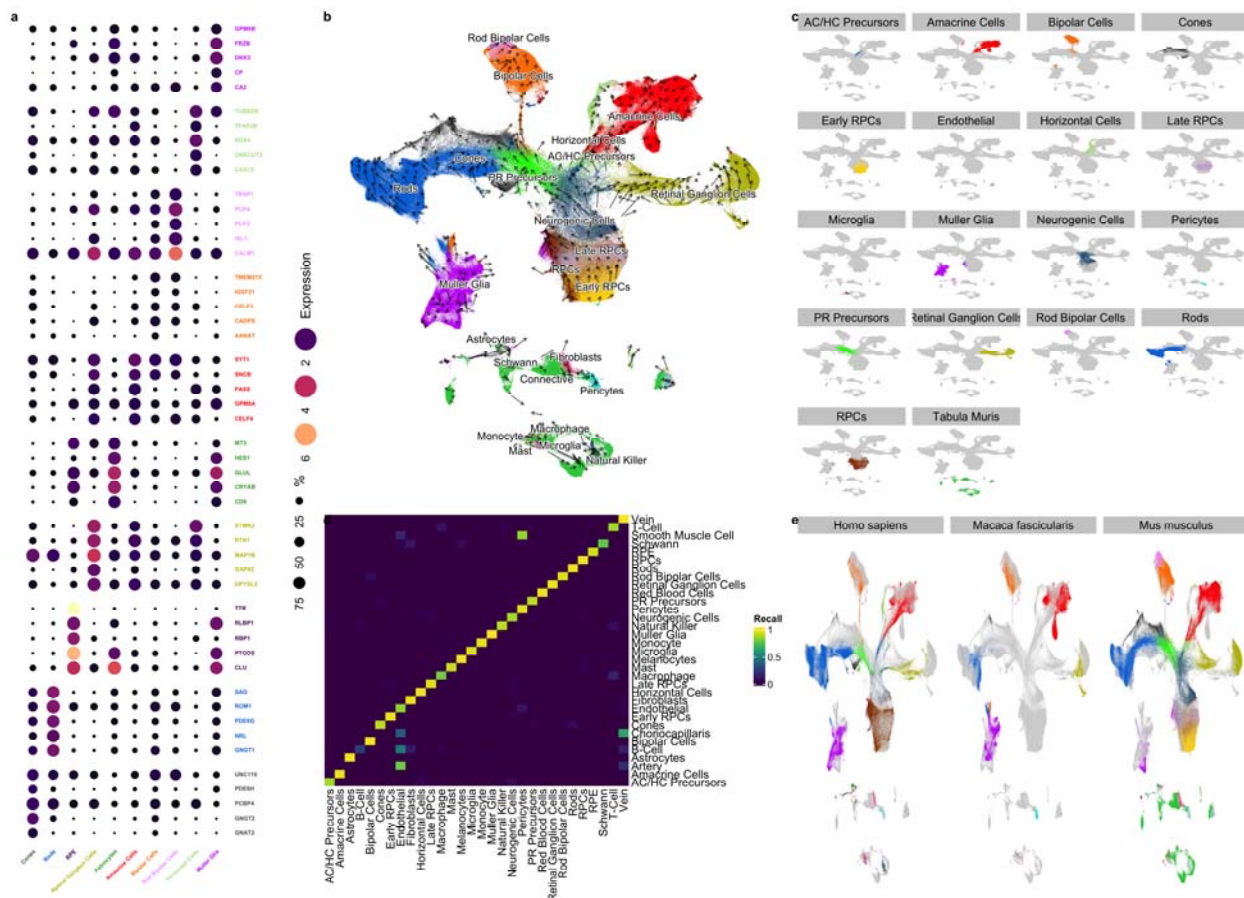
164 0.8.0) adapted from scArches that allows one to build a reference model and query or project the new cells onto it.⁴⁰
165 We refer to this projectable model as “scVI-projection”, and the previous scVI model as “scVI-standard”. We built a
166 scVI-projected model trained on human cells and then projected the mouse and macaque data onto it. We then
167 compared scVI-projection against the against the previous scVI-standard across all the previously mentioned
168 parameters. Using scPOP we first saw that the scVIprojection approach generally performed better than running
169 scVI with all of the data (Figure 2d). We found the optimal overall parameters to be 5000 HVG, 8 latent dimensions,
170 and with 5 k-nearest neighbors for the cluster finding. We also varied the UMAP projection values of nearest
171 neighbors and minimum distance to qualitatively pick a 2D projection, selecting a minimum distance of 0.1 and 50
172 nearest neighbors.

173 **High accuracy xgboost ML model built to label unknown cell types across all technologies**

174 To further study cell type specific expression patterning we needed to label the 414,106 unlabeled cells.
175 Traditionally this is done by clustering the cells, then using cell type specific markers to label the clusters. However,
176 as we had hundreds of thousands of expert labeled cells across 17 publications (Figure 1a) we built a xgboost-based
177 machine learning model that used 2/3 of the labeled cells as a training set (see methods for more details) to train a
178 cell type predictor for scEiaD input. The trained model was used to predict the cell type assignments for all cells in
179 scEiaD. In this manner we both label most cells (cells which cannot assign a cell type with a probability above 0.5
180 were left unlabeled) and correct a small number of probable mislabels in the truth set (Figure 3a).

181 As a brief case study of our ML performance, we look at the cell type labels assigned to a Shekhar et
182 al. study where they used SMART-seq on a retinal bipolar cell enriched population of cells²¹. Even though our ML
183 algorithm was trained only on droplet-based data, our algorithm labels most of this dataset as retinal bipolar cells,
184 with the next most common cell types being amacrine. The same result was found by Shekhkar et al. (Supplemental
185 Figure 3). To more generally evaluate performance we use the precision recall (PR) curve, which visualizes the
186 ability of the model to precisely label known cells at a given confidence. The area under the PR curve (AUC)
187 summarizes the effectiveness of the model across different cell types, with 1 being the highest performance. The
188 xgboost model can predict rods, bipolar cells, and Müller glia with near perfect performance (Supplemental Figure
189 4). Most of the remaining cell types can be predicted with an AUC over 0.9 Supplemental Table 2. Several of the
190 precursor cell types (photoreceptor, AC/HC, and neurogenic) were labeled with lower confidence Supplemental

191 Table 3. The next most common labels for these cell types were either other precursor cells (e.g. 100 AC/HC
 192 Precursors were labeled as Neurogenic) or the adjacent terminal cell type (e.g. 254 photoreceptor precursors were
 193 labeled as cones). Other cell types that were challenging for the model to predict were the artery, choriocapillaris,
 194 and vein. Artery, choriocapillaris, and vein are constructed from endothelial cells and we find that for all three of
 195 these, endothelial was the second most common label. Overall, our xgboost based ML model shows strong accuracy
 196 across the major cell types of the retina (overall AUC of 0.98).



197
 198 *Figure 3: a. Top genes that are differentially expressed across the major cell types of the retina (PR is short for*
 199 *Photoreceptor). Genes are colored by which cell type they are differentially expressed in. The dot size is*
 200 *proportional to the percentage of those cells that have detectable levels of the gene. The color of the dot is the log2*
 201 *scaled CPM expression. b. 2D UMAP projection of scEiaD, colored by cell type (Tabula Muris data is gray). Arrows*
 202 *are scvelo RNA velocity. Longer arrows are cells with higher velocity (relatively more unsplliced transcripts).*
 203 *c. Facet plot that demonstrates how each major cell type of the retina is contained within a distinct space.*
 204 *d. Confusion matrix of cell type prediction performance of our xgboost labeller between predicted (x axis) and*
 205 *known (y axis) using data withheld from the machine learner. Most of the cell types are indeed labeled as their true*
 206 *type. e. Faceting of 2D UMAP by species and colored by cell type demonstrate how the major cell types of the retina*
 207 *share space with like cell types, despite being from mouse, human, and macaque.*

208 **ML cell type labels result in high study diversity for each cell type**

209 After ML projection of cell type labels from the original 375966 labels onto a total of 758278 cells we have
210 substantially improved the number of studies per cell type. For example, we went from 8 human studies with labeled
211 Müller Glia to 14 after labeling (Supplemental Table 4). Overall we go from an average of 4 studies per human cell
212 type to 8 and 2 studies per mouse cell type to 8 after transferring the cell type labels.

213 As another check on the quality of the cell type assignments, we ran the cell and cluster independent
214 haystack gene search and pairwise differential expression tests between the predicted cell types (see methods for
215 further details). We show the five most differentially expressed genes for each of the major retina cell types are
216 consistent with known retinal cell markers (Figure 3a). As a simple metric to identify known and unknown genes
217 relating to the cell type specific expression we search PubMed for the number of publications with two searches per
218 gene. We expect most of the genes identified to be known in the literature. The first search is the more precise “gene
219 AND cell type” (e.g. “PDE6H AND Cones”) and the second search is the more inclusive “gene AND retina”
220 (e.g. “PDE6H AND Retina”). Of the 50 genes in (Figure 3a), 37 had one or more citations in the gene - cell type
221 search (Supplemental File celltype_markers.tsv) and 45 had one or more citation in the inclusive search. The 50
222 genes had a mean of 46 studies (with the inclusive gene by “retina” search). In contrast, 100 randomly chosen genes
223 had a mean of 2 (wilcox test $p < 1.44 \times 10^{-17}$).

224 **The scVI-based scEiaD UMAP projection blends batches and species while separating cell** 225 **types**

226 The 2D UMAP projection of the scVI-calculated batch corrected 8 latent dimensional space blends the 33
227 studies together while also maintaining distinct space for the 31 unique cell types. We also see good mixing across
228 all the droplet and well based single cell technologies (Supplemental Figure 5). We see the neurogenic and
229 progenitor populations from which the retinal cell types are derived near the center of the UMAP visualization. The
230 photoreceptor precursors are adjacent to the neurogenic population and, as demonstrated by the RNA velocity
231 dynamics, flow into the rods and cones. The amacrine and horizontal precursors (AC/HC) likewise flow from the
232 neurogenic center into the mature amacrine and horizontal cells.

233 The photoreceptors (cones and rods) of the retina which are responsible for color and low-light vision,
234 respectively, are near each other in the UMAP space. The major remaining retina cell types, by proportion in the

235 mammalian eye are the Müller Glia, which are a glial cell type which help support the neurons of the retina. Next we
236 have the neural cell types which transmit and help interpret the signals from the photoreceptors before they leave the
237 retina via the optic stalk: the amacrine cells, retinal ganglia, horizontal, and bipolar cells. All of these cells are in
238 well separated spaces in the UMAP. Finally we see across species that the major cell types overlap each other
239 (Figure 3d). The macaque retina cells are not present in the precursor/neurogenic center of the UMAP as expected
240 because only fully developed tissues were sampled in these data.

241 While by eye the UMAP 2D projection generally blends together the three different species (macaque,
242 mouse, human) in the UMAP visualization, we more rigorously tested this by changing the inputs to our xgboost cell
243 type predictor machine learning system. We trained two new models: one with only human data and one with only
244 mouse data. We then applied this model to the other two species to see whether the model trained on one species was
245 generally transferable to another species. Again, both models (human only and mouse only) both proved to be highly
246 accurate at predicting cell types in the other species with the terminal cell types (Supplemental Figure 6).

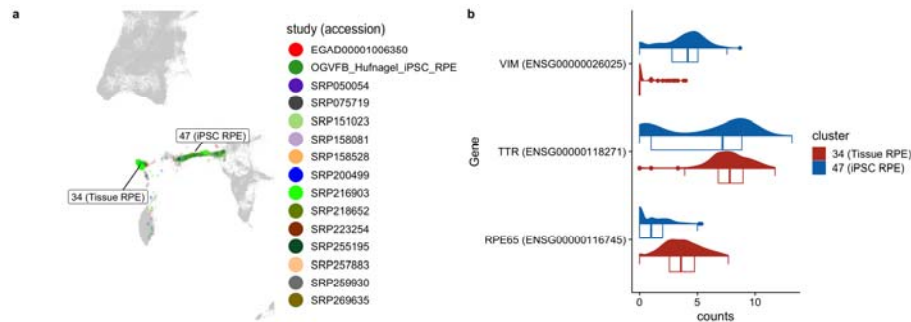
247 **Projection of outside data onto scEiaD demonstrates similarities and differences of iPSC** 248 **and organoids to primary cells**

249 The hard work of creating this resource can be leveraged and extended by the wider community with a few
250 relatively simple steps and modest compute requirements. Very briefly, if outside groups quantify their mouse or
251 human scRNA with kallisto (bustools) and the same references (see methods), they can overlay their data on top of
252 scEiaD by 1. installing scVI (version 0.9.0 or higher), 2. downloading our 13 megabyte scVI model, and 3.
253 following the Jupyter notebook on Google colab that we provide as a live demo available from
254 <https://github.com/davemcg/scEiaD>. We demonstrate the power of this approach in two ways.

255 First, a wild-type retinal organoid dataset from Kallman et al was projected onto scEiaD.⁴⁴ In the Google
256 colab notebook we show how a subset of the reads from Kallman et al.'s SRR12130660 can be processed from the
257 raw reads to a UMAP visualization in under 10 minutes. Indeed we see how organoid-based retinal cell types can be
258 detected overlapping primary cells for RPCs, photoreceptors, amacrine cells, and retinal ganglion cells
259 (Supplemental Figure 7).

