

1 **Integrating bulk and single cell RNA-seq refines transcriptomic profiles of**
2 **specific *C. elegans* neurons.**

3

4 Author list: Alec Barrett¹, Erdem Varol², Alexis Weinreb^{1,3}, Seth R. Taylor⁴,
5 Rebecca M. McWhirter⁴, Cyril Cros⁵, Manasa Basaravaju^{1,3}, Abigail Poff⁴, John
6 A. Tipps⁴, Maryam Majeed⁵, Berta Vidal⁵, Chen Wang⁵, Eviatar Yemini⁶, Emily A.
7 Bayer⁵, HaoSheng Sun⁷, Oliver Hobert^{5,8*}, David M. Miller III^{4,9*}, Marc
8 Hammarlund^{1,3*}

9

10 ¹ Department of Genetics, Yale University School of Medicine, New Haven, CT, USA

11 ² Department of Statistics, Columbia University, New York, NY, USA

12 ³ Department of Neuroscience, Yale University School of Medicine, New Haven, CT, USA

13 ⁴ Department of Cell and Developmental Biology, Vanderbilt University School of Medicine,
14 Nashville, TN, USA

15 ⁵ Department of Biological Sciences, Columbia University, New York, NY, USA

16 ⁶ Neurobiology Department, University of Massachusetts Chan Medical School, Worcester, MA,
17 USA

18 ⁷ Department of Cell, Developmental, and Integrative Biology, University of Alabama
19 Birmingham, AL, USA

20 ⁸ Howard Hughes Medical Institute, Columbia University, New York, NY, USA

21 ⁹ Program in Neuroscience, Vanderbilt University School of Medicine, Nashville, TN, USA

22 * Corresponding Authors

23

24 Marc Hammarlund: marc.hammarlund@yale.edu

25 David M. Miller, III: david.miller@vanderbilt.edu

26 Oliver Hobert: or38@columbia.edu

27 **Abstract**

28 Neuron-specific morphology and function are fundamentally tied to differences in
29 gene expression across the nervous system. We previously generated a single
30 cell RNA-seq dataset for every anatomical neuron class in the *C. elegans*
31 hermaphrodite. Here we present a complementary set of bulk RNA-seq samples
32 for 41 of the 118 neuron classes in *C. elegans*. We show that the bulk dataset
33 captures both lowly expressed and noncoding RNAs that are missed in the single
34 cell dataset, but also includes false positives due to contamination by other cell
35 types. We present an integrated analytical strategy that effectively resolves both
36 the low sensitivity of single cell RNA-seq data and the reduced specificity of bulk
37 RNA-Seq. We show that this integrated dataset enhances the sensitivity and
38 accuracy of transcript detection and quantification of differentially expressed
39 genes. We propose that our approach provides a new tool for interrogating gene
40 expression, by bridging the gap between old (bulk) and new (single cell)
41 methodologies for transcriptomic studies. We suggest that these datasets will
42 advance the goal of delineating the mechanisms that define neuronal
43 morphology and connectivity in *C. elegans*.

44 **Introduction**

45 Neurons exhibit an extraordinary range of morphological forms and
46 physiological functions. Because this diversity is largely driven by underlying
47 differences in gene expression, a key goal of neuroscience is to identify the
48 transcripts expressed in each neuron type.

49 To date, *C. elegans* is the only organism for which goal has been achieved;
50 a gene expression map of the entire nervous system at the resolution of single
51 neuron types. The adult *C. elegans* hermaphrodite contains 302 neurons divided
52 into 118 anatomically distinct neuron types. The structure, connectivity, and
53 lineage are known for each of these neurons (Brittin et al., 2021; Cook et al.,
54 2019; Moyle et al., 2021; Sulston and Horvitz, 1977; Sulston et al., 1983; White
55 et al., 1986). Recently, the *C. elegans* Neuronal Gene Expression Map &
56 Network project (CeNGEN) (Hammarlund et al., 2018) used single cell RNA
57 sequencing (scRNA-seq) technology to generate a gene expression atlas that
58 matches the single neuron resolution of the structural map of the mature *C.*
59 *elegans* nervous system (Taylor et al., 2021).

60 The CeNGEN scRNA-seq dataset was acquired with 10x Genomics
61 technology and is largely comprised of reads from poly-adenylated transcripts.
62 Thus, major classes of non-poly-adenylated transcripts, noncoding RNAs in
63 particular, are poorly represented in the CeNGEN scRNA-seq data. In addition,
64 low abundance transcripts may be under-represented in scRNA-seq data,
65 particularly in clusters with relatively few cells (Taylor *et al.*, 2021). Both
66 noncoding RNAs and low abundance transcripts are potentially important
67 mediators of neuronal fate. A description of their expression is therefore needed
68 to complement the CeNGEN scRNA-seq map of neuronal poly-adenylated
69 transcripts.

70 Here, we use FACS to isolate single neuron types for bulk RNA
71 sequencing with the goal of describing neuronal gene expression with high
72 sensitivity and specificity. We generated profiles for 41 individual neuron types
73 from the mature *C. elegans* hermaphrodite nervous system. This data set
74 samples a wide range of neuron types including motor neurons, interneurons,

75 and sensory neurons. We built sequencing libraries with random primers for
76 robust detection of both poly-adenylated and non-coding RNAs (Barrett et al.,
77 2021). Importantly, we developed a novel computational approach to integrate
78 the bulk dataset with the existing CeNGEN scRNA-seq dataset. Our new
79 analytical strategy enhanced the accuracy and sensitivity of both data sets for
80 profiles of each neuron type. The resultant integrated data set refines quantitative
81 measures of gene expression and improves accuracy of differential expression
82 calling between neuron types. These data provide a unique opportunity for future
83 studies that link gene expression to neuron function, structure, and connectivity.

84 **Methods**

85

86 **Strains**

87 Strains used for FACS isolation of individual neuron classes are listed in
88 Supplementary Table S1.

89

90 **FACS isolation for RNA-seq**

91 Labeled neuron types were isolated for RNA-seq as previously described
92 (Spencer et al., 2014; Taylor et al., 2021). Briefly, synchronized populations of L4
93 stage larvae were dissociated and labeled neuron types isolated by
94 Fluorescence Activated Cell Sorting (FACS) on a BD FACSAria III equipped with
95 a 70-micron diameter nozzle. DAPI was added to the sample (final concentration
96 of 1 mg/mL) to label dead and dying cells. For bulk RNA-sequencing of individual
97 cell types, sorted cells were collected directly into TRIzol LS. At ~15-minute
98 intervals during the sort, the sort was paused, and the collection tube with TRIzol
99 was inverted 3-4 times to ensure mixing. Cells in TRIzol LS were stored at -80C
100 for RNA extractions (see below).

101

102 **RNA extraction**

103 RNA extractions were performed as previously described (Taylor *et al.*,
104 2021). Briefly, cell suspensions in TRIzol LS (stored at -80°C) were thawed at
105 room temperature. Chloroform extraction was performed using Phase Lock Gel-
106 Heavy tubes (Quantabio) according to the manufacturer's protocol. The aqueous
107 layer from the chloroform extraction was combined with an equal volume of
108 100% ethanol and transferred to a Zymo-Spin IC column (Zymo Research).
109 Columns were centrifuged for 30 s at 16,000 RCF, washed with 400 mL of Zymo
110 RNA Prep Buffer, and centrifuged for 16,000 RCF for 30 s. Columns were
111 washed twice with Zymo RNA Wash Buffer (700 mL, centrifuged for 30 s,
112 followed by 400 mL, centrifuged for 2 minutes). RNA was eluted by adding 15 mL
113 of DNase/RNase-Free water to the column filter and centrifuging for 30 s. A 2 µL

114 aliquot was submitted for analysis using the Agilent 2100 Bioanalyzer Picochip to
115 estimate yield and RNA integrity, and the remainder was stored at -80°C.

116

117 **Bulk sequencing and mapping**

118 Each bulk RNA sample was processed for sequencing using the SoLo
119 Ovation Ultra-Low Input RNaseq kit from Tecan Genomics according to
120 manufacturer instruction, modified to optimize rRNA depletion for *C. elegans*
121 (Barrett *et al.*, 2021). Libraries were sequenced on the Illumina HiSeq 2500 with
122 150 bp paired end reads. Reads were mapped to the *C. elegans* reference
123 genome from WormBase (version WS281) using STAR (version 2.7.0) with the
124 option `--outFilterMatchNminOverLread 0.3`. Duplicate reads were removed using
125 NuDup (Tecan Genomics, version 2.3.3), and a counts matrix was generated
126 using the featureCounts tool of SubRead (version 1.6.4). FASTQC was used for
127 quality control before alignment, and four samples were removed for failing QC
128 or for a low number of reads.

129

130 **Pseudobulk aggregation of single-cell data**

131 We downloaded CeNGEN scRNA-seq dataset as a Seurat object from the
132 CeNGEN website (www.cengen.org). Cells from the same cell type and
133 biological replicate (e.g. AFD cluster, replicate eat_4) were aggregated together
134 by summation into a single pseudobulk sample if there were more than 10 cells
135 in the single cell-type-replicate. For this work, single cell clusters of neuron
136 subtypes were collapsed to the resolution of the bulk replicates (ex: VB and VB1
137 clusters in the single cell data were treated as one VB cluster).

138

139 **Sample Normalization**

140 Intra-sample normalization (gene length normalization for bulk samples)
141 was performed before integration. Inter-sample normalization (library size
142 normalization) was performed after integration. Library size normalizations were
143 performed using a TMM (trimmed mean of M-values) correction in edgeR
144 (version 3.36.0). TMM Normalizations were performed separately for each

145 integrated matrix. For differential expression (Figure 3), bulk counts were used as
146 input for integration, as edgeR uses unnormalized counts values as the input. For
147 gene detection (Figure 1), bulk sample counts were normalized to gene length
148 prior to integration, as this intra sample normalization shows improved accuracy
149 for calling gene expression (Supplementary Figure 1C).

150

151 **Integrating bulk and pseudobulk samples**

152 We integrated bulk and single cell profiles by randomly pairing bulk
153 samples and pseudobulk replicates for the same cell type, and then taking the
154 geometric mean. A value of 0.1 was added to all pseudobulk data sets to obviate
155 zero values (Equation 1). Our analysis was limited to cell types with at least 2
156 bulk samples and 2 pseudobulk replicates (supplementary table S3).

157

$$\text{Equation 1: } I = \frac{\log(\text{Bulk} + 0.1) + \log(\text{Pseudobulk} + 0.1)}{2}$$

158

159 The random pairing and integration step was performed 50 times. As an
160 example: for AFD, we began with 5 bulk samples, and 3 pseudobulk replicates.
161 For each integration, we randomly selected 3 bulk samples, and paired them with
162 3 pseudobulk replicates. Each pseudobulk replicate was then scaled to match
163 the total counts in the corresponding bulk sample. Each AFD bulk-pseudobulk
164 pairing was integrated by taking the geometric mean (with an added pseudo-
165 count of 0.1), producing 3 integrated samples. This process was repeated 50
166 times, across all cell types, producing 50 separate integrated matrices (genes x
167 integrated-replicates), sampling from all possible bulk-pseudobulk pairings
168 across all cell types.

169

170 **Ground-truth genes**

171 As an independent measure of gene expression, we used a “ground truth”
172 dataset of 160 genes for which expression in individual neuron types is known
173 with high precision across the entire nervous system. These studies used high

174 confidence fosmid fluorescent reporters, CRISPR strains or other methods
175 (Bhattacharya et al., 2019; Harris et al., 2019; Reilly et al., 2020; Stefanakis et
176 al., 2015; Taylor *et al.*, 2021; Yemini et al., 2021).

177 We also curated a list of 445 genes that are exclusively expressed outside
178 the nervous system to assess potential non-neuronal contamination in each
179 sample. This list was curated from published datasets of fluorescent reporters,
180 tissue specific RT-PCR, and transcriptomic studies available on WormBase
181 (Harris *et al.*, 2019). Genes were included if two forms of evidence both
182 suggested expression in the same non-neuronal tissue (non-overlap was allowed
183 so long as at least one tissue was consistent), and there was no evidence
184 available suggesting neuronal expression.

185 Ground truth gene expression is available in supplementary tables S5 &
186 S6.

187

188 **Comparing datasets to ground-truth**

189 When comparing bulk, single cell, and integrated data to “ground truth”
190 gene expression, a static threshold was applied to the average normalized cell
191 profile (arithmetic mean across all cells, or samples). Single cells were
192 normalized to library size prior to averaging to calculate TPM counts (Packer et
193 al., 2019). Bulk samples were normalized using the GeTMM method (Smid et al.,
194 2018), first normalizing to gene length, then to library size using a TMM
195 correction in edgeR (version 3.36.0). Each of the 50 integrated matrices were
196 separately normalized to library size, the average cell profile for each integrant
197 was calculated, then the 50 resultant genes x cell-types matrices were averaged.
198 The area under the curve (AUC) for the Receiver-Operator Characteristic (ROC)
199 and the Precision-Recall (PR) curves were calculated using the auc function with
200 the trapezoid option from the bayestestR package (version 0.11.5).

201

202 **Thresholding lowly expressed genes and noncoding genes**

203 For lowly expressed protein coding genes, and noncoding RNAs, genes
204 were called expressed in a cell type if more than 65% of replicates detect the

205 gene at or above the threshold. For lowly expressed genes, the threshold (73
206 normalized counts) was set to match the FDR (14%) for the published single cell
207 analysis (Taylor *et al.*, 2021). For noncoding RNAs, the threshold was set at 5
208 normalized counts.

209

210 **Proportion estimates**

211 Contamination estimates were performed for each bulk sample by using
212 non-negative least squares (NNLS) modeling on down-sampled and square root
213 transformed counts, averaging across 100 estimates per sample. Down-sampling
214 was performed to reduce bias against neuron types with small cluster sizes. For
215 each sample (ex: AFD replicate 1), proportions were estimated using only
216 neuronal cells for the corresponding single cell cluster (ex: AFD), and identified
217 non-neuronal clusters (Glia, Excretory, Hypodermis, Intestine, Muscle-
218 mesoderm, Pharynx, and Reproductive) For each iteration, all 8 single cell
219 clusters were down sampled to 30 cells each, and average TPM counts were
220 calculated using the arithmetic mean for each gene in the 30 cells. Gene level
221 variance was calculated using the averaged TPM values, and low variance
222 genes were removed. Bulk sample counts and single cell TPMs were square root
223 transformed before the NNLS calculation. NNLS estimates across all 100
224 iterations were averaged for the final estimate. NNLS calculations were
225 performed using the nnls package in R (version 1.4).

226

227 **Correlating gene expression to non-neuronal contaminants**

228 Each gene was correlated to non-neuronal contamination across all
229 samples using Spearman's correlation test. High correlation to any contaminant
230 was used to indicate that the gene is likely detected because of contamination,
231 not expression in the target neuron. For genes passing an expression threshold >
232 2 normalized counts in at least 2 sample, their highest correlation value to any
233 contaminant tissue was collected, and cutoffs were determined by fitting a
234 gaussian mixture model using the normalmixEM2comp function in mixtools
235 (version 1.2.0), fitting 2 gaussian distributions to the distribution of highest

236 contaminant correlations. Cutoffs were selected to exclude 98% of the predicted
237 contaminant distribution.

238

239 **Differential expression and harmonic mean p combination**

240 Differential expression was performed using the quasi-likelihood F-test
241 approach in edgeR (glmQLFit and glmQLFTest functions). Each integrant
242 dataset was fit and tested separately. P-values across integrated tests were
243 treated as dependent, and were combined using the harmonic mean p approach,
244 using the harmonicmeanp package in R (version 3.0) (Wilson, 2019). LogFC
245 values were combined by taking the arithmetic mean across integrated tests.
246 Consensus values were obtained by counting the number of iterations where a
247 gene was called differentially expressed (P-value < 0.05). In the bulk dataset,
248 genes were called differentially expressed if they had a P-value less than 0.05,
249 and an absolute logFC greater than 2. In the integrated dataset, genes were
250 called differentially expressed if they had a consensus value of at least 40 (P-
251 value < 0.05 in 40 out of 50 separate tests), and an absolute average logFC
252 greater than 2.

253 We used edgeR to perform pairwise differential expression analysis on
254 each of the 50 integrated datasets separately, resulting in 50 edgeR comparisons
255 per neuron pair. As these comparisons are not fully independent, we combined
256 p-values across all 50 tests using the harmonic mean p procedure (Wilson,
257 2019). We also generated a consensus value based on how often a gene was
258 called differentially expressed in the individual integrated comparisons ($p < 0.05$).

259

260 **Ground-truth for differential expression**

261 We adapted the binary ground-truth expression matrix to provide a ground
262 truth for continuous differential expression analysis. For all neuron-neuron pairs,
263 we subset the ground-truth genes to genes that are expressed in one of the two
264 cells, and genes expressed in neither cell. We reasoned that genes called
265 expressed in one cell but not the other in the ground truth data should predict
266 differential expression when comparing continuous data from the two neurons.

267 We also reasoned that genes called unexpressed in both cells in the ground truth
268 data should not be called differentially expressed when comparing continuous
269 data. However, genes called expressed in both cell types in the binary ground-
270 truth data are likely to be a mix of genes that are truly differentially expressed (eg
271 low expression vs high expression), and genes that are not differentially
272 expressed. Therefore, genes expressed in both cell types in the binary ground-
273 truth data were excluded from this analysis. These ground-truth sets for
274 differential expression were designed in a directional manner. For example, when
275 comparing OLQ and PVD, we generated two sets of ground truth genes, and a
276 separate TPR, FPR, and FDR are calculated for OLQ and PVD. For OLQ, the
277 true genes are the genes called expressed in OLQ but not PVD in the ground-
278 truth matrix. The false genes are the genes called unexpressed in both neurons
279 *and* the genes called expressed in PVD alone (we expect those genes to be
280 enriched in PVD, and thus if they are called enriched in OLQ they would be
281 labeled false positives). Thus, we first calculate the genes enriched in OLQ, and
282 compare them to what we expect to see enriched in OLQ, and we separately
283 compare genes enriched in PVD to the genes that we expect to see in PVD.

284 Accuracy scores were calculated by adding up all true positive (TP)
285 events, and all true negative (TN) events, and dividing by the total number of
286 ground truth genes used (Equation 2).

287

288 Equation 2: Accuracy =
$$\frac{TP + TN}{TP + TN + FP + FN}$$

289

290 Matthew's Correlation Coefficient (MCC) is a metric for evaluating binary
291 true/false classifications that is robust to imbalanced datasets (Chicco and
292 Jurman, 2020; Jurman et al., 2012; Matthews, 1975) (Equation 3). This is useful
293 for evaluating differential expression performance as the ground truth dataset is
294 heavily biased towards actual false values.

295

296 Equation 3: MCC =
$$\frac{(TP * TN) - (FP * FN)}{\sqrt{(TP + FP) * (TP + FN) * (TN + FN) * (TN + FP)}}$$

297 **Results**

298

299 **Bulk sequencing of individual neuron types**

300 The model organism *C. elegans* is uniquely suitable for the task of defining
301 gene expression in the nervous system at high resolution and genome scale
302 (**Fig. 1**). *C. elegans* is the first metazoan with a completely sequenced genome
303 (Consortium, 1998) and the only animal for which we know every cell division
304 that gives rise to the adult body plan (*i.e.*, “cell lineage”) (Sulston and Horvitz,
305 1977; Sulston *et al.*, 1983), as well as the anatomy of each neuron and all of its
306 connections with other cells (Brittin *et al.*, 2021; Cook *et al.*, 2019; Moyle *et al.*,
307 2021; Varshney *et al.*, 2011; White *et al.*, 1986). The entire *C. elegans*
308 hermaphrodite nervous system contains 302 neurons with 118 anatomically-
309 defined neuron classes, each comprised of relatively few cells, ranging from 1 to
310 13 neurons (White *et al.*, 1986). Most of these neuron classes are either a
311 bilateral pair of anatomically similar cells (70 classes) or single neurons (26
312 classes) with unique morphological and functional characteristics. The rich array
313 of distinct neuron classes in *C. elegans*, combined with the fact that these types
314 are invariant among individuals, means that each neuron class can be analyzed
315 in depth to reveal the genetic programs that define neuronal diversity.

316 We previously generated a gene expression atlas for the entire *C.*
317 *elegans* nervous system at the resolution of single neuron types. We completed
318 this atlas with single-cell techniques by adopting the strategy of using FACS to
319 enrich for specific groups of neurons for a series of scRNA-seq experiments.
320 However, the description of gene expression in this atlas is incomplete (Taylor *et*
321 *al.*, 2021); (1) Lowly-expressed genes, particularly in clusters with few cells, may
322 not be detected; and (2) Non-poly adenylated transcripts are excluded (Taylor *et*
323 *al.*, 2021).

324 To address these limitations and to provide a broader description of gene
325 expression across the nervous system, we used a bulk RNA sequencing strategy
326 to profile different neuron types. We used a series of *C. elegans* strains, each of
327 which uses one or more fluorescent markers to label an individual neuron type

328 for isolation by FACS. 21 individual neurons could be uniquely marked with a
329 single neuron-specific promoter. For an additional 19 neuron types, we used an
330 intersectional strategy involving different colored fluorophores to label each target
331 neuron, and for 1 neuron we collected some samples with one fluorophore, and
332 other samples with the intersection of two fluorescent markers (Table S1). For
333 example, we used *flp-22::GFP* and *unc-47::mCherry* to mark the single neuron
334 AVL (Figure 1A).

335 For each strain, we used FACS to isolate neurons from synchronized
336 populations of hermaphrodites at the L4 stage, by which time all neurons have
337 been born and are terminally differentiated (Sulston and Horvitz, 1977). Labeled
338 cells were collected in TRIzol LS for RNA extraction (Figure 1B). We isolated a
339 wide range of cells (~700 – 90,000) in each sample across neuron types. Multiple
340 biological replicates (e.g., separately grown cultures) were generated for each
341 neuron class. In total, we sequenced 160 samples across 41 neuron types
342 (Figure 1A; Table S2). The 41 neurons that we profiled sample a wide range of
343 anatomical locations (head ganglia, ventral cord, mid-body and tail neurons,
344 pharyngeal neurons) functional modalities (sensory, inter- and motor neurons),
345 neurotransmitter usage (glutamatergic, GABAergic, cholinergic, aminergic) and
346 lineage history (Figure S1A). (A few of these bulk neuron profiles have been
347 previously described, Taylor et al., 2021.)

348 We used a ribodepletion strategy combined with random priming for cDNA
349 synthesis. This approach optimized whole transcript coverage for each gene and
350 also captured non-polyadenylated RNAs (see Methods) (Barrett *et al.*, 2021).
351 The resultant datasets comprise a high-resolution view of RNA expression
352 across the *C. elegans* nervous system. A distribution of neuron-specific data sets
353 for the first two principal components shows separation between sensory
354 neurons (especially ciliated sensory neurons) vs motor/interneurons, a result
355 consistent with patterns observed for scRNA-seq data on the same neuron
356 classes (Figure 1B) (Taylor *et al.*, 2021).

357

358 **A strategy for integrating bulk and single-cell data to improve gene**
359 **detection accuracy**

360 Bulk RNA-seq and scRNA-seq datasets have complementary strengths
361 and weaknesses. Bulk RNA-seq can enhance sequencing depth and gene
362 detection, capture non-polyadenylated transcripts, and result in uniform coverage
363 of the transcript body (Barrett *et al.*, 2021). Bulk RNA-seq data are typically
364 contaminated, however, with transcripts from non-target cell types which can limit
365 specificity for some genes. By contrast, scRNA-seq datasets allows for high
366 specificity in gene detection, as contaminating cells can be identified post-hoc,
367 but can show reduced transcript sensitivity, especially for low abundance cell
368 types (Taylor *et al.*, 2021).

369 Recent studies have exploited the strengths of these complementary
370 approaches, i.e., the depth of bulk RNA sequencing and the specificity afforded
371 by scRNA-seq, for downstream analysis. These approaches primarily focused on
372 the problem of deconvolution, seeking to infer cell-type expression profiles from
373 tissue level bulk samples, using scRNA-seq references as a guide (Newman *et*
374 *al.*, 2019; Wang *et al.*, 2021a; Wang *et al.*, 2021b; Zhu *et al.*, 2018). By contrast,
375 our dataset contains bulk RNA-seq reads for individual cell types isolated by
376 FACS, exactly matching cell types identified as scRNA-seq clusters. Thus, our
377 data present an opportunity to directly integrate bulk and scRNA-seq profiles for
378 individual cell types, with the goal of combining both datasets to increase depth
379 and accuracy.

380 We constructed pseudobulk samples from the scRNA-seq data for the
381 subset of overall neuron types represented in our bulk RNA-seq data set. Each
382 pseudobulk sample was generated by aggregating scRNA-seq data from
383 individual biological replicates for each annotated cell type. For example, for the
384 AFD cluster, we generated 3 pseudobulk samples, each containing cells from a
385 different single cell experiment, with cell numbers ranging from 27 to 141, and
386 total read counts across all genes ranging from 28,781 to 126,778 (Table S3).
387 We adopted the approach of generating separate pseudobulk data sets for
388 scRNA-seq data from independent single cell experiments because biological

389 replicates have been shown to improve the accuracy of differential expression
390 analysis of scRNA-seq datasets (Crowell et al., 2020; Squair et al., 2021;
391 Thurman et al., 2021).

392 For integrating bulk and scRNA-seq data sets, we adopted the
393 straightforward approach of calculating the geometric mean for each transcript of
394 randomly paired bulk and pseudobulk replicates. This pairing was performed
395 across 50 iterations to sample all possible bulk-pseudobulk arrangements and
396 averaged for comparison to the ground truth genes (see Methods for details,
397 Figure 2A).

398

399 **Integrating bulk and single-cell data improves gene detection** 400 **accuracy.**

401 Accurately detecting gene expression (distinguishing between true signal
402 vs noise) is a central goal for RNA-seq experiments. We first set out to assess
403 our bulk datasets by comparison to ground truth genes (see Methods,
404 Supplementary Table S5) (Taylor *et al.*, 2021). We also used published
405 expression data to curate a list of 445 ground truth genes in non-neuronal cells
406 that are likely not expressed in neurons (Supplementary Table S6).

407 For the bulk, scRNA-seq, and integrated datasets, expression calling was
408 performed by setting a single threshold at the average normalized counts values
409 for each cell type. Thus, all genes in all cell types that meet or exceed the
410 threshold are called “expressed”, and all genes in all cells that fall below the
411 threshold are called “unexpressed”. These binary expression values were then
412 compared to the ground-truth datasets for neuronal and non-neuronal cells. This
413 treatment determined that the bulk samples show a high (FPR) (False Positive
414 Rate) versus combined ground truth genes for neuron and non-neuronal cells
415 across all thresholds (Figure 2B-D). These results suggest that bulk data set
416 contains non-neuronal transcripts from a low level of contaminating cells in the
417 FACS preparation. By contrast, the clustering algorithms used to generate the
418 scRNA-seq data (before pseudobulk aggregation) effectively exclude unwanted
419 cell types and thus result in fewer false positives in the scRNA-seq data.