260 Finally, we demonstrate how the RPE derived from the Bharti iPSC differentiation process express
261 canonical RPE markers of TTR and RPE65 but end up in a slightly different position in the UMAP projectionFigure

262 4a). Differential expression analysis (see methods) identified that vimentin is nearly exclusively expressed in the
263 iPSC-based RPE Figure 4b). This result is not surprising as Hunt et al. demonstrated that cultured and proliferating
264 RPE express high levels of vimentin⁴⁵ and this same trend is seen in bulk RNA-seq from stem cell RPE and primary
265 RPE Supplemental Figure 8).⁴⁶



266

267 *Figure 4: a. RPE distribution colored by study demonstrates how the most RPE are in two locations. iPSC-based*
268 *RPE we provided are located more enriched in cluster 47. Tissue RPE more enriched in cluster 34 b. Violin plot of*
269 *two functional RPE markers (TTR, RPE65) and vimentin (proliferating RPE marker)*

270 Methods

271 Reproducibility

272 The set of Snakemake pipelines that takes in the raw fastq sequence and outputs the scEiaD is at
273 <https://github.com/davemcg/scEiad>.⁴⁷ The publication commit is #9d3c5e1 Furthermore, the repository has been
274 deposited at Zenodo under 10.5281/zenodo.4638117 We will briefly discuss the pipeline choices, programs and
275 algorithms, and versions below. For the R packages, we provide package versions as the supplementary file
276 “R_session_info.txt”

277 Quantification of gene counts

278 Gene quantification is handled by the SnakeQUANT snakefile. First, we generated multiple quantification
279 indices to facilitate calculation of RNA velocity. For each droplet technology, we generated a separate set of
280 transcript sequences that contain both exonic and intronic sequences using the “get_velocity_files” function from the
281 R package BUSpaRse. The reference transcriptome annotation used to build sets of transcript sequences were the
282 Gencode “gencode.vM25.annotation.gtf.gz” and “gencode.v35.annotation.gtf.gz” for mouse and human,
283 respectively.^{48,49} Because the Macaca Fascicularis genome is less well annotated we used the Ensembl release 101

284 genome and transcriptome annotation “Macaca_mulatta.Mmul_101.gtf.gz.”⁵⁰ A single set of exonic transcript
285 sequences was created for each species for all well technologies. A kallisto quantification index was generated for
286 each of these fastas using kallisto index (0.46.2).⁵¹ Well based samples were quantified using kallisto quant. For the
287 relatively few single ended samples, the params “-single -l 200 -s 30” were used for kallisto quant. Otherwise the “-
288 bias” flag was added. For the droplet-based samples, we adapted the bustools workflow for generating spliced and
289 unspliced count matrices. (0.39.4).⁵²

290 **Intersection of gene names between mouse, macaque, and human**

291 To facilitate comparison of gene expression across species, where possible we converted mouse and
292 macaque gene ids and names to human ones. We downloaded a mapping of orthologous genes between human,
293 mouse, and macaque using the Ensembl BioMart web browser in November 2020. We identified 15,759 human
294 genes that could be directly mapped to mouse and macaque orthologs. Genes present in mouse or macaque that were
295 not found in human were not used for HVG gene selection, but were retained and used for differential gene
296 expression.

297 **Custom macaque reference quantification**

298 As we noticed that several retina marker genes (e.g. NRL and CRX) had very low expression in the
299 macaque data we quantified the scRNA data twice: once with the Ensembl reference and again with the same
300 Gencode human reference used for the human data. We compared the gene-level counts for each cell and replaced
301 the macaque gene count with the human counts if the human counts were greater than the macaque counts and, to
302 prevent genes with very few total counts from being used, we required the counts greater than the first quartile of
303 non-zero macaque gene expression.

304 **Remove empty droplets and further QC.**

305 After bustools count, we used R (3.6.2) to remove empty droplets. The BUSpaRse package was used to
306 input the bustools counts mtx file. The DropletUtils package with the “barcodeRanks” function was used to
307 automatically detect the inflection point in the barcode count ranks that delineates the likely empty droplets.⁵³ We
308 then removed cells with percent mitochondrial reads of >10%. After merging the individual count matrices into one
309 sparse matrix, we created a Seurat version 3 object and removed cells with fewer than 200 detected unique genes,
310 and for the droplet data, more than 3000 detected genes (these are likely to be doublets).³²

311 **Normalization and batch effect correction**

312 The following steps (normalization through benchmarking are handled by the SnakePOP pipeline) We
313 tested several gene count normalization approaches as we were not certain which would produce an optimal
314 outcome: standard (default Seurat, library size normalization, then log transform), sqrt (same, but with sqrt
315 normalization), libSize (omit the log or sqrt normalization), scran, SCT from Seurat, and for scVI, no normalization
316 (counts).⁵⁴ Our R implementation of the normalization approaches as well as how we constructed the Seurat v3
317 object can be found in the supplementary file “make_seurat_obj_functions.R”

318 **Batch normalization under a grid search procedure**

319 We tested scArches, bbknn, insct, magic, scVI, CCA, scanorama, harmony, fastMNN, combat, none against
320 2000 HVGs, the different gene count normalization procedures discussed above, and both 8 and 30 outputted batch
321 corrected latent dimensions. The latent dimensions are the input for clustering, the 2D UMAP visualization, and the
322 xgboost machine learning to transfer cell type labels to unlabeled cells. We were unable to run every method
323 successfully with every normalization method. Magic could not complete with the standard or libSize normalization.
324 CCA could not complete with the libSize normalization. We also tried the DESC, liger, and Conos batch corrections
325 methods but were unable to get them to work reliably so they were dropped. The batch correction step
326 implementation can be found the supplementary file “merge_methods.R” and in the github repo
327 (https://github.com/davemcg/scEiaD/blob/master/src/merge_methods.R).

328 **Clustering and UMAP**

329 Louvain-Jaccard clustering against the batch corrected latent dimensions used the Seurat implementation.⁵⁵
330 We tried the k-nearest neighbors (knn) parameters 5 and 7 (where 5 gives more clusters than 7). We also used the
331 leiden algorithm as implemented by PARC with a resolution of 0.6 and 0.8 (higher results in more clusters).^{56,57}
332 These two resolutions were chosen as they roughly gave the same number of clusters at the Seurat Louvain-Jaccard
333 approach with a knn 7.

334 The UMAP visualization was calculated with the Seurat “RunUMAP” using the uwot R package.⁵⁸ We
335 tried min.dist parameters of 0.001, 0.1, and 0.3 and tried n.neighbors across 15, 30, 100, and 500. A smaller min.dist
336 value gives “tighter” groupings while a higher number of n.neighbors uses a larger number of near cells to calculate
337 the global positioning.

338 **Benchmarking and scPOP**

339 We wrote the scPOP R package to unify the LISI and Silhouette metrics from Harmony and kBet,
340 respectively, along with NMI and ARI.^{36,59} LISI and Silhouette require a dense matrix, which is a problem for our
341 data as a 766615 cell by 8 latent dimension dense matrix cannot fit in our largest available memory node (1.5 TB).
342 We down-sampled the dataset to ~100,000 cells, taking care to keep all rarer cell types for the LISI and Silhouette
343 benchmarking.

344 To merge these metrics into a balanced single score, we Z scale each and sum them. scPOP produces both
345 tables and visualizations allowing the user to quickly see both the interplay of batch mixing and cluster/cell type
346 separation and the overall performance. If a user wishes to prioritize batch mixing or cluster/cell type separation we
347 let the user provide a custom batch/cluster-cell type scaling value (1 is the default).

348 **Multi-step doublet removal**

349 To identify probable doublets (more than one cell in a droplet) we ran DoubletDetect and scrublet and
350 calculated the distribution of DoubletDetect and scrublet scores across all clusters and removed clusters with a score
351 in both metrics greater than 4 standard deviations above the mean.^{60,61} This removed another 23,457 cells, leaving
352 766,615 in total.

353 **xgboost based cell type model**

354 In order to identify cell types for the 361,456 unlabeled cells we designed a custom xgboost based cell type
355 classifier. We took labeled data and split it into training (2/3) and test (1/3) sets, stratified by cell type. The input
356 features used to train the model are the scVI latent dimensions, the total number of reads in each cell, the number
357 genes detected in each cell, and the percent mitochondrial gene expression of each cell. We additionally generated
358 features using the age of each sample by group sample into three developmental categories (Early Development,
359 Late Development, and Adult) and then generated a one-hot encoded feature for each category. In order to speed up
360 training times, we used the gpu implementation of the xgboost algorithm from the the xgboost python library. The
361 model was trained using default parameters. The trained model had an overall macro and micro AUC score of 0.98
362 and 0.99, respectively. This model was then used to identify labels for all cells. Unlabeled data was pre-processed
363 identically to training data and fed into model to generate a vector of label probabilities for each cell. We selected
364 the highest label probability for each cell, and required a minimum probability of 0.5 to assign a label to a cell.

365 For the organism specific xgboost ML we followed the above procedure, except that we combined Early
366 RPCs, Late RPCs, and RPCs into one category and did not attempt to predict the non-retina cell types
367 (e.g. fibroblasts) as there were very few labeled cells across all three organisms.

368 **Marker gene identification**

369 To identify marker genes across the CellType (predict) and cluster groups, we used the scran findmarkers
370 (wilcox test) along with the singleCellHaystack algorithm. The scran findmarkers test runs a wilcox test in a
371 pairwise manner (e.g. Rods vs all other cell types). It returns an overall p-value (and FDR) that assesses how well
372 the gene is at separating the group of interest from all other cells. It also returns for each pair-wise comparison an
373 area under the curve (AUC) score, where 1 is a perfect power to distinguish and 0 is no power. The
374 singleCellHaystack algorithm uses a Kullback-Leibler divergence (D KL) measurement of the scVI lower
375 dimensional space to identify genes with non-random distribution. A higher D KL score represents a gene with
376 “specific” expression in the lower dimensional space and is used to calculate a FDR corrected p value against the
377 full distribution of D KL values. We filtered to keep genes with scran FDR < 1, a mean AUC > 0.2, and a log10(D
378 KL FDR) < -10000. No more than 50 genes for each cell type were retained (sorted by mean AUC).

379 **Calculation of RNA velocity**

380 RNA velocity calculations were with the velocraptor wrapping of the scVelo python library.⁶² From the
381 anndata objects generated by our Snakemake pipeline we calculated velocity across all genes. Genes without
382 detectable velocity were dropped. The scVI generated latent space (instead of PCA) was used to calculate first and
383 second order moments. The calculated moments were used to estimate RNA velocity. Differential velocity was
384 tested between celltypes using pairwise wilcox rank sum tests.

385 **Conclusion**

386 **Limitations**

387 The scVI model is first built the on human data. The mouse and macaque data are then projected (or
388 queried) onto it with the scVI implementation of the scArches method. While this system works very well to
389 integrate information between these three species, this approach may not scale to more distantly related species.

390 Another limitation is that discrepancies between cell type labels between different labs makes certain transitioning
391 cell type labels a bit imprecise. One example is how the rods and photoreceptor precursor labels partially overlap.
392 Though we attempted to ameliorate the issue by removing cell type labels in large disagreement with the consensus,
393 some disagreements could propagate into our machine labeled cells type assignments. These issues may reflect
394 labeling continuous processes with discrete labels.

395 While scEiaD distinguishes the major cell types very well, some of the cell types contain many “sub types”
396 - notably the amacrine cells have a huge variety in morphology with a recent study identifying over sixty different
397 types of amacrine cells in mouse.²⁹ At this time our batch corrected pan retina cell space does not precisely resolve
398 these sub cell types with high resolution. We are actively working to “sub cluster” the cell types so we can robustly
399 and reliably identify the high diversity of retinal cell types across the entire retina.

400 **The scEiaD is a unique ocular resource that provides a highly diverse, large N dataset with** 401 **a relatively small amount of compute power**

402 We have assembled the largest ocular single cell transcriptome database to date. The rapid of advancement
403 of algorithms to batch correct and process data continue to reduce the computational requirements to handle huge
404 numbers of cells. The scVI batch correction step on around one million cells runs within a few hours on a GPU and
405 150GB of memory. This places this crucial step within the capabilities of a moderately powerful computer or a cloud
406 compute node. Further downstream processing can largely be done a computers with 64+ Gb of memory and a few
407 hundred GB of disk space. We believe that our efforts can be replicated in any other tissue / system with a large
408 number of independent studies by a small number of computational scientists following our general approach. We
409 provide the completed analysis as both Scanpy (h5ad) and Seurat objects at <https://github.com/davemcg/scEiaD>.

410 **Benchmarking and quantitation of integration performance is crucial for meta-atlas** 411 **studies**

412 We originally intended to use the Seurat CCA method to integrate the datasets. However, the long run-times
413 of CCA and poor integration of our known cell types led us to benchmark more methods and parameters. After
414 adding more integration methods we first attempted “hand-assess” the integration results by using the UMAP 2D
415 projection view. This proved to scale poorly and this led use to curate some of the more useful benchmarking
416 algorithms (NMI, ARI, silhouette, and LISI) that roughly matched our “hand-assessed” results into the scPOP R

417 package. While we chose the scVI algorithm, we strongly suggest any other groups attempting a similar meta-atlas
418 construction chose a quantifiable set of criteria so optimal methods and parameters can be picked.

419 **Transfer of cell type labels from a smaller number of studies onto the remaining cells is a** 420 **powerful way to increase diversity in a meta-atlas**

421 We first hand curated over 350,000 published cell type labels across 24 publications. With a xgboost
422 algorithm using the latent dimensions, cell age classification (developing or matured), and the UMAP coordinates,
423 we can very accurately label the remaining cells. Many other cell type labeling algorithms and systems exist for
424 those groups less willing or able to tune a machine learning algorithm. For example, the developers of scVI also
425 have a cell type label projection algorithm called scANVI. Whichever approach you use, taking a smaller number of
426 high quality labels and projecting them onto the remaining cells is a powerful way to leverage community
427 knowledge across a huge diverse dataset.

428 **Projection allows community knowledge to be leveraged by all**

429 Many retina atlases have been published to date. We argue that we have created the first atlas that is
430 generally useful because 1. our dataset/atlas is several times larger than any other published set, 2. our data is
431 available via download in several forms at <https://github.com/davemcg/scEiaD>, and crucially 3. we provide a
432 Google colab/Jupyter notebook which exactly lays out how to use scVI to project (or query) outside data onto our
433 scEiaD with minimum compute resources. We demonstrate concretely how this can work by showing how iPSC
434 RPE can be queried onto the reference dataset to demonstrate both similarities and dissimilarities in their
435 transcriptomes.