420 Interestingly, at relatively low precision (or high FPR), the bulk data
421 approached a TPR (True Positive Rate) of 100% (Figure 2B-D). By contrast, the
422 scRNA-seq pseudobulk data peak at a 91.9% TPR, suggesting that the single
423 cell dataset fails to detect some genes. Together, this analysis indicates that bulk
424 and single-cell approaches both afford robust approximations of gene
425 expression, but that they have different disadvantages: the bulk approach is
426 prone to contaminating data from other cell types, whereas the single-cell
427 approach is limited in detection.

428 Measured against the neuronal ground truth genes, the integrated dataset
429 shows a similar sensitivity to the bulk data at low thresholds, while matching the
430 scRNA-seq ratio of specificity and sensitivity across most thresholds and
431 improving on the scRNA-seq performance for some thresholds (Figure 2B-C).
432 The scRNA-seq data still outperforms the integrated dataset for non-neuronal
433 ground truth genes, but the integrated dataset performs nearly as well at
434 thresholds above 10 normalized counts (Figure 2D). Together these results show
435 that geometric mean integration of bulk RNA-seq and scRNA-seq datasets
436 combines the strengths of both approaches, providing high sensitivity and high
437 specificity across a wide range of thresholds.

438

439 **Integration of bulk and single-cell data enhances the accuracy of** 440 **differential expression analysis.**

441 To determine the effect of integration on the accuracy of differential
442 expression analysis, we compared differential expression (DE) analysis of our
443 bulk vs integrated data sets. For both cases, we performed DE analysis for all
444 possible pairwise combinations of different neuron types (595 in total). Genes
445 were called differentially expressed in bulk data for p-values < 0.05 and an
446 absolute value log₂ fold-change (logFC) > 2, (i.e. 4-fold enrichment) in either cell
447 type. Genes were scored as differentially expressed in the integrated data if they
448 were called significant in at least 40/50 iterations (consensus ≥ 40) and had an
449 average absolute value logFC > 2 (see Methods).

450 We scored the accuracy of differential expression of bulk and integrated
451 data by comparison to neuronal ground truth data. For each pair of neuron types
452 A and B, the ground truth data give rise to one of four possible outcomes for
453 each gene: (1) expressed in both neurons A and B; (2) not expressed in either
454 neuron; (3) expressed in A only; (4) expressed in B only. We assessed accuracy
455 in a directional fashion, such that we examined separately genes called
456 expressed only in A and genes called expressed only in B. For example, for
457 genes called expressed only in A, true positives are ground truth genes with
458 expression only in A, whereas false positives include ground truth genes with
459 expression only in B, as well as ground truth genes that are not expressed in
460 either cell. (Ground truth genes expressed in both neurons A and B were
461 excluded as they could correspond to genes that are not truly differentially
462 expressed between the two cell types). Non-neuronal ground truth genes were
463 used to calculate a separate FPR.

464 We calculated TPR, FPR, and FDR (False Discovery Rate) values for
465 every pair of neurons in both the bulk and integrated datasets. In addition, we
466 calculated Accuracy scores (total true calls / all possible calls, see Methods), and
467 the Matthew's Correlation Coefficient (MCC) (the Pearson product-moment
468 correlation coefficient of the observed and expected results, see Methods)
469 (Chicco and Jurman, 2020) (Figure S3C-F). These results indicate that the
470 integrated dataset is more accurate overall than the bulk dataset (Figure 3A). In
471 addition, on for each neuron-neuron pair, integration results in more improvement
472 than degradation in differential expression accuracy (mean = 0.026, 95.conf.int \pm
473 0.003) (Figure 3B). A similar relationship was observed for MCC scores.
474 Specifically, the number of comparisons with MCC scores near 0 was lower in
475 the integrated data set (figure 3C), which represents the expected performance
476 of a coin toss (Chicco and Jurman, 2020). The difference in MCC scores for each
477 neuron-neuron pair also showed higher scores in the Integrated dataset (mean =
478 0.089, 95.conf.int \pm 0.010). Together, these analyses indicate that integration
479 improves the accuracy of differential gene expression.

480 Next, we examined whether integration could improve differential
481 expression analysis even when scRNA-seq data are limited. The lowest
482 abundance single cell clusters show reduced gene detection (Taylor *et al.*, 2021),
483 suggesting that they might not perform as well for integration. Of the 41 neuron
484 types for which we performed bulk sequencing, PVD and OLQ were the neuron
485 types with the fewest cells per cluster in the single cell dataset (62 cells and 85
486 cells, respectively). In the bulk data, a majority of the genes expected to be
487 enriched in PVD from the neuronal ground truth dataset are correctly called, but
488 none of the expected OLQ genes are called enriched. For example, the gene
489 *gar-1* is expected to be enriched in OLQ but is instead enriched in PVD in the
490 bulk RNA-seq data. After integration, *gar-2* is called enriched in OLQ, and all but
491 one gene that was enriched in PVD or showed mild enrichment in PVD now
492 show mild enrichment towards OLQ, though only *gar-2* passes both the logFC
493 and significance cutoffs (Figure 3F). Considering all true positive genes for both
494 PVD and OLQ, we see a modest increase in the TPR for this comparison (Figure
495 3G), along with a sharp drop in the FPR for neuronal ground-truth genes (Figure
496 3H), and non-neuronal ground-truth genes (Figure 3I). Similar results were
497 observed for other comparisons (Figure S3H-K, although there are also rare
498 instances in which integration decreased the TPR (Figure S3J). Thus, integration
499 with scRNA-seq data improves the accuracy of differential gene expression in
500 bulk RNA samples, even when scRNA-seq data are limited.

501 Non-neuronal contamination in FACS-isolated neuronal Bulk RNA-seq
502 samples varies between samples, and between cell types (Figure S4A-B). This
503 variance could lead to non-neuronal genes being erroneously called significantly
504 enriched in some neuron-neuron comparisons. Most neuron-neuron comparisons
505 in both the bulk and integrated datasets show low but detectable false positive
506 rates for non-neuronal ground truth genes (Figure S3Gi-ii). In addition, some
507 neuron-neuron comparisons in the bulk dataset show low specificity scores for
508 non-neuronal ground truth genes, suggesting that differences in non-neuronal
509 contamination are influencing differential expression calling (Figure S3Giii). The
510 integrated dataset shows much higher specificity scores for the same neuron-

511 neuron pairs, and modest specificity improvements overall (Figure S3Giv). When
512 comparing I5 and BAG neurons in the bulk analysis, 26.3% of non-neuronal
513 ground truth genes are called enriched in either I5 or BAG. In the integrated
514 analysis, only 4.5% of non-neuronal genes are called enriched in either cell type.
515 We conclude that our analysis of the systematic pairwise differential expression
516 among all cell types shows that integration improves differential expression by
517 reducing false positives, both for genes expressed in the nervous system and
518 non-neuronal genes, while maintaining the overall true positive rate.

519

520

521 **Bulk sequencing powers detection of low-abundance transcripts**

522 scRNA-seq analysis of the *C. elegans* neuronal transcriptome generated a
523 map of protein coding gene expression for a total of 128 transcriptionally distinct
524 neuron types. However, this map contains some false negatives—ground truth
525 genes that are known to be expressed in the neuron type but are not detected in
526 the scRNA-seq data. Two factors that contribute to these dropouts are low gene
527 expression and small cluster size (clusters with few neurons tend to detect fewer
528 genes) (Mereu et al., 2020; Taylor *et al.*, 2021).

529 We tested whether bulk RNA-seq data might provide this missing
530 information. We collected a minimum of 701 cells per bulk sample (Table S1),
531 and sequenced each sample to high depth, suggesting that even low-expressed
532 genes might be represented in bulk data. A comparison of protein coding genes
533 between bulk and single-cell data showed a mean Spearman coefficient of 0.612
534 (95.conf.int \pm 0.027), with a sharp drop off in the Spearman coefficient for the
535 smallest single cell clusters (Figure 4A). (This analysis used all protein-coding
536 genes detected in a minimum of 3 cells in the single cell dataset.) This result
537 matches previous analysis of the scRNA-seq data, which showed that gene
538 detection is reduced for clusters with < 500 cells (Taylor *et al.*, 2021). Together
539 these results indicate that bulk data contain gene expression information that is
540 missing from scRNA-seq clusters that contain few cells.

541 Although bulk sequencing typically includes lowly-expressed genes, at
542 least some of them may represent false positives derived from non-neuronal
543 tissue contamination. Since these genes are typically not included in the scRNA-
544 seq data, the integration strategy described above does not ameliorate this
545 problem. Previous studies have shown that correlations between gene
546 expression and tissue level proportion estimates can be used to deconvolve the
547 profiles of multiple tissues from one mixed bulk profile (Wang *et al.*, 2021a). We
548 utilized a similar approach to enrich for genes that are truly expressed in our cell
549 types of interest. First, we estimated contamination in each bulk sample using a
550 non-negative least squares regression (NNLS). We used 100 bootstraps to
551 reduce bias against lowly abundant single cell clusters (see Methods,
552 Supplementary Figure S4A-B). We then calculated per-gene Spearman
553 correlations to each contaminant type (e.g., the correlation of *pgl-1* to
554 reproductive cell contamination across all samples). We validate this approach
555 by observing that contaminant correlations for non-neuronal ground-truth genes
556 are higher than the contaminant correlations for all other protein coding genes
557 (Figure S4C). Using the highest correlation per gene, we modeled this data as a
558 mixture of two Gaussian distributions, one distribution of low contamination
559 correlation scores representing truly expressed neuronal genes, and a second
560 distribution of higher contamination correlation scores representing genes likely
561 present due to contamination from non-neuronal tissues. (Figure S4D). Setting a
562 threshold which removes all genes with a contaminant correlation higher than 0.3
563 excludes 98% of the predicted contaminant distribution profile.

564 Using this decontaminated data, we tested our detection of poorly
565 represented genes. We first interrogated the expression of all genes that are
566 detectable in scRNA-seq experiments, by virtue of being called expressed in at
567 least one cell type (by thresholding on the proportion of cells detecting the gene,
568 see Methods). We tested whether our decontaminated bulk data might provide
569 evidence for expression in additional neuron types. Using a minimum normalized
570 count threshold in the bulk data to match the FDR of “threshold 2” from the
571 published single cell analysis (Taylor *et al.*, 2021), we detected 5 to 169 genes

572 per cell type that were missed in that single cell cluster (mean = 36.9, 95.conf.int
573 ± 9.4) (Figure 4B). Plotting the number of newly detected genes against the
574 single cell cluster size reveals that bulk sequencing detects more protein coding
575 genes for cell types with low coverage in the single cell dataset vs cell types with
576 larger numbers of cells in each cluster (Figure 4C). We used GO term
577 enrichment to evaluate genes called expressed in bulk that were missing in the
578 scRNA-seq data. Most cell types show enrichment for neuron-associated terms,
579 chiefly neuropeptide signaling (Figure 4D, S3). Several cell types also show
580 enrichment for synaptic signaling, dendritic morphology, and receptor regulator
581 activity. Thus, we detect genes in the bulk dataset that are missing from some
582 single cell clusters with the greatest improvement biased towards clusters with
583 low coverage in the scRNA-seq dataset.

584 Next, we tested whether decontaminated bulk data might yield expression
585 information about genes that were undetected in the scRNA-seq dataset.
586 Thresholding the scRNA-seq data results in 3,567 protein coding genes that are
587 identified as not expressed in all cell types, including non-neuronal tissues (see
588 Methods). Additionally, 873 protein coding genes were excluded from analysis in
589 the scRNA-seq dataset because they were detected in fewer than 3 of the
590 100,955 cells sequenced. We combined these gene sets to generate a list of
591 4,440 ‘unexpressed’ genes that were not detected in the single cell analysis
592 (Supplementary table S7).

593 To examine expression of these unexpressed genes in the bulk data, we
594 first ‘decontaminated’ the data by removing genes with strong correlations to any
595 contaminants as described above. We used the non-neuronal ground-truth genes
596 to set a minimum normalized counts threshold for calling expression, which was
597 set to a non-neuronal FPR of 0%. Using this threshold on the remaining
598 decontaminated unexpressed genes, we detected between 9 and 150 protein
599 coding genes per cell type (mean = 25.9, 95.conf.int = ± 7.8) (Figure 4E). Using
600 ADL as an example, we performed Tissue Enrichment Analysis on the 150 new
601 genes (Angeles-Albores et al., 2016). The most enriched term is “ADL genes”, as
602 expected, followed by the “amphid sensillum” and “lateral ganglion”, structures

603 that include the ADL neuron (Figure 4F) (Inglis et al., 2007). Thus, these results
604 suggest that our analysis of bulk data reveals truly expressed genes that were
605 not detected by scRNA-seq.

606

607 **Bulk RNA-seq reveals both broadly expressed and neuron-specific** 608 **noncoding RNAs**

609 A significant benefit of our bulk RNA-seq approach is its sensitivity to non-
610 poly-adenylated transcripts, which include many species of non-coding RNA
611 (Barrett *et al.*, 2021). However, we do not have a ground-truth data set of non-
612 coding genes to evaluate accuracy. In addition, most non-coding RNAs are
613 expressed at lower levels than protein coding genes, making it unreasonable to
614 apply a static threshold using the protein coding FDR (Figure S5A). Thus, we
615 opted to apply a uniform threshold for “expressed” genes and selected the
616 criteria of > 5 normalized counts in at least 65% of samples within a cell type. We
617 again used gene level correlation to contamination estimates as a procedure to
618 eliminate genes that were likely detected due to contamination from other tissues
619 in the bulk samples. First, we estimated contamination for each sample using a
620 bootstrapped NNLS regression (see Methods, Supplementary figure S4A-B), and
621 then calculated per-gene Spearman correlations to each contaminant type. We
622 applied a threshold on the gene level correlation to contamination estimates for
623 each sample by fitting a Gaussian mixture model to the maximum correlation
624 score for each gene. We selected a cutoff of 0.23, which excludes 98% of the
625 estimated contamination distribution (Figure 5A). With these thresholds, an
626 average of 603 noncoding RNAs were identified as “expressed” per cell type (95
627 CI \pm 54.5). By RNA type, we detected 23.0 ± 1.7 lincRNAs, 55.6 ± 7.1
628 pseudogenes, 62.6 ± 12.5 tRNAs, 49.3 ± 2.1 snRNAs, 148.9 ± 2.4 snoRNAs, and
629 266.6 ± 39.1 uncategorized ncRNAs per cell type (Figure 5B).

630 Next, we sought to identify noncoding RNAs with broad expression across
631 multiple neuron types. This approach detected 266 non-coding genes that are
632 called expressed in > 90% of neuron classes defined by bulk RNA-seq (Figure
633 5C, D). These broadly expressed noncoding RNAs, include 128 (48%) snoRNAs

634 and 37 (13.9%) snRNAs, both tenfold greater than the expected proportion
635 assuming a random distribution (Fisher's exact test, P-value < 0.01) (Figure
636 S5B). In contrast, pseudogenes and otherwise uncategorized ncRNAs were
637 significantly depleted (P-value < 0.001). These results indicate that snoRNAs and
638 snRNAs are widely expressed, which matches studies showing broad expression
639 of many snoRNAs and snRNAs in other systems (Fafard-Couture et al., 2021;
640 Isakova et al., 2020), and is consistent with their key roles in rRNA processing
641 and splicing (Bratkovič et al., 2019; Valadkhan, 2013; Wassarman and Steitz,
642 1992).

643 We also sought to identify cell-type-specific noncoding RNAs. We
644 calculated tissue specificity scores for each noncoding RNA called expressed in
645 at least one cell type using the Preferential Expression Measure (PEM) score
646 (Huminiecki et al., 2003; Kryuchkova-Mostacci and Robinson-Rechavi, 2016).
647 We called these genes cell-type specific according to three criteria: (1) Called
648 expressed in \geq one cell type (see above); (2) PEM score > 0.65; (3) > 2
649 normalized counts in a maximum of 10/41 cell types. Using these thresholds, we
650 identified 561 cell-type-specific noncoding RNAs (Figure 5E). By RNA type, 347
651 (61.8%) of cell type-specific noncoding RNA genes are uncategorized ncRNAs,
652 186 (33.2%) are pseudogenes, 15 (2.6%) are tRNAs, 8 (1.4%) are lincRNAs, 3
653 (0.5%) are snoRNAs, and 2 (0.3%) are snRNAs (Figure S5C). We observed
654 significant enrichment of pseudogenes, and a subtle but significant depletion of
655 ncRNAs, snoRNAs, and tRNAs (P-value < 0.01). Clustering by genes and cell
656 type modalities revealed clear enrichment for noncoding RNAs in individual
657 neuron types (Figure 5F). The number of specific noncoding RNAs per cell type
658 ranged from 0 (PVC) to 120 (ADL), with a mean of 14 (\pm 8.5) (Supplementary
659 Table S8). These data reveal a wide diversity of noncoding RNA expression
660 across the nervous system and open the door to in depth studies of noncoding
661 RNA contributions to individual neuron function.

662

663 **Discussion**

664 In this work, we present bulk RNA-seq data for 41 neuron classes or about
665 1/3 of all known neuron types in the *C. elegans* nervous system (Figure 1A-B).
666 We describe a new method of integrating these bulk RNA-seq data with
667 previously obtained single-cell RNA-seq data (Taylor et al., 2021) that improves
668 gene detection accuracy for both data sets (Figure 1D-F). Integrated data sets
669 also outperform the original bulk samples in accurately calling differential gene
670 expression across all pairwise comparisons (Figure 3), with a clear reduction in
671 false positives (Figure S3D, G). With the rapid growth of scRNA-seq atlases that
672 complement bulk RNA-seq datasets for individual tissues, our results offer a
673 timely and useful opportunity to improve the accuracy of cell and tissue-specific
674 transcriptional profiles. Furthermore, our computational integration approach is
675 general and can be applied to combine additional sequencing modalities to
676 further incorporate complementary gene expression signals to amplify the depth
677 of sequencing.

678 In addition to enhancing the accuracy of differential gene expression, the
679 integrated bulk RNA-seq dataset detects lowly expressed protein coding genes
680 that were not detected by scRNA-seq (Figure 4B-C,E) and thus could reveal new
681 drivers of neuron-specific traits. Because our library construction methods were
682 designed to capture non-polyadenylated transcripts, our bulk RNA-seq data set
683 detects noncoding RNAs that were not revealed by previous scRNA-seq results
684 (Barrett et al., 2021; Taylor et al., 2021) (Figure 5B). Some of these noncoding
685 RNAs are broadly expressed in the nervous system (Figure 5C-D) which is
686 suggestive of shared functions across different types of neurons. Interestingly, a
687 subset of non-coding RNAs are expressed in a limited number of neuron types
688 (Figure 5E-F) pointing to potentially important roles in determining key neuron-
689 specific functions. In addition, the bulk RNA-seq dataset contains transcript
690 information across the gene body, which might yield information about mRNA
691 splicing that is not found in the scRNA-seq dataset.

692 Overall, our approach achieves a comprehensive representation of all
693 classes of transcripts expressed in individual neuron types. These data can now

694 drive analysis of mechanisms that control gene expression across the genome in
695 individual neuron types, and also support identification of differentially expressed
696 genes that define neuron-type specific differences in morphology and function.
697 Public access to these data (described below) will enable further analysis into the
698 regulation and function of differential gene expression in *C. elegans* neurons.

699

700 **Supplementary Tables**

701 Supplementary tables S1-8 are available on figshare
702 (<https://doi.org/10.6084/m9.figshare.19522096.v1>).

703 **Data Availability**

704 Bulk raw data are in the process of being posted at GEO, and the linking
705 information will be posted to the CeNGEN website when available. Single cell
706 raw data are available at Gene Expression Omnibus (GEO)
707 (<https://www.ncbi.nlm.nih.gov/geo>, GEO: GSE136049). Counts data and
708 additional supporting files can be downloaded from the CeNGEN website
709 (<https://www.cengen.org>) and code is available at GitHub
710 (<https://www.github.com/cengenproject>).

711 **DECLARATION OF INTERESTS**

712 The authors declare no competing interests.

713

714

715

Figure legends

716

717

Figure 1: Single neuron bulk RNA-seq via targeted marker expression

718

and FACS isolation: A) Labeling, tissue dissociation, and FACS-enrichment

719

schemes for capturing individual neuron types. Intersecting *flp-22::GFP* and *unc-*

720

47::mCherry markers uniquely label AVL for isolation by FACS from dissociated

721

L4 stage larval cells. RNA from this pool of AVL-enriched cells was used for bulk

722

RNA sequencing (see Methods). B) PCA plot showing all bulk RNA-seq

723

replicates labeled by cell type and colored according to functional modality;

724

Sensory neurons (blue), motor neurons (green), interneurons (red), and CAN

725

neurons (purple).

726

727

Supplementary Figure 1: Bulk RNA sequencing encompasses a broad

728

range of neuron types and correlates with scRNA-seq results. A) Number of

729

cell types sequenced per functional modality. B) Heatmap of Spearman

730

Correlations between average single cell RNA-seq (row) and Bulk RNA-seq

731

(column) profiles for each neuron type. For each row, correlations were

732

calculated for genes called expressed in that single cell cluster (from single cell

733

thresholding) (Taylor *et al.*, 2021).

734

735

Figure 2: Integrating bulk RNA-seq and scRNA-seq data sets

736

improves gene detection accuracy. A) Individual pseudobulk scRNA-seq

737

replicates and bulk RNA-seq samples from the same neuron type (NSM neuron

738

samples illustrated) are randomly paired and integrated (50X for each neuron

739

type) using the geometric mean (see Methods) to generate 50 integrated

740

matrices (genes x integrated-replicate). The average integrated profile was used

741

to call gene expression. Pairwise neuron-neuron differential expression (edgeR)

742

was performed for each of the 50 integrated matrices which were then combined

743

to generate consensus sets of differentially expressed genes. Bulk RNA-seq

744

datasets are used to identify genes that are not detected in scRNA-seq data,

745 including noncoding RNAs and lowly expressed mRNAs. B) Receiver Operator
746 Characteristic (ROC) curve for bulk, single-cell, and integrated datasets
747 compared to neuronal ground-truth genes. The x-axis shows the False Positive
748 Rate (FPR), and the y-axis shows the true positive rate (TPR). C) Precision-
749 Recall (PR) curve for bulk, single-cell, and integrated datasets compared to
750 neuronal ground-truth genes. The x-axis shows the Precision (1 – False
751 Discovery Rate/FDR), and the y-axis shows the TPR (Recall). D) The non-
752 neuronal FPR across a range of thresholds for bulk, single-cell, and integrated
753 datasets compared to non-neuronal ground-truth genes. The x-axis shows the
754 \log_{10} -transformed threshold used for each point; the y-axis shows the FPR. A
755 pseudocount of 1 was added for the \log_{10} -transformation. Each point represents
756 a static threshold applied to all genes in all samples (e.g., expressed ≥ 10
757 normalized counts); Bulk RNA-seq data (green), scRNA-seq (blue), average
758 integrated data (red).

759

760 **Supplementary Figure 2: Intra-sample normalization improves the**
761 **FPR for non-neuronal genes in bulk RNA-seq samples:** The non-neuronal
762 FPR across a range of thresholds for bulk RNA-seq datasets with different
763 normalizations compared to non-neuronal ground-truth genes. The x-axis shows
764 the \log_{10} transformed threshold for each point, the y-axis shows the FPR. Each
765 point represents a static threshold applied to all genes in all samples (e.g.,
766 expressed ≥ 10 normalized counts). Bulk data with only inter-sample
767 normalization using TMM factors (trimmed mean of M-values, used by edgeR)
768 (green) vs bulk data with both intra-sample and inter-sample normalization
769 (GeTMM) (red). AUC = Area Under Curve. A pseudocount of 1 was added for the
770 \log_{10} -transformation.

771

772 **Figure 3: Integrated samples show improved accuracy in detecting**
773 **differentially expressed genes.** A) Density histograms of the accuracy score for
774 all pairwise differential expression comparisons in bulk RNA-seq (blue) vs
775 integrated (orange) datasets. B) Density histogram of the difference (integrated

776 minus bulk) in the accuracy score for each pairwise differential expression
777 comparison, vertical dashed line at 0 represents no difference between the
778 datasets. C) Density histograms of the Matthew's Correlation Coefficient (MCC)
779 score for all pairwise differential expression comparisons in the bulk RNA-seq
780 (blue) vs integrated (orange) datasets. D) Density histogram of the difference
781 (integrated minus bulk) in the MCC score for each pairwise differential
782 expression comparison, vertical dashed line at 0 represents no difference
783 between the datasets. E) Volcano plot for the differential expression profile of the
784 bulk RNA-seq PVD samples vs OLQ samples. Dots represent individual genes.
785 X-axis is log₂ fold change (logFC), and the Y-axis is -log₁₀(P-value). Grey dots
786 are genes that are not called significant, and black dots are genes that pass
787 significance thresholds (P-value < 0.05, and |logFC| > 2, red lines). F) Volcano
788 plot for the differential expression profile of the Integrated PVD samples vs OLQ
789 samples. X-axis is the log₂ fold change (logFC), and the Y-axis is the -
790 log₁₀(harmonic mean P value) (p.hmp). Grey dots are genes that are not called
791 significant, and black dots are genes that pass significance thresholds (P-value <
792 0.05 in ≥80% of edgeR runs across all 50 integrations, and |logFC| > 2).
793 Magenta squares mark genes expected to be enriched in PVD from the neuronal
794 ground-truth dataset, and orange triangles denote genes expected to be enriched
795 in OLQ. *gar-1* and *gar-2* are expected to be enriched in OLQ. G) Bar plot
796 showing the differential expression True Positive Rate (TPR) for genes expected
797 to be expressed in OLQ or PVD but not both. H) Bar plot showing the differential
798 expression false positive rate (FPR) for genes expected to be expressed in
799 neither OLQ nor PVD, and genes that were called enriched in the wrong neuron
800 type. I) Bar plot showing the differential expression FPR for genes expected to be
801 expressed only in non-neuronal tissues.