436 **Supplemental Information**

437 **Tables**

CellType	HS Published	MF Published	MM Published	HS Transferred	MF Transferred	MM Transferred
AC/HC Precursors	1,450	0	0	2,779	2	291
Amacrine Cells	13,430	19,246	35,837	21,179	36,308	53,958
Artery	167	0	0	0	0	0
Astrocytes	1,102	0	39	1,699	64	326
B-Cell	509	0	8,109	336	27	76
Bipolar Cells	5,012	7,617	15,536	9,855	14,342	32,573
Choriocapillaris	225	0	0	0	0	0
Cones	3,242	694	4,705	6,844	1,107	18,509
Endothelial	420	295	3,678	1,048	370	3,926
Fibroblasts	1,632	0	221	2,037	16	3,477

CellType	HS Published	MF Published	MM Published	HS Transferred	MF Transferred	MM Transferred
Horizontal Cells	5,282	163	1,719	9,033	259	8,501
Macrophage	559	0	0	586	0	91
Mast	101	0	0	216	0	177
Melanocytes	259	0	0	323	1	3
Microglia	425	154	25	829	206	978
Monocyte	311	0	25	369	0	454
Muller Glia	19,303	13,763	4,368	26,508	25,025	9,325
Natural Killer	168	0	0	221	0	4
Neurogenic Cells	2,166	0	8,314	5,255	12	28,811
Pericytes	457	2,459	20	583	3,703	1,228
PR Precursors	2,009	0	5,122	6,540	1	16,107
Red Blood Cells	63	0	1,737	706	6	1,842
Retinal Ganglion Cells	6,663	9,762	27,887	10,308	12,495	41,112
Rod Bipolar Cells	392	0	9,837	600	0	11,838
Rods	31,292	1,481	22,887	69,339	14,004	58,129
RPCs	17,701	0	0	35,696	191	5,403
RPE	369	0	0	573	27	1,501
Schwann	288	0	0	501	5	284
Smooth Muscle Cell	45	0	0	0	0	0
T-Cell	1,075	0	4,335	1,238	24	49
Unlabelled	101,950	53,641	235,058	1,119	1,074	6,144
Vein	490	0	0	1,021	0	37
Early RPCs	0	0	27,962	898	6	35,080
Late RPCs	0	0	21,362	318	0	34,444

438 *Supplemental Table 1: Counts for cell type labels. Published are the author created labels from the*
 439 *published datasets. Transferred are the cell labels that were transferred by our xgboost-based machine learning*
 440 *model onto the entire scEiaD dataset.*

AUC	Cell Type	Study
0.76	AC/HC_Precurs	All
0.71	AC/HC_Precurs	SRP151023
0.83	AC/HC_Precurs	SRP223254
0.99	Amacrine Cells	All
0.20	Amacrine Cells	E-MTAB-7316
1.00	Amacrine Cells	EGAD00001006350
0.97	Amacrine Cells	SRP050054
0.98	Amacrine Cells	SRP075719
0.98	Amacrine Cells	SRP151023
0.93	Amacrine Cells	SRP158081
1.00	Amacrine Cells	SRP158528
1.00	Amacrine Cells	SRP194595
0.72	Amacrine Cells	SRP222001
1.00	Amacrine Cells	SRP222958
0.98	Amacrine Cells	SRP223254
1.00	Amacrine Cells	SRP255195
0.97	Astrocytes	All
1.00	Astrocytes	E-MTAB-7316
0.97	Astrocytes	EGAD00001006350
0.94	Astrocytes	SRP050054
1.00	Astrocytes	SRP255195
0.97	B Cell	All
0.81	B Cell	EGAD00001006350
1.00	B Cell	SRP131661
0.41	B Cell	SRP218652
0.91	B Cell	SRP257883
0.98	Bipolar Cells	All
1.00	Bipolar Cells	E-MTAB-7316
1.00	Bipolar Cells	EGAD00001006350
0.97	Bipolar Cells	SRP050054
0.99	Bipolar Cells	SRP075719

AUC	Cell Type	Study
0.97	Bipolar Cells	SRP151023
0.94	Bipolar Cells	SRP158081
1.00	Bipolar Cells	SRP158528
1.00	Bipolar Cells	SRP194595
1.00	Bipolar Cells	SRP222001
1.00	Bipolar Cells	SRP222958
0.98	Bipolar Cells	SRP223254
1.00	Bipolar Cells	SRP255195
0.94	Cones	All
1.00	Cones	E-MTAB-7316
1.00	Cones	EGAD00001006350
0.96	Cones	SRP050054
0.71	Cones	SRP151023
0.88	Cones	SRP158081
0.97	Cones	SRP158528
1.00	Cones	SRP194595
1.00	Cones	SRP222001
0.47	Cones	SRP222958
0.97	Cones	SRP223254
0.90	Cones	SRP255195
0.98	Early RPCs	All
0.99	Early RPCs	SRP158081
0.97	Endothelial	All
0.98	Endothelial	EGAD00001006350
0.99	Endothelial	SRP050054
1.00	Endothelial	SRP131661
0.94	Endothelial	SRP158528
1.00	Endothelial	SRP194595
0.00	Endothelial	SRP218652
1.00	Endothelial	SRP222958
0.59	Endothelial	SRP255195
0.97	Fibroblasts	All
1.00	Fibroblasts	EGAD00001006350
0.95	Fibroblasts	SRP050054
0.99	Fibroblasts	SRP131661
0.53	Fibroblasts	SRP218652
0.99	Fibroblasts	SRP257883
0.96	Horizontal Cells	All
1.00	Horizontal Cells	E-MTAB-7316
1.00	Horizontal Cells	EGAD00001006350
0.97	Horizontal Cells	SRP050054
0.99	Horizontal Cells	SRP151023
0.79	Horizontal Cells	SRP158081
0.92	Horizontal Cells	SRP158528
1.00	Horizontal Cells	SRP222001
1.00	Horizontal Cells	SRP222958
0.98	Horizontal Cells	SRP223254
1.00	Horizontal Cells	SRP255195
0.97	Late RPCs	All
0.97	Late RPCs	SRP158081
0.70	Macrophage	All
0.58	Macrophage	SRP218652
0.97	Macrophage	SRP257883
0.83	Mast	All
0.92	Mast	EGAD00001006350
0.96	Mast	SRP218652
0.97	Melanocytes	All
1.00	Melanocytes	EGAD00001006350
0.00	Melanocytes	SRP218652
0.98	Melanocytes	SRP257883
0.96	Microglia	All
0.95	Microglia	E-MTAB-7316
1.00	Microglia	EGAD00001006350
0.96	Microglia	SRP050054
0.96	Microglia	SRP158528
1.00	Microglia	SRP194595
1.00	Microglia	SRP222001

AUC	Cell Type	Study
0.90	Microglia	SRP222958
1.00	Microglia	SRP255195
0.97	Monocyte	All
0.99	Monocyte	EGAD00001006350
0.99	Muller Glia	All
1.00	Muller Glia	E-MTAB-7316
1.00	Muller Glia	EGAD00001006350
1.00	Muller Glia	SRP050054
1.00	Muller Glia	SRP075719
0.80	Muller Glia	SRP151023
0.74	Muller Glia	SRP158081
1.00	Muller Glia	SRP158528
1.00	Muller Glia	SRP194595
1.00	Muller Glia	SRP222001
0.99	Muller Glia	SRP222958
0.99	Muller Glia	SRP223254
1.00	Muller Glia	SRP255195
0.89	Natural Killer	All
0.97	Natural Killer	EGAD00001006350
0.81	Neurogenic Cells	All
0.76	Neurogenic Cells	SRP151023
0.85	Neurogenic Cells	SRP158081
0.63	Neurogenic Cells	SRP223254
0.95	Pericytes	All
0.97	Pericytes	EGAD00001006350
0.98	Pericytes	SRP158528
0.14	Pericytes	SRP218652
0.98	Pericytes	SRP257883
0.69	Photoreceptor Precursors	All
0.48	Photoreceptor Precursors	SRP151023
0.91	Photoreceptor Precursors	SRP158081
0.40	Photoreceptor Precursors	SRP223254
0.98	Red Blood Cells	All
1.00	Red Blood Cells	SRP131661
0.96	Red Blood Cells	SRP158081
0.78	Red Blood Cells	SRP257883
0.99	Retinal Ganglion Cells	All
0.10	Retinal Ganglion Cells	E-MTAB-7316
1.00	Retinal Ganglion Cells	EGAD00001006350
1.00	Retinal Ganglion Cells	SRP050054
0.98	Retinal Ganglion Cells	SRP151023
0.93	Retinal Ganglion Cells	SRP158081
1.00	Retinal Ganglion Cells	SRP158528
1.00	Retinal Ganglion Cells	SRP222958
0.98	Retinal Ganglion Cells	SRP223254
1.00	Retinal Ganglion Cells	SRP255195
0.99	Rod Bipolar Cells	All
1.00	Rod Bipolar Cells	EGAD00001006350
1.00	Rod Bipolar Cells	SRP075719
0.99	Rods	All
1.00	Rods	E-MTAB-7316
1.00	Rods	EGAD00001006350
0.97	Rods	SRP050054
0.96	Rods	SRP151023
0.97	Rods	SRP158081
0.96	Rods	SRP158528
1.00	Rods	SRP194595
1.00	Rods	SRP222001
1.00	Rods	SRP222958
0.97	Rods	SRP223254
0.96	Rods	SRP255195
0.98	RPCs	All
0.99	RPCs	SRP151023
0.99	RPCs	SRP223254
0.94	RPE	All
1.00	RPE	EGAD00001006350
0.81	Schwann	All

AUC	Cell Type	Study
0.38	Schwann	SRP218652
1.00	Schwann	SRP257883
0.97	T Cell	All
0.98	T Cell	EGAD00001006350
0.99	T Cell	SRP131661
0.05	T Cell	SRP218652
0.99	T Cell	SRP257883
0.91	Vein	All
1.00	Vein	SRP257883

441 *Supplemental Table 2: Area under the precision recall curve (AUC) for each cell type, split by study. The*

442 *“All” study is the AUC score across all cells within the cell type*

CellType	CellType_predict	Count	Ratio
AC/HC Precursors	AC/HC Precursors	1,248	0.86
AC/HC Precursors	Neurogenic Cells	100	0.07
AC/HC Precursors	Horizontal Cells	64	0.04
Amacrine Cells	Amacrine Cells	65,815	0.96
Artery	Endothelial	129	0.79
Artery	Vein	32	0.20
Astrocytes	Astrocytes	1,103	0.97
B-Cell	Endothelial	235	0.48
B-Cell	B-Cell	119	0.24
B-Cell	Vein	65	0.13
B-Cell	Fibroblasts	54	0.11
B-Cell	Macrophage	14	0.03
Bipolar Cells	Bipolar Cells	27,417	0.97
Choriocapillaris	Vein	140	0.64
Choriocapillaris	Endothelial	78	0.36
Cones	Cones	7,771	0.90
Cones	Rods	452	0.05
Cones	PR Precursors	265	0.03
Early RPCs	Early RPCs	26,654	0.96
Early RPCs	Neurogenic Cells	607	0.02
Endothelial	Endothelial	781	0.84
Endothelial	Pericytes	94	0.10
Endothelial	Vein	27	0.03
Fibroblasts	Fibroblasts	1,585	0.96
Fibroblasts	Vein	48	0.03
Horizontal Cells	Horizontal Cells	6,902	0.96
Late RPCs	Late RPCs	20,162	0.95
Late RPCs	Early RPCs	496	0.02
Late RPCs	Neurogenic Cells	458	0.02
Macrophage	Macrophage	443	0.80
Macrophage	T-Cell	83	0.15
Macrophage	Mast	14	0.03
Mast	Mast	90	0.92
Mast	RPCs	3	0.03
Melanocytes	Melanocytes	252	0.97
Microglia	Microglia	560	0.94
Monocyte	Monocyte	316	0.95
Monocyte	Microglia	7	0.02
Muller Glia	Muller Glia	36,671	0.98
Natural Killer	Natural Killer	150	0.90
Natural Killer	T-Cell	11	0.07
Natural Killer	B-Cell	4	0.02
Neurogenic Cells	Neurogenic Cells	9,061	0.87
Neurogenic Cells	RPCs	392	0.04
Neurogenic Cells	Late RPCs	339	0.03
Neurogenic Cells	Early RPCs	225	0.02
Neurogenic Cells	PR Precursors	223	0.02
Pericytes	Pericytes	2,771	0.95
Pericytes	Vein	96	0.03

CellType	CellType_predict	Count	Ratio
PR Precursors	PR Precursors	6,377	0.90
PR Precursors	Cones	254	0.04
PR Precursors	Neurogenic Cells	230	0.03
PR Precursors	Rods	147	0.02
Red Blood Cells	Red Blood Cells	673	0.97
Retinal Ganglion Cells	Retinal Ganglion Cells	42,766	0.97
Rod Bipolar Cells	Rod Bipolar Cells	9,634	0.94
Rod Bipolar Cells	Bipolar Cells	589	0.06
Rods	Rods	53,731	0.97
Rods	PR Precursors	1,477	0.03
RPCs	RPCs	17,135	0.97
RPE	RPE	343	0.93
RPE	RPCs	13	0.04
Schwann	Schwann	220	0.77
Schwann	Fibroblasts	28	0.10
Schwann	Melanocytes	16	0.06
Schwann	Pericytes	12	0.04
Smooth Muscle Cell	Pericytes	38	0.84
Smooth Muscle Cell	Endothelial	7	0.16
T-Cell	T-Cell	923	0.86
T-Cell	Macrophage	91	0.09
T-Cell	Natural Killer	29	0.03
Vein	Vein	481	0.98

443 *Supplemental Table 3: Counts of cell type labels with our xgboost machine learning system (PredCellType)*
 444 *and the published cell type labels (TrueCellType). Ratio is calculated as CellType that were labeled as*
 445 *CellType_predict. Ratio < 0.05 were filtered out from view in the table.*