802

803 **Supplementary Figure 3** A) Table showing an example ground-truth
804 matrix for OLQ and PVD neurons. Here we expect Gene a to be differentially
805 enriched in OLQ over PVD, so it would be considered a positive ground-truth for
806 OLQ and would be used to calculate the TPR for OLQ vs PVD. All other genes

807 shown are expected not to be enriched in OLQ and would thus be used as
808 negative ground-truth, to calculate the FPR and FDR for OLQ vs PVD. When
809 calculating the ground-truth for PVD vs OLQ, we expect Gene b to be enriched in
810 PVD, and so it is treated as a positive ground-truth gene, and all other genes
811 shown are treated as negative ground-truth. B) Example heatmap showing the
812 MCC score for directional OLQ and PVD differential expression. In the OLQ row,
813 we use edgeR to compare genes enriched in OLQ vs expected enrichment using
814 the ground truth data. In the PVD row, we perform the same function, looking
815 instead for enrichment in PVD. Thus, we have 595 neuron-neuron comparisons,
816 with two entries for each pair. For OLQ vs PVD, we have an OLQ entry showing
817 the scores for genes enriched in OLQ, and a PVD entry showing the scores for
818 genes enriched in PVD. C-G) Heatmaps and density plots, showing scores for
819 differential expression compared to neuronal ground-truth genes (C-F) and non-
820 neuronal genes (G) across all neuron types. C) Recall, D) Specificity (1-FPR), E)
821 Accuracy, F) MCC score, and G) non-neuronal specificity. i) Heatmap of the
822 score for the Bulk samples. ii) Heatmap of the score for the Integrated samples.
823 iii) Heatmap of the difference in the scores (Integrated minus Bulk). iv) Density
824 plot for the difference in the scores, black line at 0 indicates no difference
825 between integrated and bulk comparisons. H-J) Bar plots showing neuronal
826 ground-truth TPR and FPR, and the non-neuronal FPR, for four pairs of neurons.
827 All bar graphs of TPR and FPR are shown for both directions of the comparison.
828

829 **Figure 4: Bulk RNA-seq samples detect protein coding genes that are**
830 **not detected in scRNA-seq clusters:** A) Scatter plot showing the relationship
831 between the size of a scRNA-seq cluster (i.e., the number of cells in the cluster)
832 and the Spearman correlation between the average bulk RNA-seq profile and the
833 average scRNA-seq for all protein coding genes. Each dot represents one cell
834 type. Red dashed line shows a Michaelis-Menten fit (see Methods), $g_{max} =$
835 0.675 , $\beta = 29.507$. Blue dashed lines show the 97.5% confidence interval of
836 the fit. B) Bar plot showing the number of protein coding genes detected per cell
837 type in the bulk dataset. Genes plotted are: 1) called unexpressed in the

838 corresponding single cell cluster; 2) have a maximum correlation to any
839 contaminant tissue less than 0.3; and 3) are expressed above 73 normalized
840 counts in the average bulk profile for that cell type. C) Scatter plot showing the
841 relationship between the size of a scRNA-seq cluster and the number of
842 additional protein coding genes detected per cell type (as defined in panel B).
843 Each dot represents one cell type. Red dashed line shows an exponential decay
844 fit (see Methods), $M = 140.2$, $m = 26.5$, $\alpha = 89.1$. Blue dashed lines show the
845 97.5% confidence interval of the fit. D) GO enrichment analysis for protein coding
846 genes detected in bulk IL1 samples that were not detected in the IL1 scRNA-seq
847 cluster. GO enrichment performed using WormBase. E) Bar plot showing the
848 number of protein coding genes detected per cell type in the bulk dataset.
849 Restricted to genes that are never called expressed in any scRNA-seq cluster,
850 have a contaminant correlation less than 0.3, and are expressed above 16
851 normalized counts (determined by setting the non-neuronal FPR threshold to 0).
852 F) Tissue enrichment analysis for protein coding genes detected in the ADL bulk
853 samples but never called expressed in any scRNA-seq cluster (Angeles-Albores
854 *et al.*, 2016).

855

856 **Supplementary Figure 4** A-B) Scatter plots with a linear fit showing the
857 relationship between the \log_{10} transformed single cell cluster size and the
858 estimated neuronal proportion of each bulk sample. Estimates were made using
859 an NNLS regression (non-negative least squares, see Methods). A) Estimates
860 with all single cells in each cluster. Neuronal proportion = $0.081 * \log_{10}(\text{sc_size})$
861 $+ 0.149$. $R^2 = 0.05489$, $p = 0.001666$ B). Estimates taken from the average
862 Neuronal proportion estimate across 100 bootstraps, down-sampled to 30 cells
863 for all clusters before each bootstrap. Neuronal proportion = $0.029 * \log_{10}(\text{sc_size})$
864 $+ 0.268$. $R^2 = 0.003752$, $p = 0.2079$. C) Density plot of the gene
865 level correlation to contaminant estimates. Only the highest correlation per gene
866 is used. Distribution for all protein coding genes (red) vs distribution for non-
867 neuronal ground-truth protein coding genes (blue). D) Density plot of the gene
868 level correlation to contaminant estimates for all genes that are detected in single

869 cell but called unexpressed in one of the 41 cell types covered by bulk
870 sequencing. Only the highest correlation per gene is used. Blue and black
871 dashed lines represent a Gaussian mixture model, used to threshold against
872 contaminant genes. Red line at 0.3 indicates the cutoff, all protein coding genes
873 with a maximum correlation above 0.3 were removed from analysis. E-H) GO
874 term enrichment plots for genes called expressed in the bulk dataset which were
875 called unexpressed in the corresponding scRNA-seq cluster for neurons OLL (E),
876 RIS (F), PVD (G) and PVM (H).

877

878 **Figure 5: Bulk analysis reveals noncoding RNA expression pattern:** A)

879 Density plot showing the distribution of gene level correlation to contaminant
880 estimates (purple), values plotted are the highest correlation per gene. Genes
881 plotted were called expressed in at least one cell type. Blue and black dashed
882 lines represent a gaussian mixture model, used to threshold against contaminant
883 genes. All noncoding genes with a maximum correlation above 0.22 (vertical red
884 line) were removed from analysis. B) Stacked bar graph showing the number of
885 noncoding RNAs called expressed in each neuron type. Colors represent RNA
886 classes. Genes were called expressed in a cell if they were detected above 5
887 normalized counts in greater than 65% of samples for that cell. C) Bar plot
888 showing number of cell types in which each noncoding RNA is detected. The x
889 axis shows the number of cells, and the y axis shows the number of genes
890 detected in that many cells. Genes to the right of the red line are called
891 expressed in more than 90% of the sequenced cell types. D) Heatmap of log
892 transformed GeTMM values of the pan-neuronal genes identified in panel C,
893 columns are annotated by neuron modality. E) Histogram showing the
894 distribution of Preferential Enrichment Measure (PEM) scores per gene, a metric
895 for cell type specificity. Genes are considered cell type specific if they have a
896 PEM greater than 0.65 (red line) and are expressed above 2 normalized counts
897 in fewer than 10 cell types. F) Heatmap of average normalized counts per cell
898 type, for genes considered cell type specific, columns are annotated by neuron
899 modality, and rows are grouped by RNA class.

900

901 **Supplementary Figure 5** A) Density plot showing the relative expression
902 of noncoding RNAs (purple) and protein coding RNAs (orange), x-axis is
903 maximum normalized counts per gene. B) Pie chart showing proportions of
904 classes of pan-neuronal noncoding RNAs. C) Pie chart showing proportions of
905 cell type specific noncoding RNAs. D) Box plot showing the number of cell type
906 specific noncoding RNAs per cell type, grouped by neuron modality.

907

908 **Supplementary Table S1:** All cell types sorted for bulk RNA-seq
909 experiments, with the strain names and allele information.

910 **Supplementary Table S2:** Replicate metadata for bulk RNA-seq
911 experiments, with replicate names, strain names, and the number of cells
912 collected.

913 **Supplementary Table S3:** All cell types used for integrating bulk and
914 single cell RNA-seq data, with the number of replicates in the bulk and single cell
915 datasets for each cell type.

916 **Supplementary Table S4:** Metadata for each single cell replicate,
917 including the replicate name, the total UMI counts in the replicate, the number of
918 individual cells included in the replicate, the cell type, and the experimental
919 replicate name.

920 **Supplementary Table S5:** Ground Truth expression for 160 genes in the
921 *C. elegans* nervous system using fosmid and CRISPR/Cas reporter lines (see
922 methods).

923 **Supplementary Table S6:** Ground Truth expression for 445 genes that are
924 expressed exclusively outside the *C. elegans* nervous system, curated from
925 published data (see methods).

926 **Supplementary Table S7:** Genes called unexpressed in all single cell
927 clusters.

928 **Supplementary Table S8:** An annotated heatmap of highly specific
929 noncoding RNA genes and their log₁₀ transformed expression values in each
930 cell type.

931

References

932

- 933 Angeles-Albores, D., N. Lee, R.Y., Chan, J., and Sternberg, P.W. (2016). Tissue enrichment
934 analysis for *C. elegans* genomics. *BMC Bioinformatics* 17, 366. [10.1186/s12859-016-](https://doi.org/10.1186/s12859-016-1229-9)
935 [1229-9](https://doi.org/10.1186/s12859-016-1229-9).
- 936 Barrett, A., McWhirter, R., Taylor, S.R., Weinreb, A., Miller, D.M., III, and Hammarlund,
937 M. (2021). A head-to-head comparison of ribodepletion and polyA selection approaches
938 for *Caenorhabditis elegans* low input RNA-sequencing libraries. *G3*
939 *Genes|Genomes|Genetics* 11. [10.1093/g3journal/jkab121](https://doi.org/10.1093/g3journal/jkab121).
- 940 Bhattacharya, A., Aghayeva, U., Berghoff, E.G., and Hobert, O. (2019). Plasticity of the
941 Electrical Connectome of *C. elegans*. *Cell* 176, 1174-1189.e1116.
942 [10.1016/j.cell.2018.12.024](https://doi.org/10.1016/j.cell.2018.12.024).
- 943 Bratkovič, T., Božič, J., and Rogelj, B. (2019). Functional diversity of small nucleolar
944 RNAs. *Nucleic Acids Research* 48, 1627-1651. [10.1093/nar/gkz1140](https://doi.org/10.1093/nar/gkz1140).
- 945 Brittin, C.A., Cook, S.J., Hall, D.H., Emmons, S.W., and Cohen, N. (2021). A multi-scale
946 brain map derived from whole-brain volumetric reconstructions. *Nature* 591, 105-110.
947 [10.1038/s41586-021-03284-x](https://doi.org/10.1038/s41586-021-03284-x).
- 948 Chicco, D., and Jurman, G. (2020). The advantages of the Matthews correlation
949 coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC*
950 *Genomics* 21, 6. [10.1186/s12864-019-6413-7](https://doi.org/10.1186/s12864-019-6413-7).
- 951 Consortium, T.C.e.S. (1998). Genome sequence of the nematode *C. elegans*: a platform
952 for investigating biology. *Science* 282, 2012-2018. [10.1126/science.282.5396.2012](https://doi.org/10.1126/science.282.5396.2012).
- 953 Cook, S.J., Jarrell, T.A., Brittin, C.A., Wang, Y., Bloniarz, A.E., Yakovlev, M.A., Nguyen,
954 K.C.Q., Tang, L.T.H., Bayer, E.A., Duerr, J.S., et al. (2019). Whole-animal connectomes of
955 both *Caenorhabditis elegans* sexes. *Nature* 571, 63-71. [10.1038/s41586-019-1352-7](https://doi.org/10.1038/s41586-019-1352-7).
- 956 Crowell, H.L., Soneson, C., Germain, P.-L., Calini, D., Collin, L., Raposo, C., Malhotra, D.,
957 and Robinson, M.D. (2020). muscat detects subpopulation-specific state transitions from
958 multi-sample multi-condition single-cell transcriptomics data. *Nature Communications*
959 11, 6077. [10.1038/s41467-020-19894-4](https://doi.org/10.1038/s41467-020-19894-4).
- 960 Fafard-Couture, É., Bergeron, D., Couture, S., Abou-Elala, S., and Scott, M.S. (2021).
961 Annotation of snoRNA abundance across human tissues reveals complex snoRNA-host
962 gene relationships. *Genome Biology* 22, 172. [10.1186/s13059-021-02391-2](https://doi.org/10.1186/s13059-021-02391-2).
- 963 Hammarlund, M., Hobert, O., Miller, D.M., 3rd, and Sestan, N. (2018). The CeNGEN
964 Project: The Complete Gene Expression Map of an Entire Nervous System. *Neuron* 99,
965 430-433. [10.1016/j.neuron.2018.07.042](https://doi.org/10.1016/j.neuron.2018.07.042).
- 966 Harris, T.W., Arnaboldi, V., Cain, S., Chan, J., Chen, W.J., Cho, J., Davis, P., Gao, S., Grove,
967 C.A., Kishore, R., et al. (2019). WormBase: a modern Model Organism Information
968 Resource. *Nucleic Acids Research* 48, D762-D767. [10.1093/nar/gkz920](https://doi.org/10.1093/nar/gkz920).
- 969 Huminiecki, L., Lloyd, A.T., and Wolfe, K.H. (2003). Congruence of tissue expression
970 profiles from Gene Expression Atlas, SAGEmap and TissueInfo databases. *BMC Genomics*
971 4, 31. [10.1186/1471-2164-4-31](https://doi.org/10.1186/1471-2164-4-31).
- 972 Inglis, P.N., Ou, G., Leroux, M.R., and Scholey, J.M. (2007). The sensory cilia of
973 *Caenorhabditis elegans*. *WormBook*, 1-22. [10.1895/wormbook.1.126.2](https://doi.org/10.1895/wormbook.1.126.2).