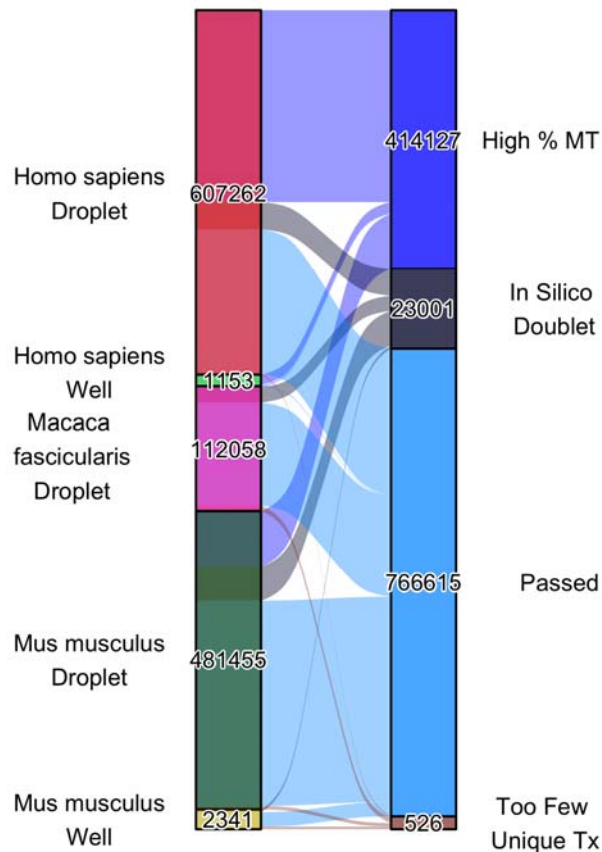
CellType	HS Studies (published)	MF Studies (published)	MM Studies (published)	HS Studies (transferred)	MF Studies (transferred)	MM Studies (transferred)
AC/HC Precursors	2	0	0	4	1	3
Amacrine Cells	8	1	5	11	1	15
Artery	1	0	0	0	0	0
Astrocytes	3	0	2	10	1	7
B-Cell	3	0	1	8	1	8
Bipolar Cells	8	1	4	11	1	14
Choriocapillaries	1	0	0	0	0	0
Cones	8	1	3	11	1	16
Endothelial	6	1	3	11	1	10
Fibroblasts	3	0	2	9	1	11
Horizontal Cells	7	1	3	10	1	11
Macrophage	2	0	0	8	0	4
Mast	2	0	0	9	0	3
Melanocytes	3	0	0	7	1	2
Microglia	6	1	2	13	1	10
Monocyte	1	0	1	9	0	8
Muller Glia	8	1	4	14	1	13
Natural Killer	1	0	0	5	0	2
Neurogenic Cells	2	0	1	5	1	7
Pericytes	3	1	1	9	1	7
PR Precursors	2	0	1	4	1	4

CellType	HS Studies (published)	MF Studies (published)	MM Studies (published)	HS Studies (transferred)	MF Studies (transferred)	MM Studies (transferred)
Red Blood Cells	1	0	2	11	1	12
Retinal Ganglion Cells	6	1	4	11	1	16
Rod Bipolar Cells	1	0	2	7	0	9
Rods	8	1	3	13	1	15
RPCs	2	0	0	10	1	10
RPE	2	0	0	7	1	7
Schwann	2	0	0	8	1	3
Smooth Muscle Cell	1	0	0	0	0	0
T-Cell	3	0	1	7	1	7
Unlabelled	15	1	18	11	1	16
Vein	1	0	0	5	0	3
Early RPCs	0	0	1	4	1	8
Late RPCs	0	0	1	3	0	6

446 *Supplemental Table 4: Counts for number of studies with cell types labels before and after cell type label*

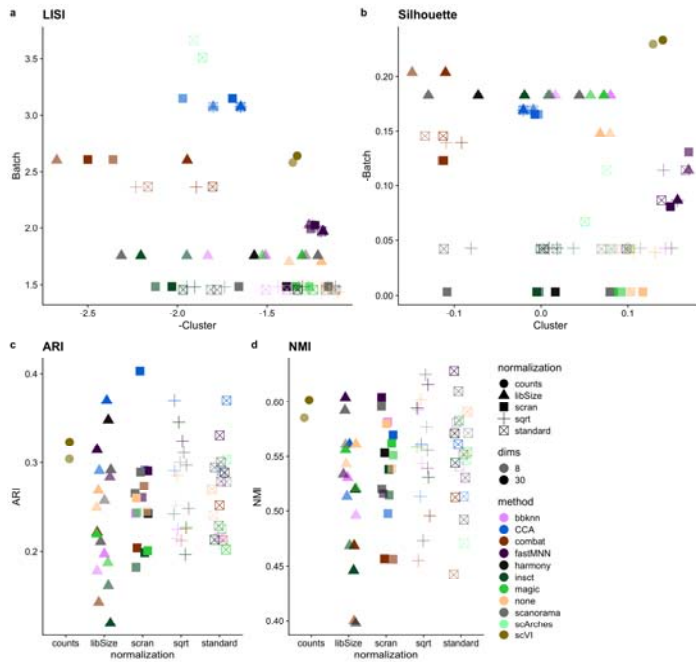
447 *transfer*

448 **Figures**



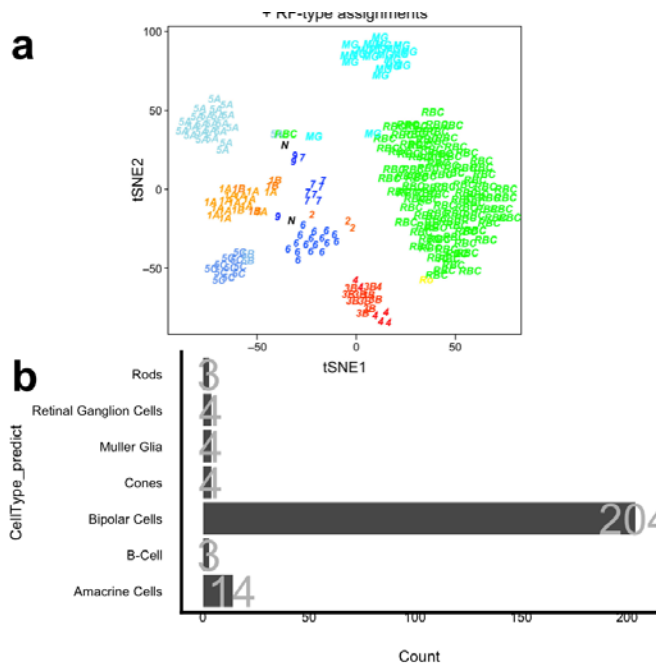
449

450 *Supplemental Figure 1: The left bar delineates the number of cells for each organism - technology combination. The*
 451 *right bar specifies the number of each cells in each post QC category. In silico doublets were identified with scrublet*
 452 *and DoubletDetector.*



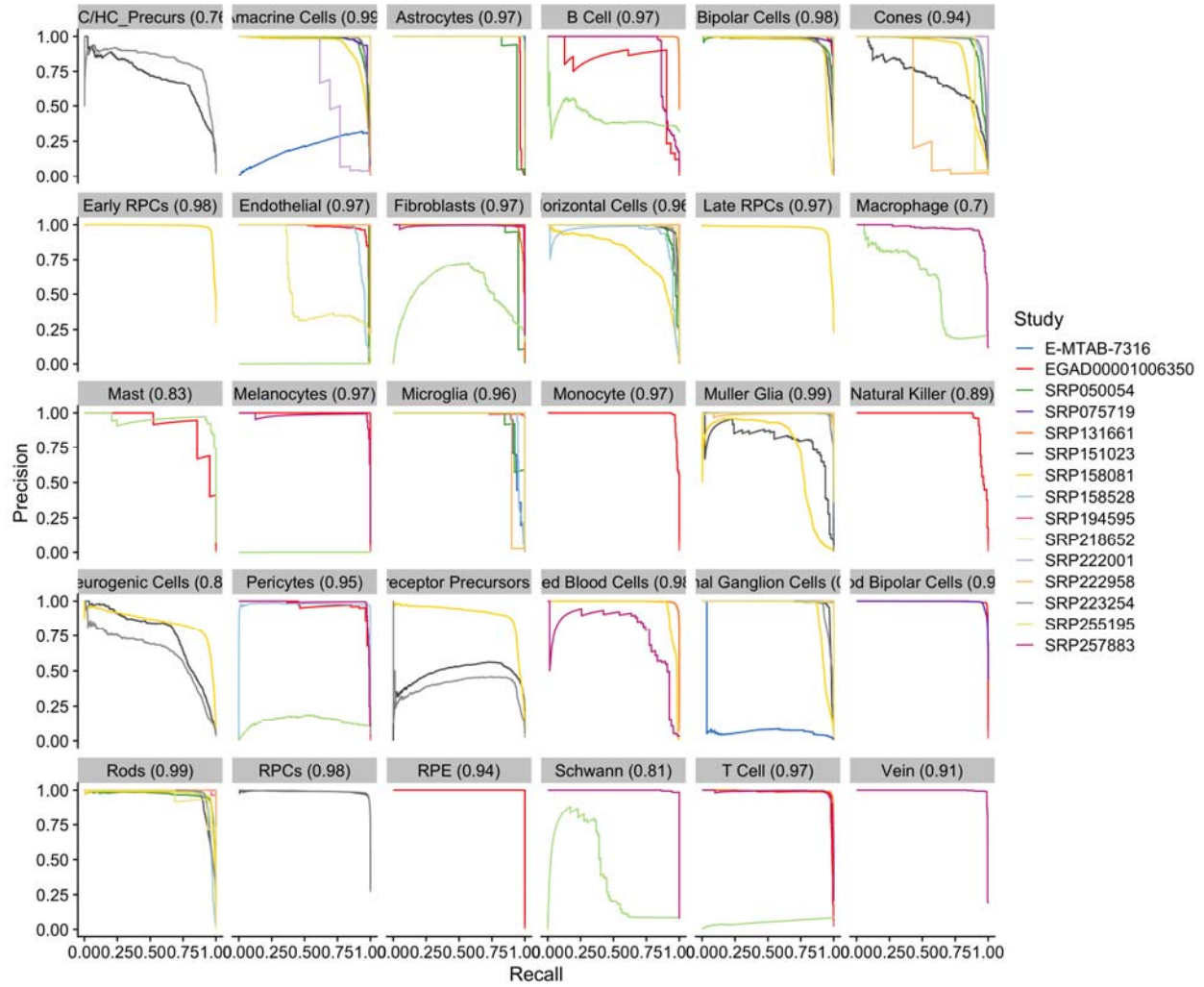
453

454 *Supplemental Figure 2: Performance of the various batch correction tools across various benchmarking metrics.*
 455 *For the LISI and Silhouette plots in A, B higher (y-axis) means better batch mixing and further to the right (x-axis)*
 456 *means better cluster purity. For the ARI and NMI metrics (which reflects how well cluster matches with cell type) in*
 457 *C, D, higher means a better score.*



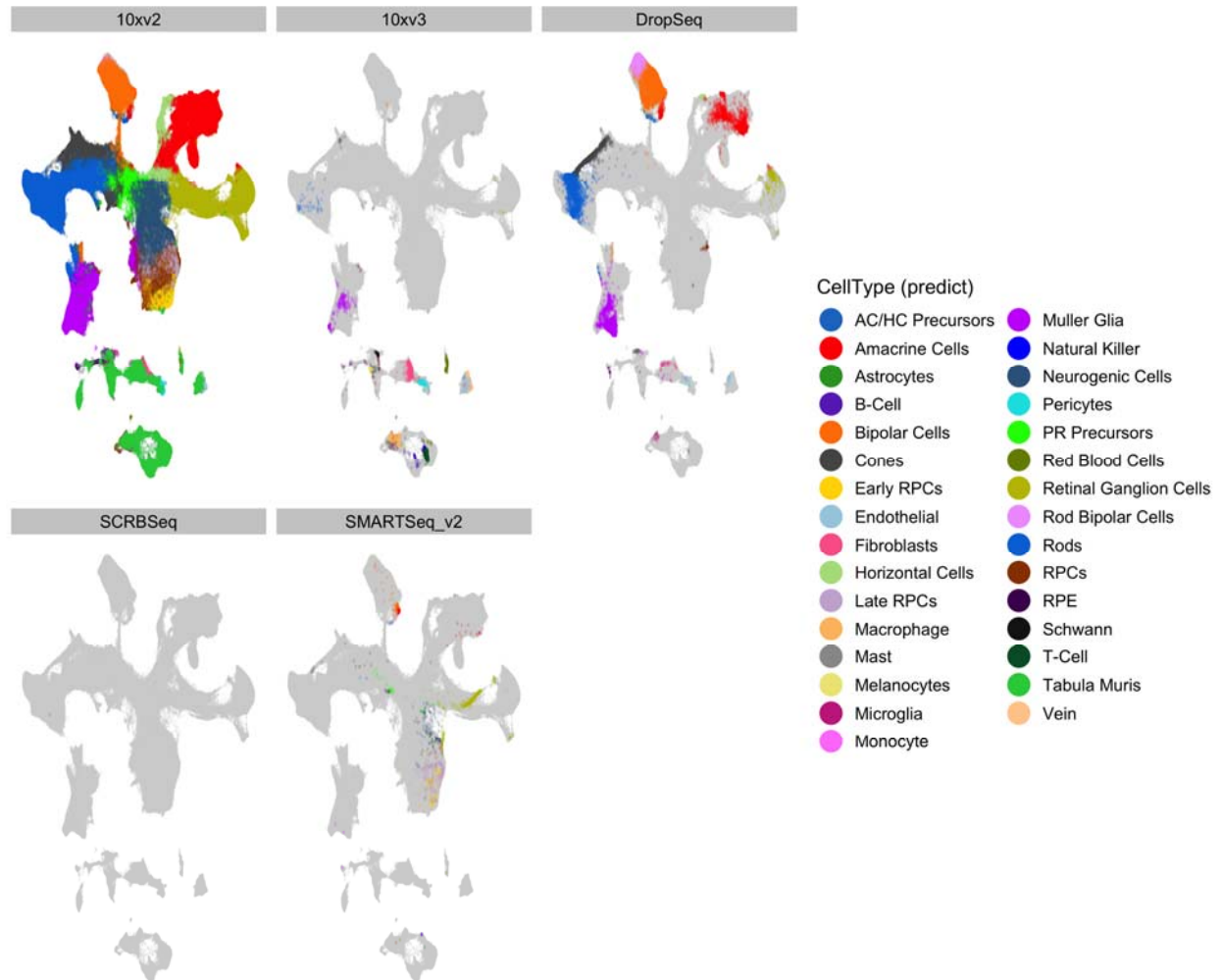
458

459 *Supplemental Figure 3: Our xgboost ML properly labels this Shekhar et al. RBC FAC sorted population as enriched*
 460 *in RBC*



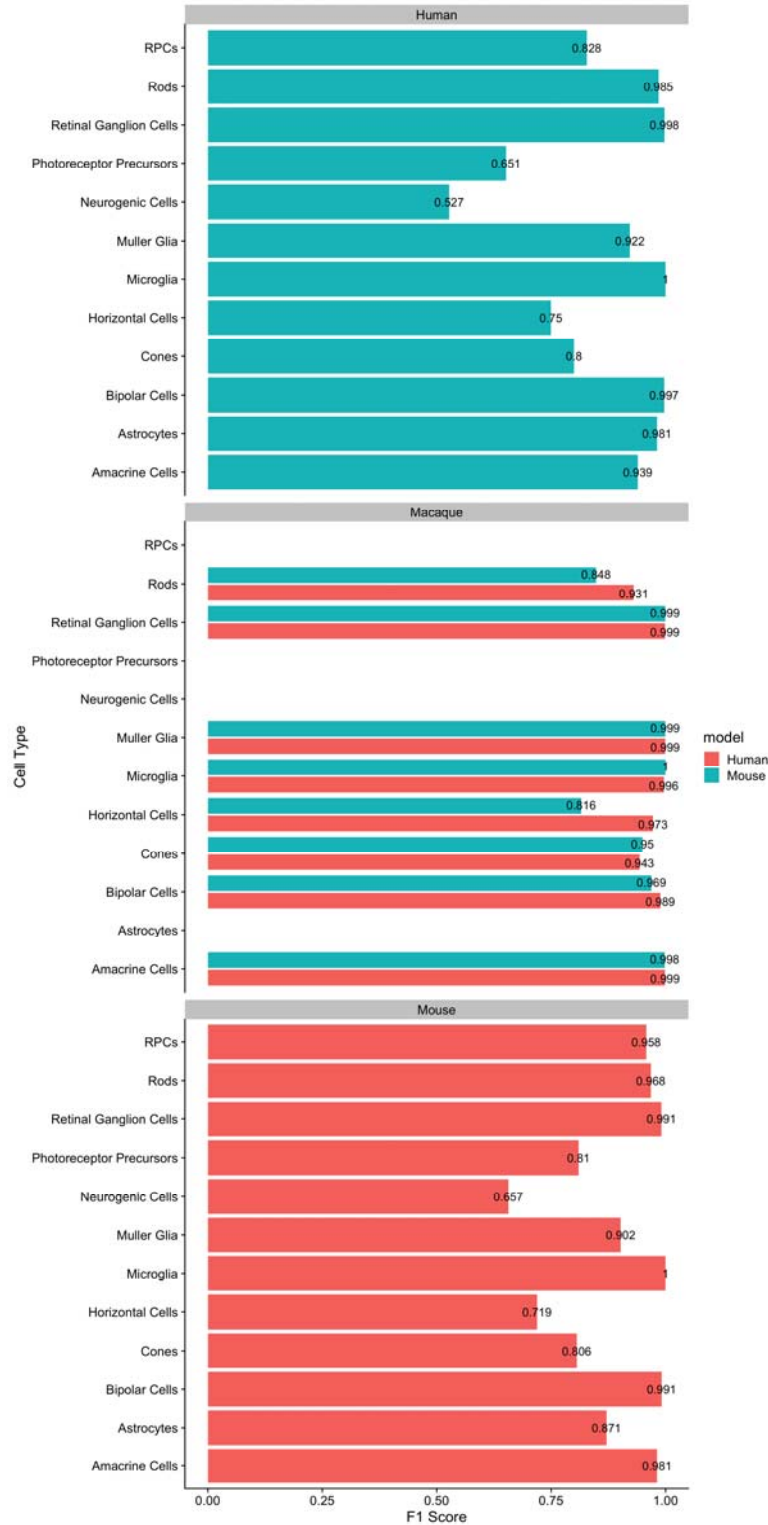
461

462 *Supplemental Figure 4: Precision recall curves for our xgboost cell type predictor model across each cell type*
 463 *predicted*



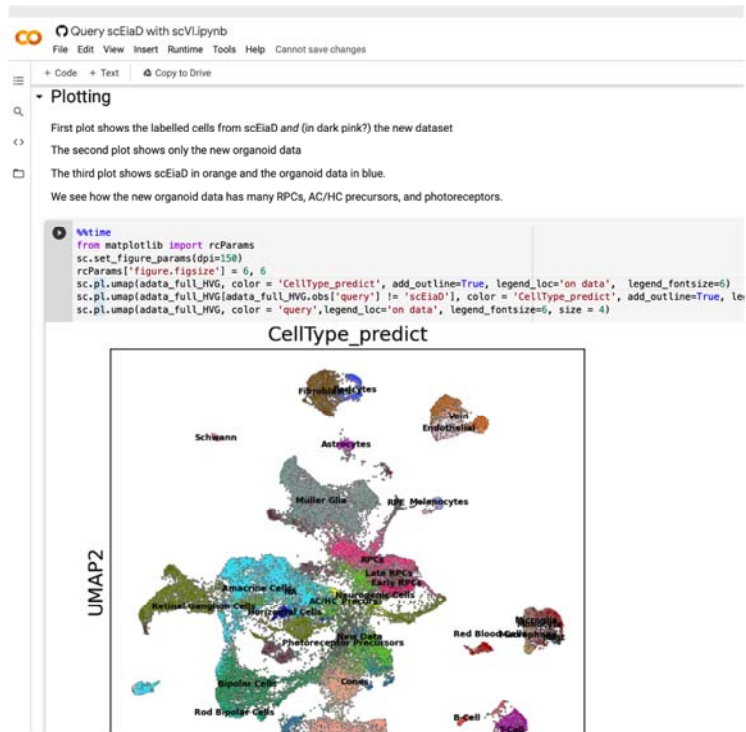
464

465 *Supplemental Figure 5: Distribution of cell types across the 2D UMAP confirms that the cell types are being*
466 *properly placed despite the wide variety of single sequencing platforms present.*



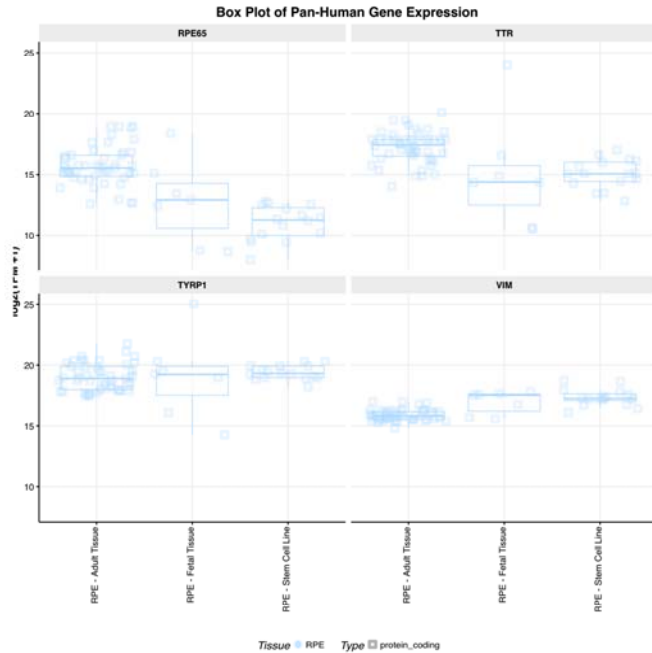
467

468 *Supplemental Figure 6: F1 scores (1 is perfect) of cell type prediction when using a human or mouse based xgboost*
 469 *cell type prediction model on other organisms (human, macaque, mouse).*



470

471 *Supplemental Figure 7: Screen shot of Google colab notebook that demonstrates how to integrate external data into*
472 *the scEiaD resource. We see in the screenshot how the organoid dataset contains many of the retinal cell types.*



473

474 *Supplemental Figure 8: Screenshot of eyeIntegration bulk RNA-seq meta-analysis of vimentin expression in different*
475 *RPE tissue sources*

476 Works Cited

477 1. Masland RH. The Neuronal Organization of the Retina. *Neuron*. 2012;76(2):266-280.

478 doi:[10.1016/j.neuron.2012.10.002](https://doi.org/10.1016/j.neuron.2012.10.002)

479 2. Annies M, Kröger S. Isoform Pattern and AChR Aggregation Activity of Agrin Expressed by Embryonic Chick
480 Retinal Ganglion Neurons. *Molecular and Cellular Neuroscience*. 2002;20(3):525-535.

481 doi:[10.1006/mcne.2002.1125](https://doi.org/10.1006/mcne.2002.1125)

482 3. Hagstrom SA, Neitz M, Neitz J. Cone pigment gene expression in individual photoreceptors and the chromatic
483 topography of the retina. *J Opt Soc Am A Opt Image Sci Vis*. 2000;17(3):527-537. doi:[10.1364/josaa.17.000527](https://doi.org/10.1364/josaa.17.000527)

484 4. Trimarchi JM, Stadler MB, Roska B, et al. Molecular heterogeneity of developing retinal ganglion and amacrine
485 cells revealed through single cell gene expression profiling. *Journal of Comparative Neurology*. 2007;502(6):1047-

486 1065. doi:[10.1002/cne.21368](https://doi.org/10.1002/cne.21368)

- 487 5. Wahlin KJ, Lim L, Grice EA, Campochiaro PA, Zack DJ, Adler R. A method for analysis of gene expression in
488 isolated mouse photoreceptor and Müller cells. *Mol Vis*. 2004;10:366-375.
- 489 6. Svensson V, da Veiga Beltrame E, Pachter L. A curated database reveals trends in single-cell transcriptomics.
490 *Database*. 2020;2020(baaa073). doi:[10.1093/database/baaa073](https://doi.org/10.1093/database/baaa073)
- 491 7. Macosko EZ, Basu A, Satija R, et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells
492 Using Nanoliter Droplets. *Cell*. 2015;161(5):1202-1214. doi:[10.1016/j.cell.2015.05.002](https://doi.org/10.1016/j.cell.2015.05.002)
- 493 8. Buenaventura DF, Corseri A, Emerson MM. Identification of Genes With Enriched Expression in Early
494 Developing Mouse Cone Photoreceptors. *Invest Ophthalmol Vis Sci*. 2019;60(8):2787-2799. doi:[10.1167/iovs.19-](https://doi.org/10.1167/iovs.19-26951)
495 [26951](https://doi.org/10.1167/iovs.19-26951)
- 496 9. Clark BS, Stein-O'Brien GL, Shiau F, et al. Single-Cell RNA-Seq Analysis of Retinal Development Identifies
497 NFI Factors as Regulating Mitotic Exit and Late-Born Cell Specification. *Neuron*. Published online May 22, 2019.
498 doi:[10.1016/j.neuron.2019.04.010](https://doi.org/10.1016/j.neuron.2019.04.010)
- 499 10. Cowan CS, Renner M, De Gennaro M, et al. Cell Types of the Human Retina and Its Organoids at Single-Cell
500 Resolution. *Cell*. 2020;182(6):1623-1640.e34. doi:[10.1016/j.cell.2020.08.013](https://doi.org/10.1016/j.cell.2020.08.013)
- 501 11. Dharmat R, Kim S, Liu H, Fu S, Li Y, Chen R. Epigenetic adaptation prolongs photoreceptor survival during
502 retinal degeneration. *bioRxiv*. Published online September 18, 2019:774950. doi:[10.1101/774950](https://doi.org/10.1101/774950)
- 503 12. Fadl BR, Brodie SA, Malasky M, et al. An optimized protocol for retina single-cell RNA sequencing. *Mol Vis*.
504 2020;26:705-717. Accessed March 26, 2021. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7553720/>
- 505 13. Hu Y, Wang X, Hu B, et al. Dissecting the transcriptome landscape of the human fetal neural retina and retinal
506 pigment epithelium by single-cell RNA-seq analysis. *PLoS Biol*. 2019;17(7):e3000365.
507 doi:[10.1371/journal.pbio.3000365](https://doi.org/10.1371/journal.pbio.3000365)
- 508 14. Lehmann GL, Hanke-Gogokhia C, Hu Y, et al. Single-cell profiling reveals an endothelium-mediated
509 immunomodulatory pathway in the eye choroid. *Journal of Experimental Medicine*. 2020;217(e20190730).
510 doi:[10.1084/jem.20190730](https://doi.org/10.1084/jem.20190730)