- 974 Isakova, A., Fehlmann, T., Keller, A., and Quake, S.R. (2020). A mouse tissue atlas of
975 small noncoding RNA. *Proceedings of the National Academy of Sciences* *117*, 25634-
976 25645. [10.1073/pnas.2002277117](https://doi.org/10.1073/pnas.2002277117).
- 977 Jurman, G., Riccadonna, S., and Furlanello, C. (2012). A Comparison of MCC and CEN
978 Error Measures in Multi-Class Prediction. *PLOS ONE* *7*, e41882.
979 [10.1371/journal.pone.0041882](https://doi.org/10.1371/journal.pone.0041882).
- 980 Kryuchkova-Mostacci, N., and Robinson-Rechavi, M. (2016). A benchmark of gene
981 expression tissue-specificity metrics. *Briefings in Bioinformatics* *18*, 205-214.
982 [10.1093/bib/bbw008](https://doi.org/10.1093/bib/bbw008).
- 983 Matthews, B.W. (1975). Comparison of the predicted and observed secondary structure
984 of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) - Protein Structure* *405*, 442-
985 451. [https://doi.org/10.1016/0005-2795\(75\)90109-9](https://doi.org/10.1016/0005-2795(75)90109-9).
- 986 Mereu, E., Lafzi, A., Moutinho, C., Ziegenhain, C., McCarthy, D.J., Álvarez-Varela, A.,
987 Batlle, E., Sagar, Grün, D., Lau, J.K., et al. (2020). Benchmarking single-cell RNA-
988 sequencing protocols for cell atlas projects. *Nature Biotechnology* *38*, 747-755.
989 [10.1038/s41587-020-0469-4](https://doi.org/10.1038/s41587-020-0469-4).
- 990 Moyle, M.W., Barnes, K.M., Kuchroo, M., Gonopolskiy, A., Duncan, L.H., Sengupta, T.,
991 Shao, L., Guo, M., Santella, A., Christensen, R., et al. (2021). Structural and
992 developmental principles of neuropil assembly in *C. elegans*. *Nature* *591*, 99-104.
993 [10.1038/s41586-020-03169-5](https://doi.org/10.1038/s41586-020-03169-5).
- 994 Newman, A.M., Steen, C.B., Liu, C.L., Gentles, A.J., Chaudhuri, A.A., Scherer, F.,
995 Khodadoust, M.S., Esfahani, M.S., Luca, B.A., Steiner, D., et al. (2019). Determining cell
996 type abundance and expression from bulk tissues with digital cytometry. *Nature*
997 *Biotechnology* *37*, 773-782. [10.1038/s41587-019-0114-2](https://doi.org/10.1038/s41587-019-0114-2).
- 998 Packer, J.S., Zhu, Q., Huynh, C., Sivaramakrishnan, P., Preston, E., Dueck, H., Stefanik, D.,
999 Tan, K., Trapnell, C., Kim, J., et al. (2019). A lineage-resolved molecular atlas of *C.*
1000 *elegans* embryogenesis at single-cell resolution. *Science* *365*. [10.1126/science.aax1971](https://doi.org/10.1126/science.aax1971).
- 1001 Reilly, M.B., Cros, C., Varol, E., Yemini, E., and Hobert, O. (2020). Unique homeobox
1002 codes delineate all the neuron classes of *C. elegans*. *Nature* *584*, 595-601.
1003 [10.1038/s41586-020-2618-9](https://doi.org/10.1038/s41586-020-2618-9).
- 1004 Smid, M., Coebergh van den Braak, R.R.J., van de Werken, H.J.G., van Riet, J., van Galen,
1005 A., de Weerd, V., van der Vlugt-Daane, M., Bril, S.I., Lalmahomed, Z.S., Kloosterman,
1006 W.P., et al. (2018). Gene length corrected trimmed mean of M-values (GeTMM)
1007 processing of RNA-seq data performs similarly in intersample analyses while improving
1008 intrasample comparisons. *BMC Bioinformatics* *19*, 236. [10.1186/s12859-018-2246-7](https://doi.org/10.1186/s12859-018-2246-7).
- 1009 Squair, J.W., Gautier, M., Kathe, C., Anderson, M.A., James, N.D., Hutson, T.H., Hudelle,
1010 R., Qaiser, T., Matson, K.J.E., Barraud, Q., et al. (2021). Confronting false discoveries in
1011 single-cell differential expression. *Nature Communications* *12*, 5692. [10.1038/s41467-
1012 021-25960-2](https://doi.org/10.1038/s41467-021-25960-2).
- 1013 Stefanakis, N., Carrera, I., and Hobert, O. (2015). Regulatory Logic of Pan-Neuronal Gene
1014 Expression in *C. elegans*. *Neuron* *87*, 733-750. [10.1016/j.neuron.2015.07.031](https://doi.org/10.1016/j.neuron.2015.07.031).
- 1015 Sulston, J.E., and Horvitz, H.R. (1977). Post-embryonic cell lineages of the nematode,
1016 *Caenorhabditis elegans*. *Dev Biol* *56*, 110-156. [10.1016/0012-1606\(77\)90158-0](https://doi.org/10.1016/0012-1606(77)90158-0).

1017 Sulston, J.E., Schierenberg, E., White, J.G., and Thomson, J.N. (1983). The embryonic cell
1018 lineage of the nematode *Caenorhabditis elegans*. *Dev Biol* *100*, 64-119. [10.1016/0012-](https://doi.org/10.1016/0012-1606(83)90201-4)
1019 [1606\(83\)90201-4](https://doi.org/10.1016/0012-1606(83)90201-4).

1020 Taylor, S.R., Santpere, G., Weinreb, A., Barrett, A., Reilly, M.B., Xu, C., Varol, E.,
1021 Oikonomou, P., Glenwinkel, L., McWhirter, R., et al. (2021). Molecular topography of an
1022 entire nervous system. *Cell* *184*, 4329-4347.e4323.
1023 <https://doi.org/10.1016/j.cell.2021.06.023>.

1024 Thurman, A.L., Ratcliff, J.A., Chimenti, M.S., and Pezzulo, A.A. (2021). Differential gene
1025 expression analysis for multi-subject single-cell RNA-sequencing studies with
1026 aggregateBioVar. *Bioinformatics* *37*, 3243-3251. [10.1093/bioinformatics/btab337](https://doi.org/10.1093/bioinformatics/btab337).

1027 Valadkhan, S. (2013). Chapter Six - The Role of snRNAs in Spliceosomal Catalysis. In
1028 *Progress in Molecular Biology and Translational Science*, G.A. Soukup, ed. (Academic
1029 Press), pp. 195-228. <https://doi.org/10.1016/B978-0-12-381286-5.00006-8>.

1030 Varshney, L.R., Chen, B.L., Paniagua, E., Hall, D.H., and Chklovskii, D.B. (2011). Structural
1031 properties of the *Caenorhabditis elegans* neuronal network. *PLoS Comput Biol* *7*,
1032 e1001066. [10.1371/journal.pcbi.1001066](https://doi.org/10.1371/journal.pcbi.1001066).

1033 Wang, J., Roeder, K., and Devlin, B. (2021a). Bayesian estimation of cell type-specific
1034 gene expression with prior derived from single-cell data. *Genome Research* *31*, 1807-
1035 1818. [10.1101/gr.268722.120](https://doi.org/10.1101/gr.268722.120).

1036 Wang, W., Yao, J., Wang, Y., Zhang, C., Tao, W., Zou, J., and Ni, T. (2021b). Improved
1037 estimation of cell type-specific gene expression through deconvolution of bulk tissues
1038 with matrix completion. *bioRxiv*, 2021.2006.2030.450493. [10.1101/2021.06.30.450493](https://doi.org/10.1101/2021.06.30.450493).

1039 Wassarman, D.A., and Steitz, J.A. (1992). INTERACTIONS OF SMALL NUCLEAR RNAS
1040 WITH PRECURSOR MESSENGER-RNA DURING INVITRO SPLICING. *Science* *257*, 1918-
1041 1925. [10.1126/science.1411506](https://doi.org/10.1126/science.1411506).

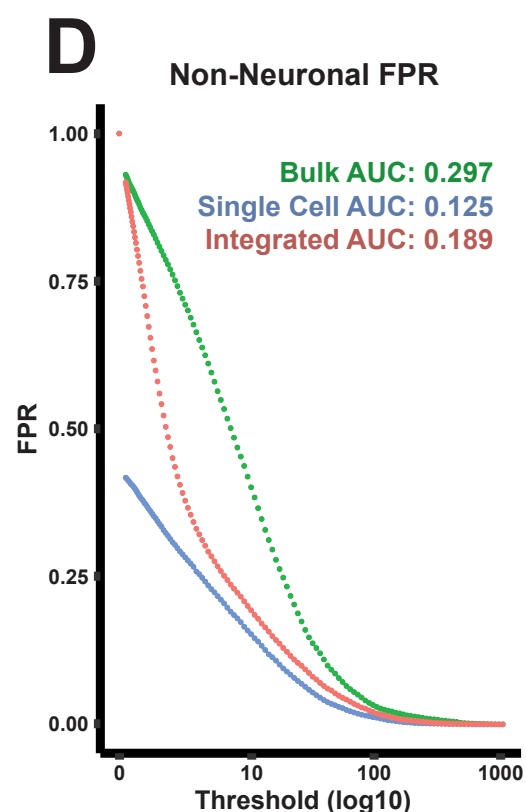
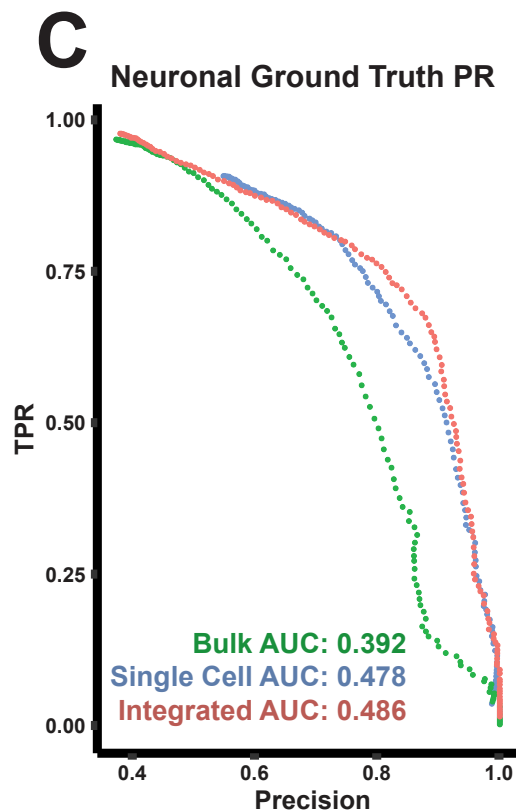
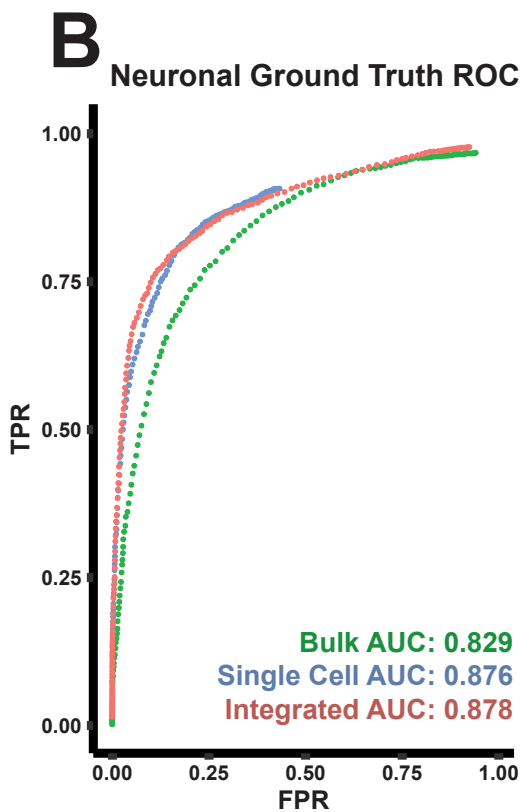
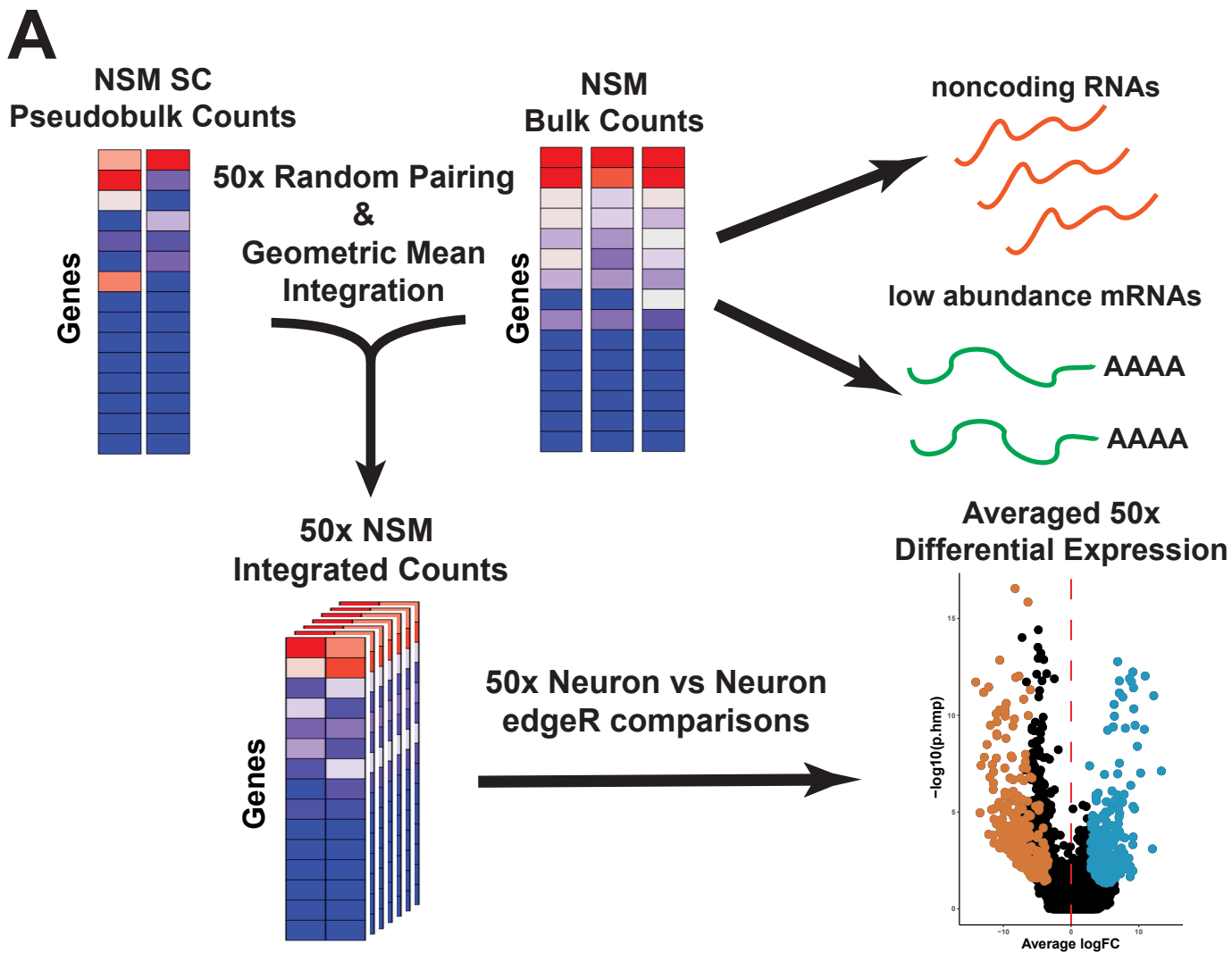
1042 White, J.G., Southgate, E., Thomson, J.N., and Brenner, S. (1986). The structure of the
1043 nervous system of the nematode *Caenorhabditis elegans*. *Philos Trans R Soc Lond B Biol*
1044 *Sci* *314*, 1-340. [10.1098/rstb.1986.0056](https://doi.org/10.1098/rstb.1986.0056).