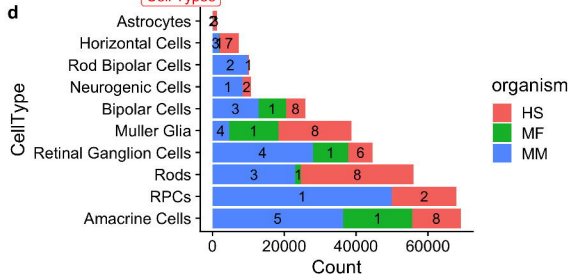
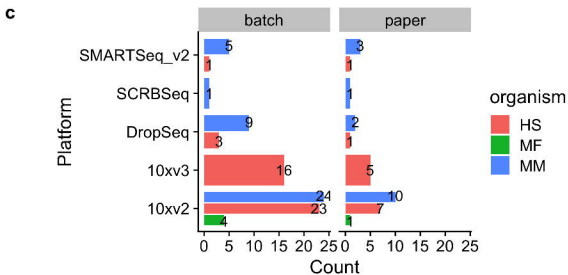
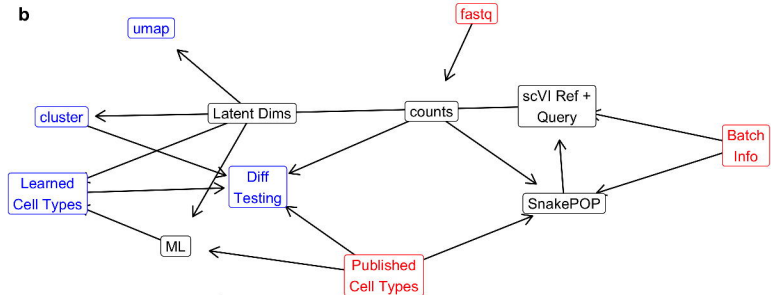
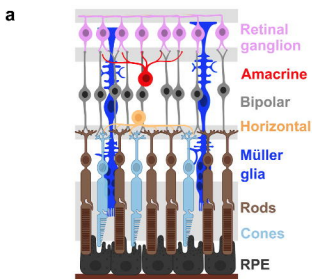
- 511 15. Lo Giudice Q, Leleu M, La Manno G, Fabre PJ. Single-cell transcriptional logic of cell-fate specification and
512 axon guidance in early-born retinal neurons. *Development*. 2019;146(17). doi:[10.1242/dev.178103](https://doi.org/10.1242/dev.178103)
- 513 16. Lukowski SW, Lo CY, Sharov AA, et al. A single-cell transcriptome atlas of the adult human retina. *EMBO J*.
514 2019;38(18):e100811. doi:[10.15252/emj.2018100811](https://doi.org/10.15252/emj.2018100811)
- 515 17. Lu Y, Shiao F, Yi W, et al. Single-Cell Analysis of Human Retina Identifies Evolutionarily Conserved and
516 Species-Specific Mechanisms Controlling Development. *Dev Cell*. 2020;53(4):473-491.e9.
517 doi:[10.1016/j.devcel.2020.04.009](https://doi.org/10.1016/j.devcel.2020.04.009)
- 518 18. Menon M, Mohammadi S, Davila-Velderrain J, et al. Single-cell transcriptomic atlas of the human retina
519 identifies cell types associated with age-related macular degeneration. *Nat Commun*. 2019;10(1):4902.
520 doi:[10.1038/s41467-019-12780-8](https://doi.org/10.1038/s41467-019-12780-8)
- 521 19. O’Koren EG, Yu C, Klingeborn M, et al. Microglial Function Is Distinct in Different Anatomical Locations
522 during Retinal Homeostasis and Degeneration. *Immunity*. 2019;50(3):723-737.e7.
523 doi:[10.1016/j.immuni.2019.02.007](https://doi.org/10.1016/j.immuni.2019.02.007)
- 524 20. Peng Y-R, Shekhar K, Yan W, et al. Molecular Classification and Comparative Taxonomics of Foveal and
525 Peripheral Cells in Primate Retina. *Cell*. 2019;176(5):1222-1237.e22. doi:[10.1016/j.cell.2019.01.004](https://doi.org/10.1016/j.cell.2019.01.004)
- 526 21. Shekhar K, Lapan SW, Whitney IE, et al. Comprehensive Classification of Retinal Bipolar Neurons by Single-
527 Cell Transcriptomics. *Cell*. 2016;166(5):1308-1323.e30. doi:[10.1016/j.cell.2016.07.054](https://doi.org/10.1016/j.cell.2016.07.054)
- 528 22. Sridhar A, Hoshino A, Finkbeiner CR, et al. Single-Cell Transcriptomic Comparison of Human Fetal Retina,
529 hPSC-Derived Retinal Organoids, and Long-Term Retinal Cultures. *Cell Rep*. 2020;30(5):1644-1659.e4.
530 doi:[10.1016/j.celrep.2020.01.007](https://doi.org/10.1016/j.celrep.2020.01.007)
- 531 23. Tran NM, Shekhar K, Whitney IE, et al. Single-Cell Profiles of Retinal Ganglion Cells Differing in Resilience to
532 Injury Reveal Neuroprotective Genes. *Neuron*. 2019;104(6):1039-1055.e12. doi:[10.1016/j.neuron.2019.11.006](https://doi.org/10.1016/j.neuron.2019.11.006)

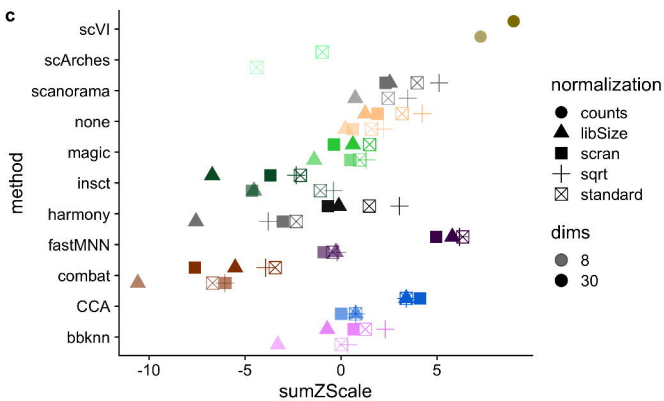
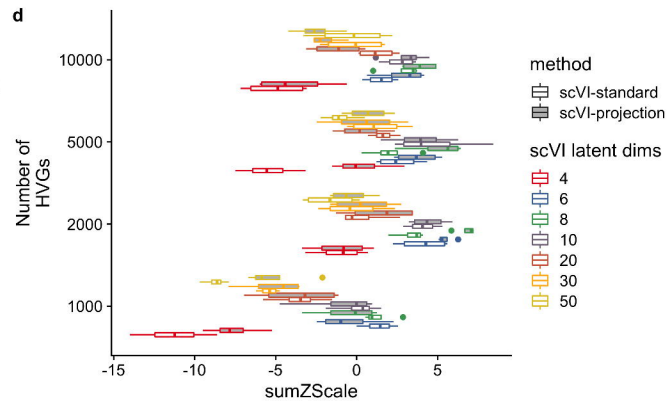
- 533 24. Voigt AP, Whitmore SS, Mulfaul K, et al. Bulk and single-cell gene expression analyses reveal aging human
534 choriocapillaris has pro-inflammatory phenotype. *Microvascular Research*. 2020;131:104031.
535 doi:[10.1016/j.mvr.2020.104031](https://doi.org/10.1016/j.mvr.2020.104031)
- 536 25. Voigt AP, Whitmore SS, Flamme-Wiese MJ, et al. Molecular characterization of foveal versus peripheral human
537 retina by single-cell RNA sequencing. *Experimental Eye Research*. 2019;184:234-242.
538 doi:[10.1016/j.exer.2019.05.001](https://doi.org/10.1016/j.exer.2019.05.001)
- 539 26. Voigt AP, Binkley E, Flamme-Wiese MJ, et al. Single-Cell RNA Sequencing in Human Retinal Degeneration
540 Reveals Distinct Glial Cell Populations. *Cells*. 2020;9(2). doi:[10.3390/cells9020438](https://doi.org/10.3390/cells9020438)
- 541 27. Voigt AP, Mulfaul K, Mullin NK, et al. Single-cell transcriptomics of the human retinal pigment epithelium and
542 choroid in health and macular degeneration. *Proc Natl Acad Sci U S A*. 2019;116(48):24100-24107.
543 doi:[10.1073/pnas.1914143116](https://doi.org/10.1073/pnas.1914143116)
- 544 28. Yan W, Peng Y-R, van Zyl T, et al. Cell Atlas of The Human Fovea and Peripheral Retina. *Sci Rep*.
545 2020;10(1):9802. doi:[10.1038/s41598-020-66092-9](https://doi.org/10.1038/s41598-020-66092-9)
- 546 29. Yan W, Laboulaye MA, Tran NM, Whitney IE, Benhar I, Sanes JR. Mouse Retinal Cell Atlas: Molecular
547 Identification of over Sixty Amacrine Cell Types. *J Neurosci*. 2020;40(27):5177-5195.
548 doi:[10.1523/JNEUROSCI.0471-20.2020](https://doi.org/10.1523/JNEUROSCI.0471-20.2020)
- 549 30. Melsted P, Boeshaghi AS, Gao F, et al. Modular and efficient pre-processing of single-cell RNA-seq. *bioRxiv*.
550 Published online July 26, 2019:673285. doi:[10.1101/673285](https://doi.org/10.1101/673285)
- 551 31. Srivastava A, Malik L, Smith T, Sudbery I, Patro R. Alevin efficiently estimates accurate gene abundances from
552 dscRNA-seq data. *Genome Biology*. 2019;20(1):65. doi:[10.1186/s13059-019-1670-y](https://doi.org/10.1186/s13059-019-1670-y)
- 553 32. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different
554 conditions, technologies, and species. *Nature Biotechnology*. 2018;36(5, 5):411-420. doi:[10.1038/nbt.4096](https://doi.org/10.1038/nbt.4096)

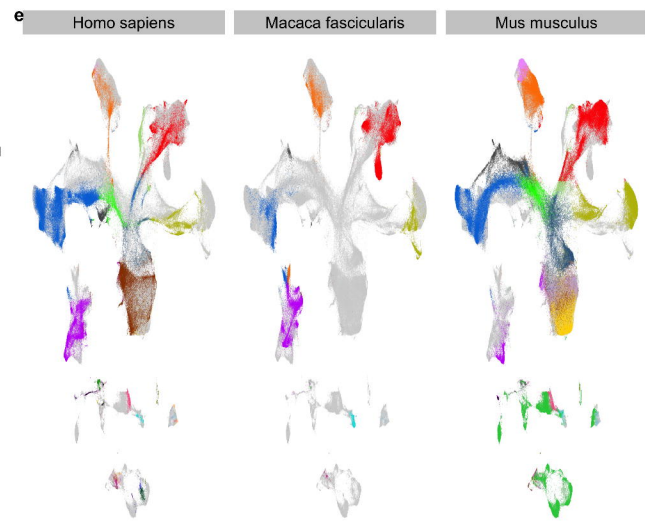
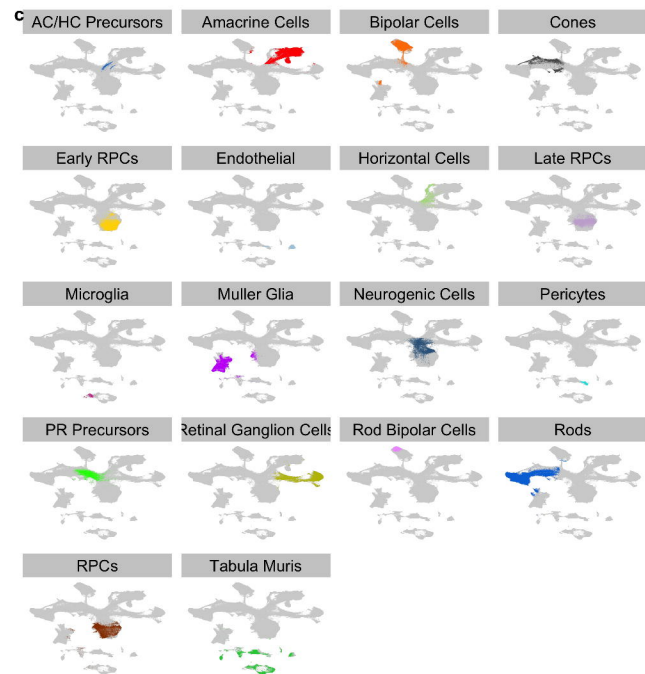
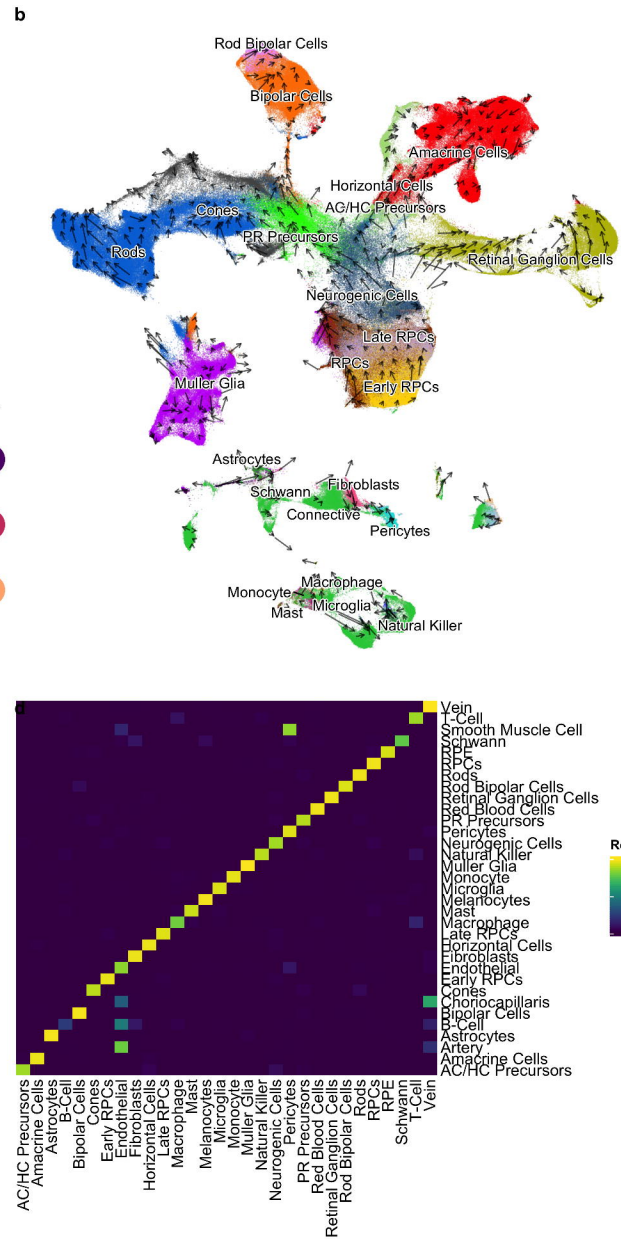
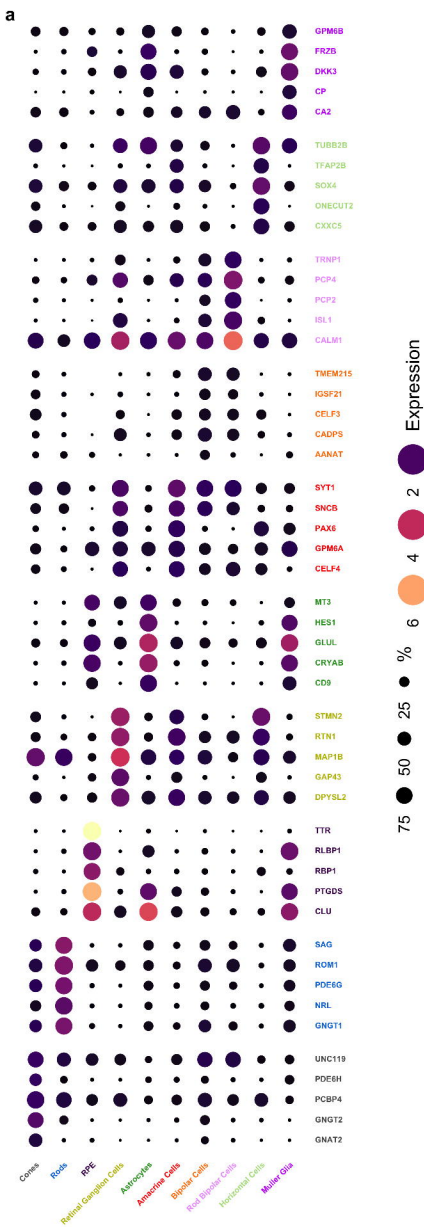
- 555 33. Haghverdi L, Lun ATL, Morgan MD, Marioni JC. Batch effects in single-cell RNA-sequencing data are
556 corrected by matching mutual nearest neighbors. *Nature Biotechnology*. 2018;36(5, 5):421-427.
557 doi:[10.1038/nbt.4091](https://doi.org/10.1038/nbt.4091)
- 558 34. Hie B, Bryson B, Berger B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama.
559 *Nature Biotechnology*. 2019;37(6, 6):685-691. doi:[10.1038/s41587-019-0113-3](https://doi.org/10.1038/s41587-019-0113-3)
- 560 35. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes
561 methods. *Biostatistics*. 2007;8(1):118-127. doi:[10.1093/biostatistics/kxj037](https://doi.org/10.1093/biostatistics/kxj037)
- 562 36. Korsunsky I, Millard N, Fan J, et al. Fast, sensitive and accurate integration of single-cell data with Harmony.
563 *Nature Methods*. 2019;16(12, 12):1289-1296. doi:[10.1038/s41592-019-0619-0](https://doi.org/10.1038/s41592-019-0619-0)
- 564 37. Liu J, Gao C, Sodicoff J, Kozareva V, Macosko EZ, Welch JD. Jointly defining cell types from multiple single-
565 cell datasets using LIGER. *Nature Protocols*. 2020;15(11, 11):3632-3662. doi:[10.1038/s41596-020-0391-8](https://doi.org/10.1038/s41596-020-0391-8)
- 566 38. Lopez R, Regier J, Cole MB, Jordan MI, Yosef N. Deep generative modeling for single-cell transcriptomics.
567 *Nature Methods*. 2018;15(12, 12):1053-1058. doi:[10.1038/s41592-018-0229-2](https://doi.org/10.1038/s41592-018-0229-2)
- 568 39. Polański K, Young MD, Miao Z, Meyer KB, Teichmann SA, Park J-E. BBKNN: Fast batch alignment of single
569 cell transcriptomes. *Bioinformatics*. 2020;36(3):964-965. doi:[10.1093/bioinformatics/btz625](https://doi.org/10.1093/bioinformatics/btz625)
- 570 40. Query to reference single-cell integration with transfer learning | bioRxiv. Accessed December 20, 2020.
571 <https://www.biorxiv.org/content/10.1101/2020.07.16.205997v1>
- 572 41. Simon LM, Wang Y-Y, Zhao Z. INSCT: Integrating millions of single cells using batch-aware triplet neural
573 networks. *bioRxiv*. Published online May 17, 2020:2020.05.16.100024. doi:[10.1101/2020.05.16.100024](https://doi.org/10.1101/2020.05.16.100024)
- 574 42. Stuart T, Butler A, Hoffman P, et al. Comprehensive Integration of Single-Cell Data. *Cell*. 2019;177(7):1888-
575 1902.e21. doi:[10.1016/j.cell.2019.05.031](https://doi.org/10.1016/j.cell.2019.05.031)
- 576 43. van Dijk D, Sharma R, Nainys J, et al. Recovering Gene Interactions from Single-Cell Data Using Data
577 Diffusion. *Cell*. 2018;174(3):716-729.e27. doi:[10.1016/j.cell.2018.05.061](https://doi.org/10.1016/j.cell.2018.05.061)

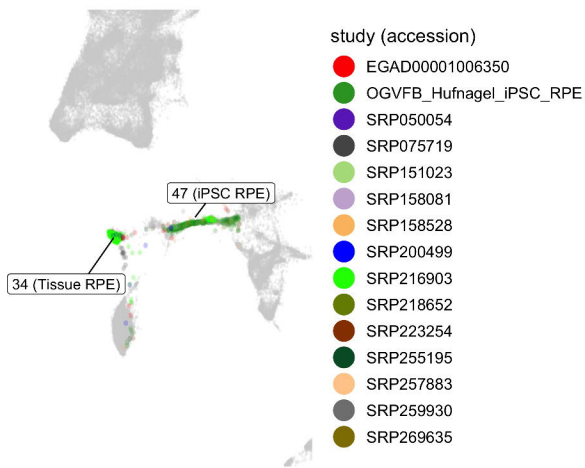
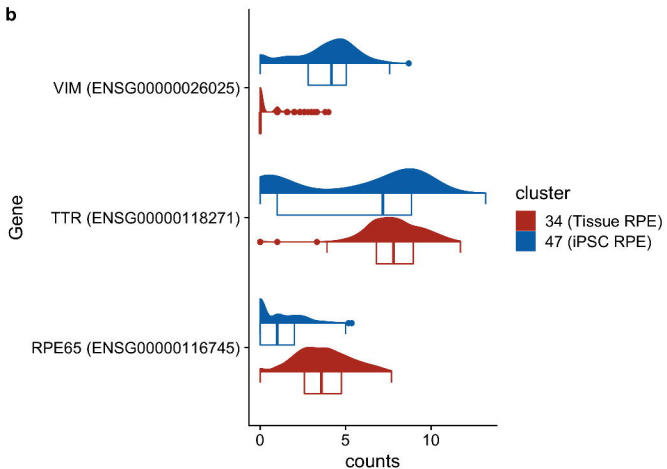
- 578 44. Kallman A, Capowski EE, Wang J, et al. Investigating cone photoreceptor development using patient-derived
579 NRL null retinal organoids. *Commun Biol.* 2020;3(1):82. doi:[10.1038/s42003-020-0808-5](https://doi.org/10.1038/s42003-020-0808-5)
- 580 45. Hunt RC, Davis AA. Altered expression of keratin and vimentin in human retinal pigment epithelial cells in vivo
581 and in vitro. *J Cell Physiol.* 1990;145(2):187-199. doi:[10.1002/jcp.1041450202](https://doi.org/10.1002/jcp.1041450202)
- 582 46. Swamy V, McGaughey D. Eye in a Disk: eyeIntegration Human Pan-Eye and Body Transcriptome Database
583 Version 1.0. *Invest Ophthalmol Vis Sci.* 2019;60(8):3236-3246. doi:[10.1167/iovs.19-27106](https://doi.org/10.1167/iovs.19-27106)
- 584 47. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics.*
585 2012;28(19):2520-2522. doi:[10.1093/bioinformatics/bts480](https://doi.org/10.1093/bioinformatics/bts480)
- 586 48. Harrow J, Frankish A, Gonzalez JM, et al. GENCODE: The reference human genome annotation for The
587 ENCODE Project. *Genome Res.* 2012;22(9):1760-1774. doi:[10.1101/gr.135350.111](https://doi.org/10.1101/gr.135350.111)
- 588 49. J H, A F, Jm G, et al. GENCODE: The reference human genome annotation for The ENCODE Project. *Genome*
589 *Res.* 2012;22(9):1760-1774. doi:[10.1101/gr.135350.111](https://doi.org/10.1101/gr.135350.111)
- 590 50. Yates AD, Achuthan P, Akanni W, et al. Ensembl 2020. *Nucleic Acids Research.* 2020;48(D1):D682-D688.
591 doi:[10.1093/nar/gkz966](https://doi.org/10.1093/nar/gkz966)
- 592 51. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol.*
593 2016;34(5):525-527. doi:[10.1038/nbt.3519](https://doi.org/10.1038/nbt.3519)
- 594 52. Melsted P, Ntranos V, Pachter L. The barcode, UMI, set format and BUStools. *Bioinformatics.*
595 2019;35(21):4472-4473. doi:[10.1093/bioinformatics/btz279](https://doi.org/10.1093/bioinformatics/btz279)
- 596 53. Lun ATL, Riesenfeld S, Andrews T, et al. EmptyDrops: Distinguishing cells from empty droplets in droplet-
597 based single-cell RNA sequencing data. *Genome Biology.* 2019;20(1):63. doi:[10.1186/s13059-019-1662-y](https://doi.org/10.1186/s13059-019-1662-y)
- 598 54. Lun ATL, McCarthy DJ, Marioni JC. A step-by-step workflow for low-level analysis of single-cell RNA-seq
599 data with Bioconductor. *F1000Res.* 2016;5:2122. doi:[10.12688/f1000research.9501.2](https://doi.org/10.12688/f1000research.9501.2)

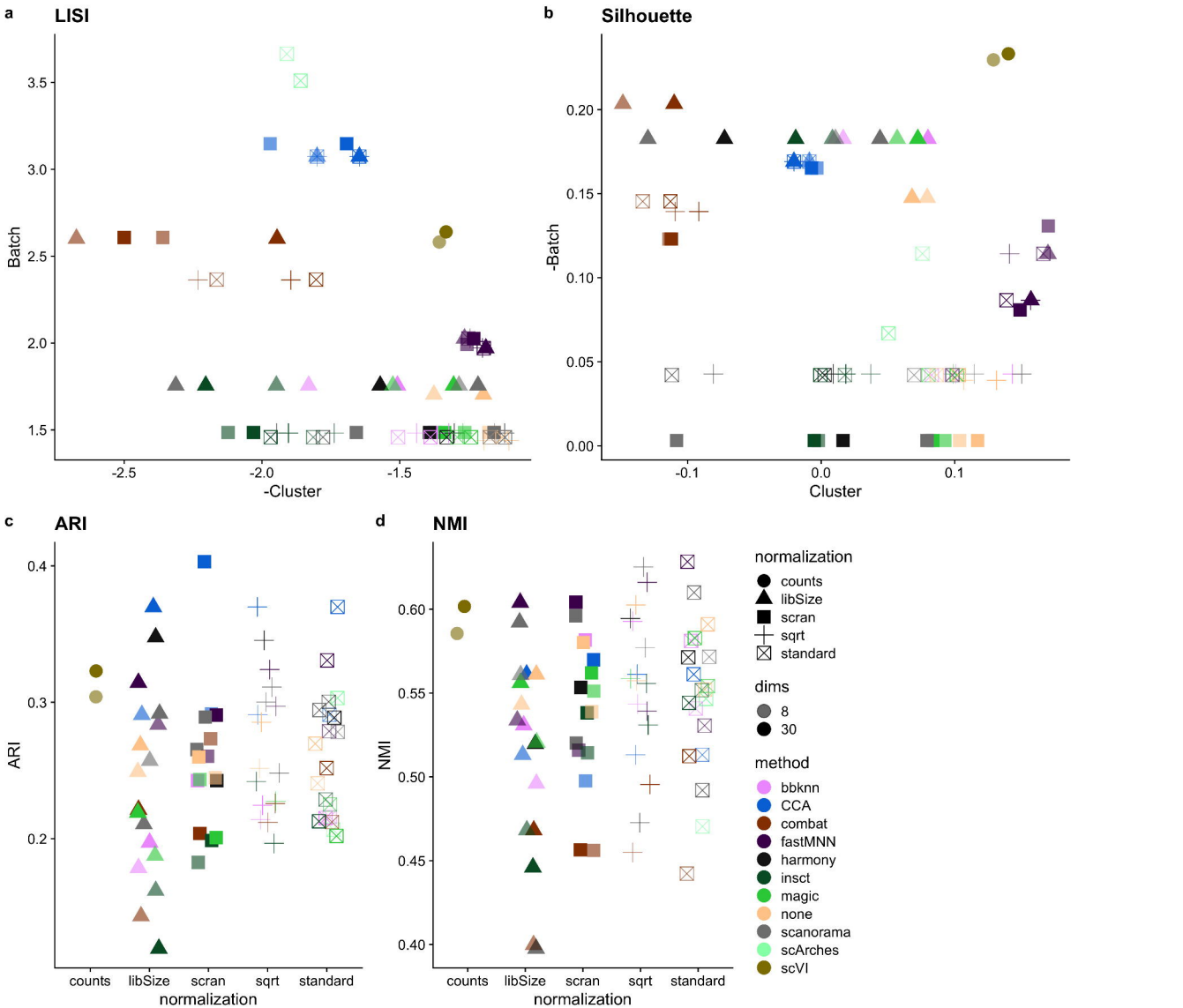
- 600 55. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat*
601 *Mech.* 2008;2008(10):P10008. doi:[10.1088/1742-5468/2008/10/P10008](https://doi.org/10.1088/1742-5468/2008/10/P10008)
- 602 56. Traag V, Waltman L, van Eck NJ. From Louvain to Leiden: Guaranteeing well-connected communities. *Sci Rep.*
603 2019;9(1):5233. doi:[10.1038/s41598-019-41695-z](https://doi.org/10.1038/s41598-019-41695-z)
- 604 57. Stassen SV, Siu DMD, Lee KCM, Ho JWK, So HKH, Tsia KK. PARC: Ultrafast and accurate clustering of
605 phenotypic data of millions of single cells. *Bioinformatics.* 2020;36(9):2778-2786.
606 doi:[10.1093/bioinformatics/btaa042](https://doi.org/10.1093/bioinformatics/btaa042)
- 607 58. McInnes L, Healy J, Melville J. UMAP: Uniform Manifold Approximation and Projection for Dimension
608 Reduction. Published September 17, 2020. Accessed December 21, 2020. <http://arxiv.org/abs/1802.03426>
- 609 59. Büttner M, Miao Z, Wolf FA, Teichmann SA, Theis FJ. A test metric for assessing single-cell RNA-seq batch
610 correction. *Nature Methods.* 2019;16(1, 1):43-49. doi:[10.1038/s41592-018-0254-1](https://doi.org/10.1038/s41592-018-0254-1)
- 611 60. Gayoso A, Shor J. *JonathanShor/DoubletDetection: Doubletdetection V3.0.* Zenodo; 2020.
612 doi:[10.5281/zenodo.4359992](https://doi.org/10.5281/zenodo.4359992)
- 613 61. Wolock SL, Lopez R, Klein AM. Scrublet: Computational Identification of Cell Doublets in Single-Cell
614 Transcriptomic Data. *Cell Systems.* 2019;8(4):281-291.e9. doi:[10.1016/j.cels.2018.11.005](https://doi.org/10.1016/j.cels.2018.11.005)
- 615 62. Bergen V, Lange M, Peidli S, Wolf FA, Theis FJ. Generalizing RNA velocity to transient cell states through
616 dynamical modeling. *Nature Biotechnology.* 2020;38(12, 12):1408-1414. doi:[10.1038/s41587-020-0591-3](https://doi.org/10.1038/s41587-020-0591-3)