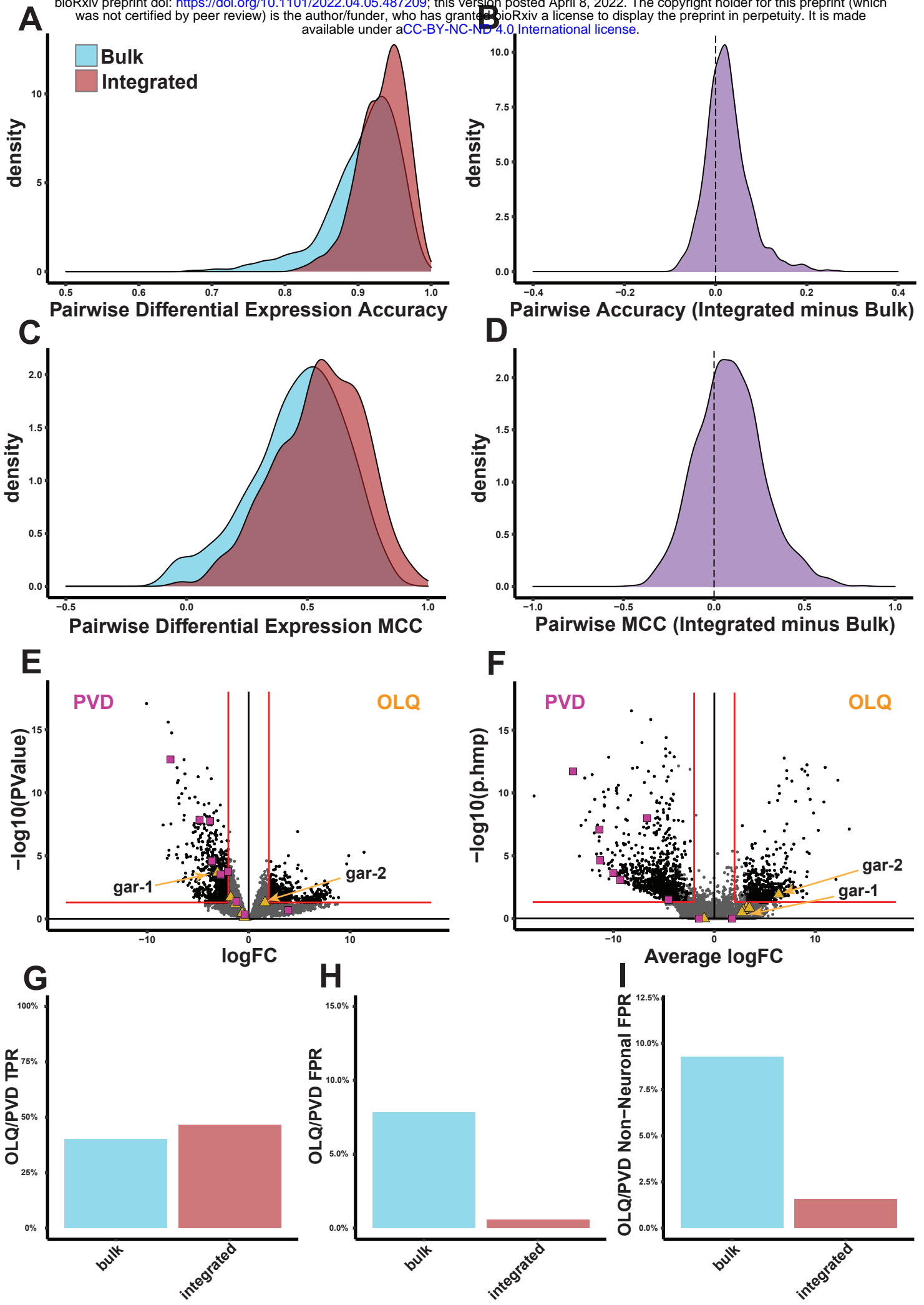
1045 Wilson, D.J. (2019). The harmonic mean *p*-value for combining dependent tests.
1046 *Proceedings of the National Academy of Sciences* *116*, 1195-1200.
1047 [10.1073/pnas.1814092116](https://doi.org/10.1073/pnas.1814092116).

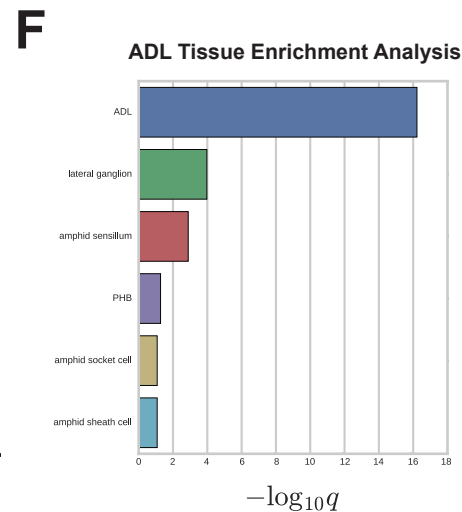
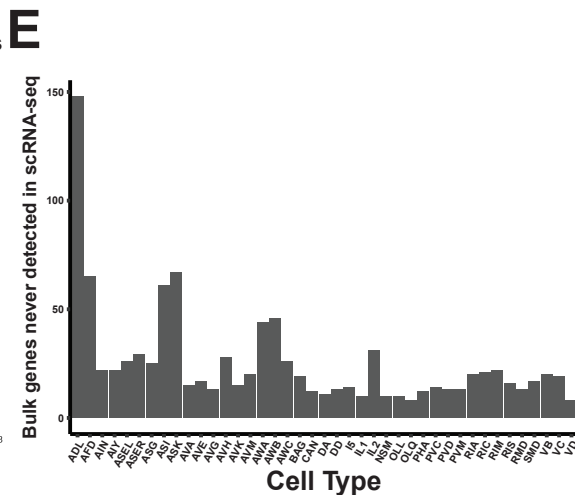
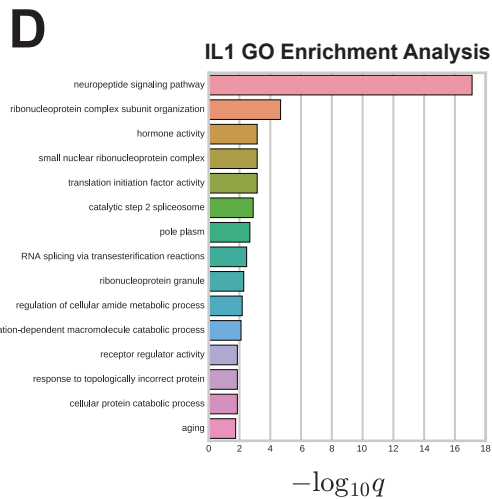
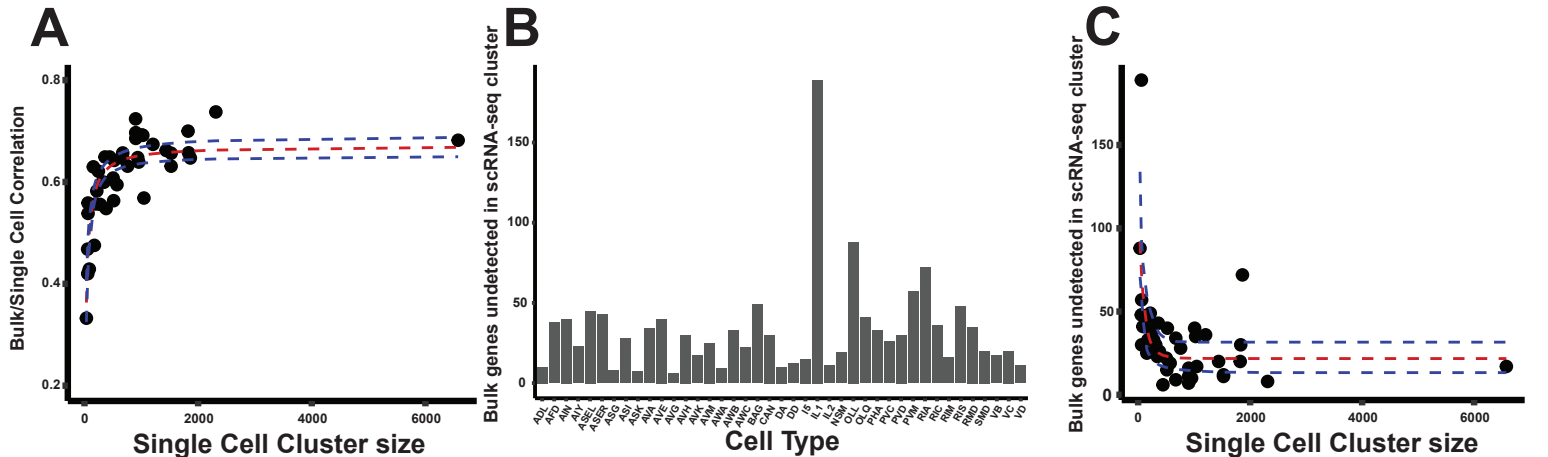
1048 Yemini, E., Lin, A., Nejatbakhsh, A., Varol, E., Sun, R., Mena, G.E., Samuel, A.D.T.,
1049 Paninski, L., Venkatachalam, V., and Hobert, O. (2021). NeuroPAL: A Multicolor Atlas for
1050 Whole-Brain Neuronal Identification in *C. elegans*. *Cell* *184*, 272-288 e211.
1051 [10.1016/j.cell.2020.12.012](https://doi.org/10.1016/j.cell.2020.12.012).

1052 Zhu, L., Lei, J., Devlin, B., and Roeder, K. (2018). A UNIFIED STATISTICAL FRAMEWORK
1053 FOR SINGLE CELL AND BULK RNA SEQUENCING DATA. *Ann Appl Stat* *12*, 609-632.
1054 [10.1214/17-AOAS1110](https://doi.org/10.1214/17-AOAS1110).

1055

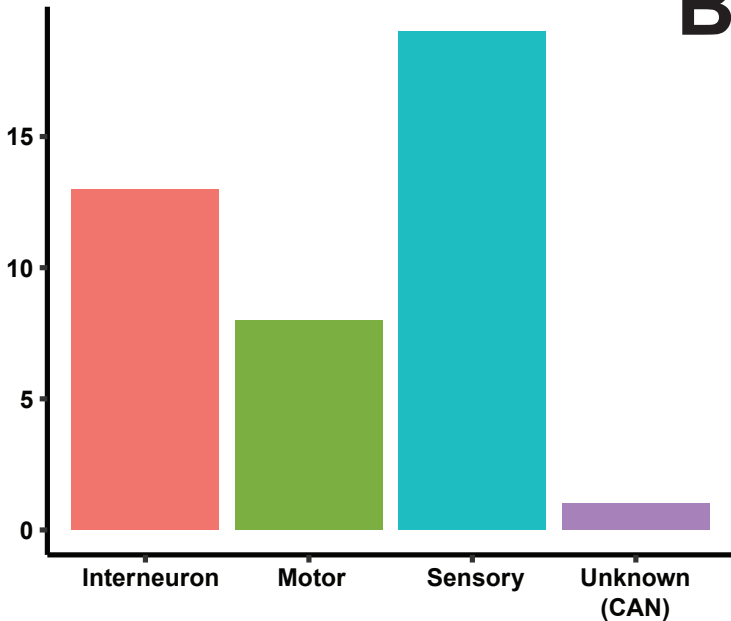
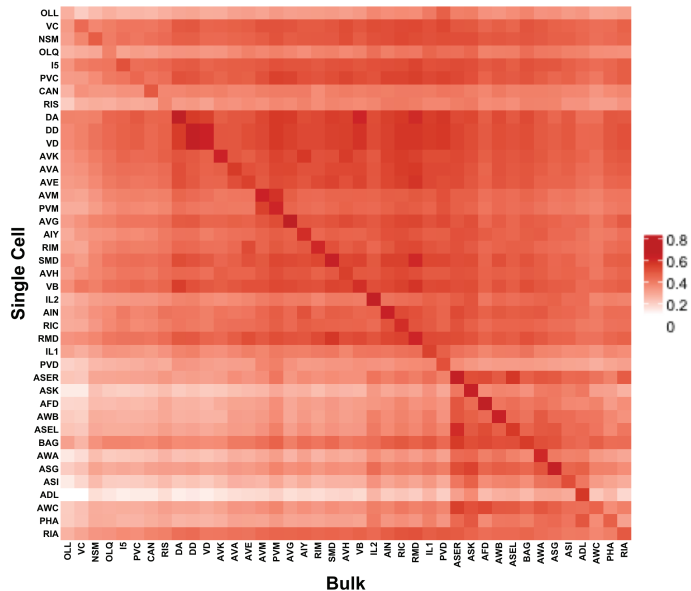