a**b****c****d**



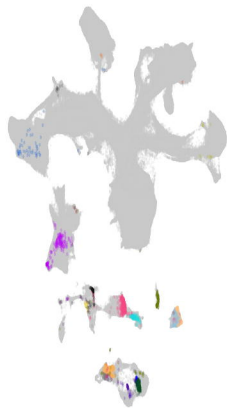
a**b**



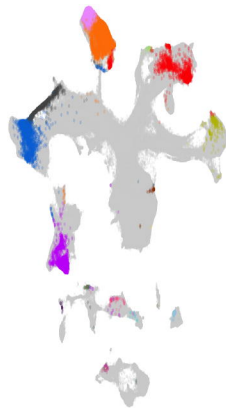
10xv2



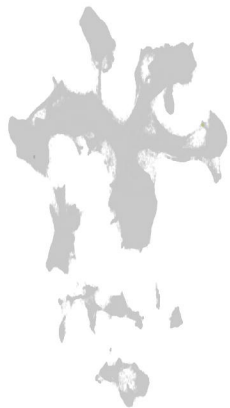
10xv3



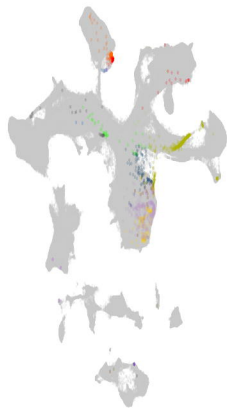
DropSeq



SCRBSeq

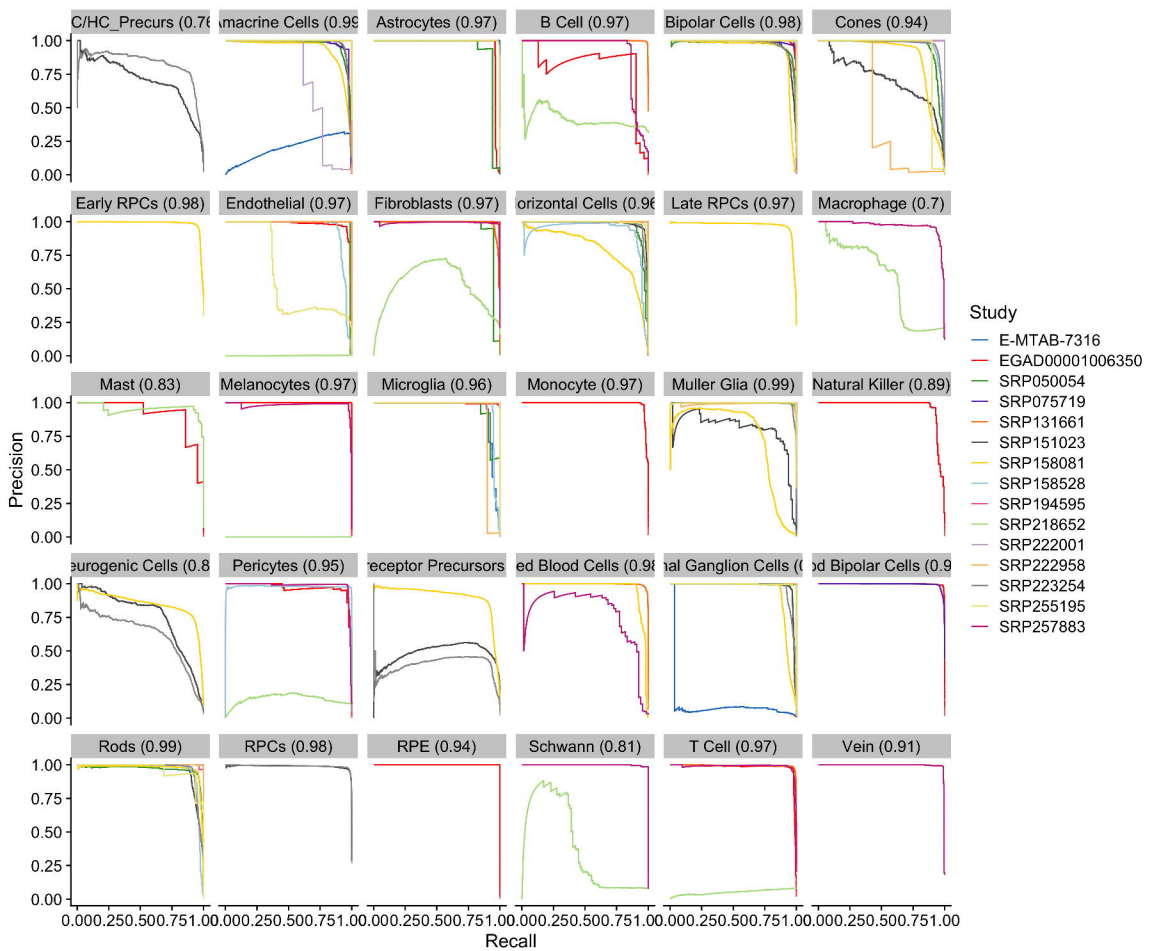


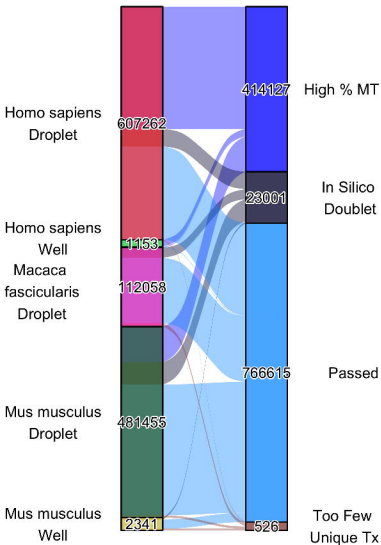
SMARTSeq_v2

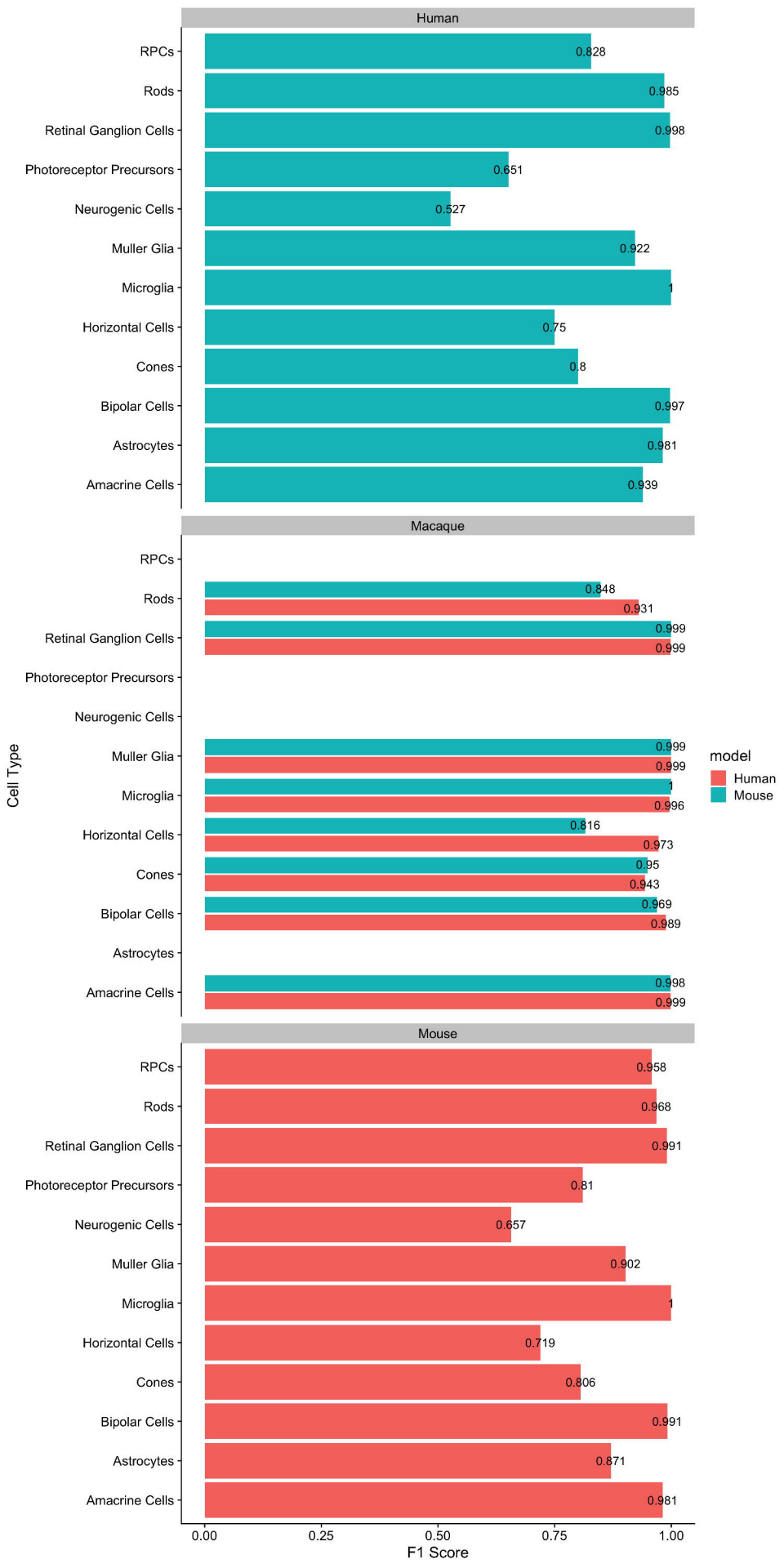


CellType (predict)











Query scEiaD with scVI.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

Plotting

First plot shows the labelled cells from scEiaD *and* (in dark pink?) the new dataset

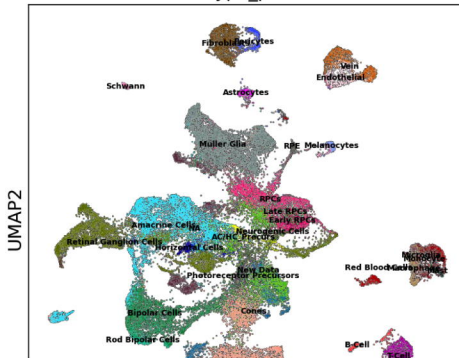
The second plot shows only the new organoid data

The third plot shows scEiaD in orange and the organoid data in blue.

We see how the new organoid data has many RPCs, AC/HC precursors, and photoreceptors.

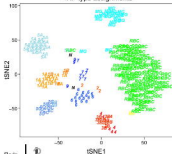
```
!time
from matplotlib import rcParams
sc.set_figure_params(dpi=150)
rcParams['figure.figsize'] = 5, 6
sc.pl.umap(adata_full_HVG, color = 'CellType_predict', add_outline=True, legend_loc='on data', legend_fontsize=6)
sc.pl.umap(adata_full_HVG[adata_full_HVG.obs['query'] != 'scEiaD'], color = 'CellType_predict', add_outline=True, legend_loc='on data', legend_fontsize=6, size = 4)
```

CellType_predict



↑ RP-type assignments

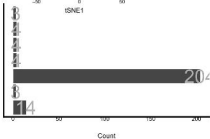
a



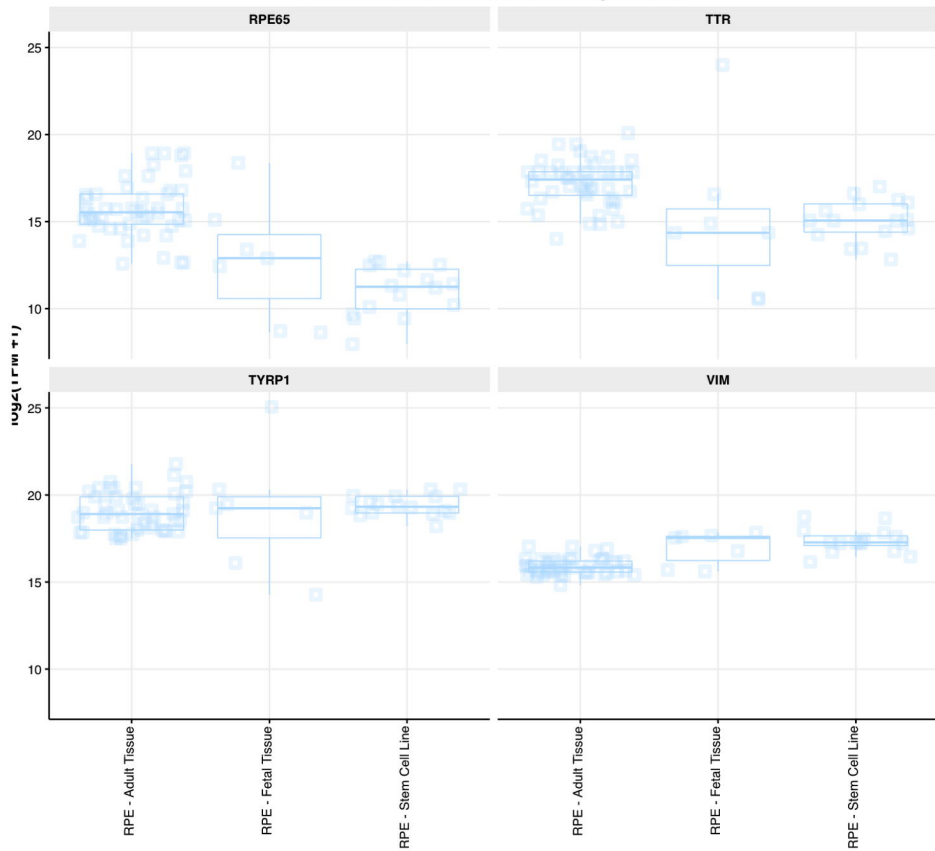
b

CellType_predist

Rods
 Retinal Ganglion Cells
 Muller Glia
 Cones
 Bipolar Cells
 B-Cell
 Amacrine Cells



Box Plot of Pan-Human Gene Expression



Tissue ● RPE ■ protein_coding