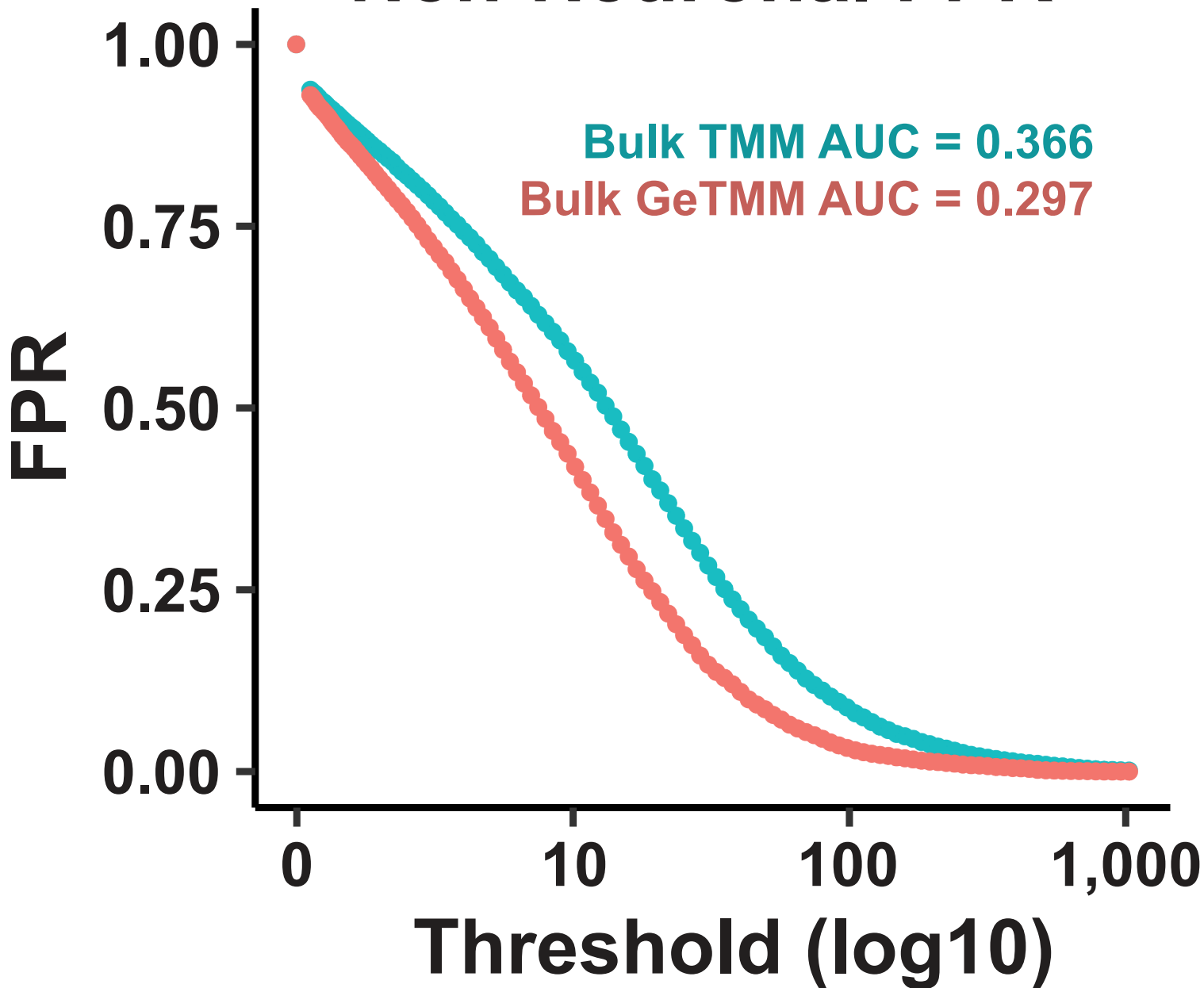


A

cell types per modality

**B****Spearman Correlation of Expressed Genes**

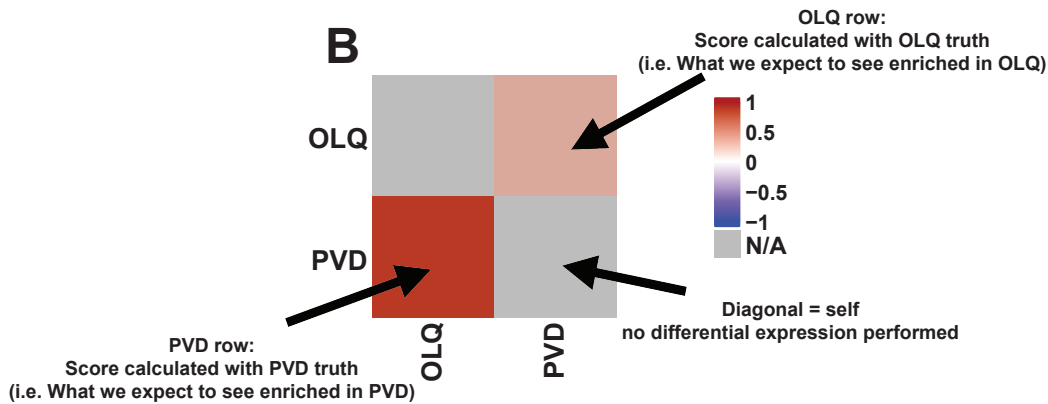
Non-Neuronal FPR



A Neuronal Ground-Truth example

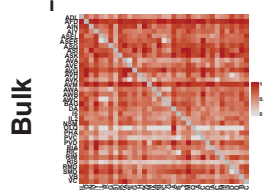
Gene	OLQ Truth	PVD Truth
a	1	0
b	0	1
c	0	0
d	0	0
e	0	0
f	0	0

B



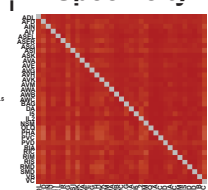
C

Recall



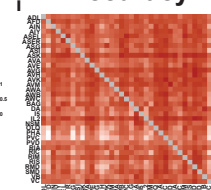
D

Specificity



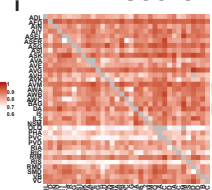
E

Accuracy

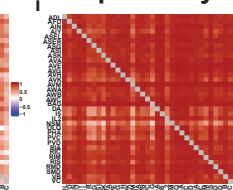


F

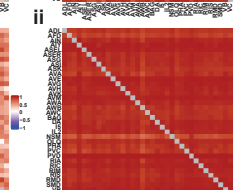
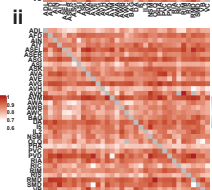
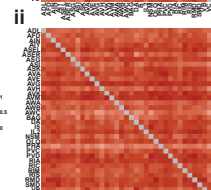
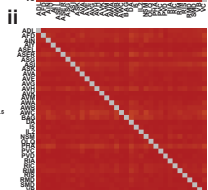
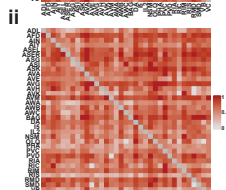
MCC score



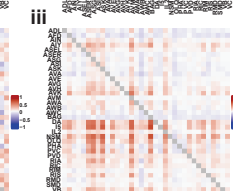
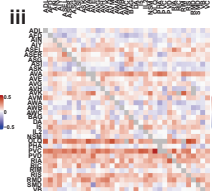
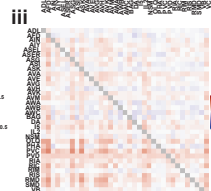
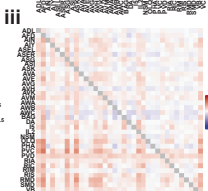
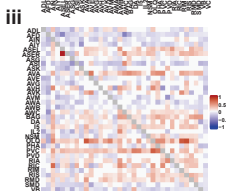
G Non-neuronal Specificity



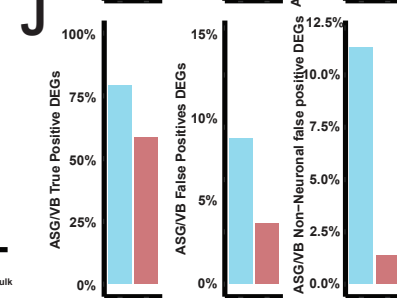
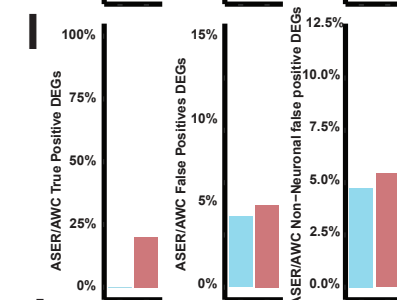
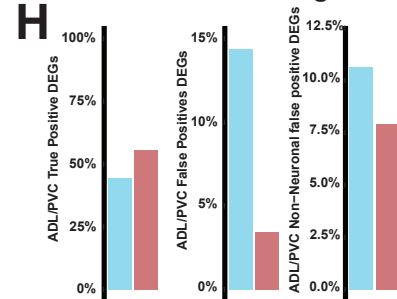
Integrated



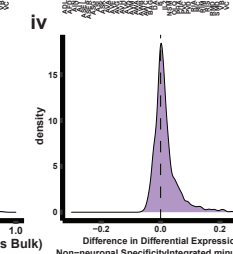
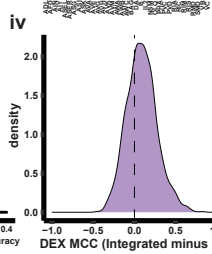
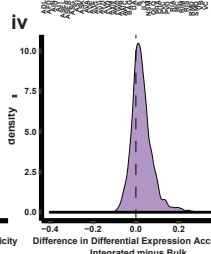
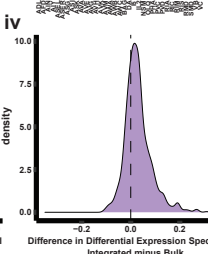
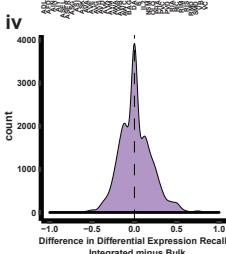
Difference



Bulk Integrated

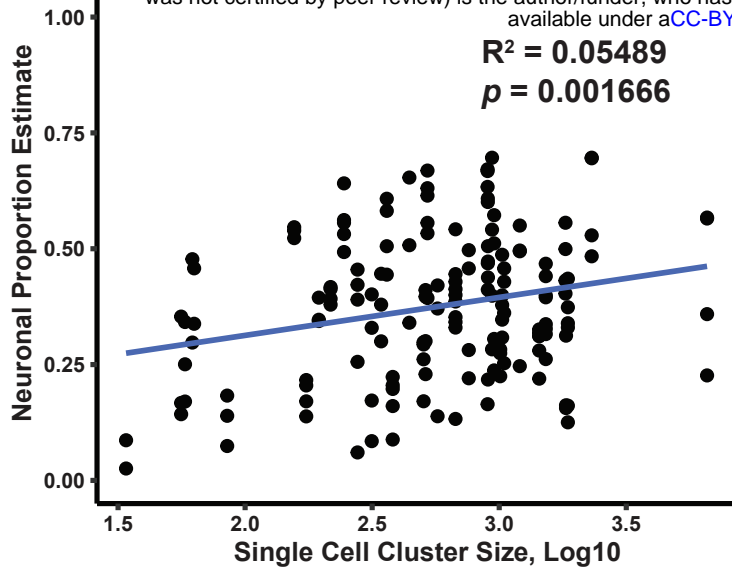
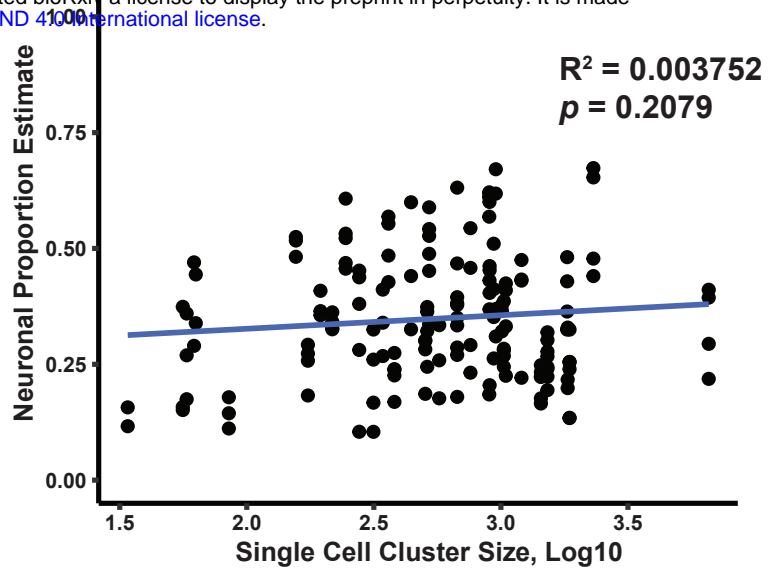
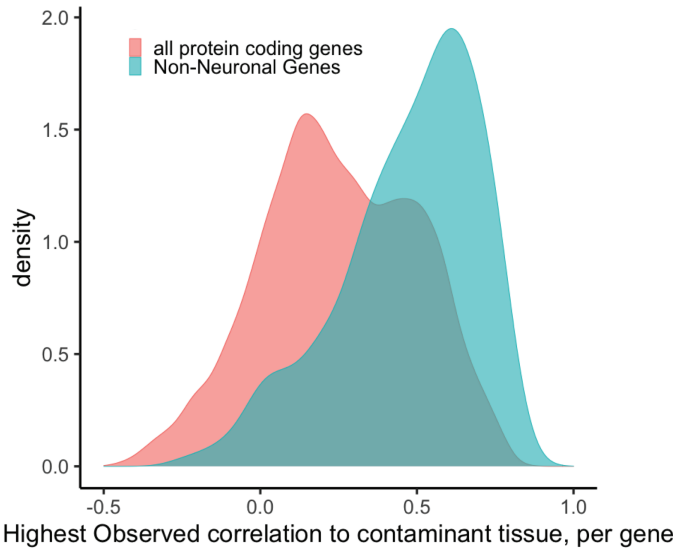
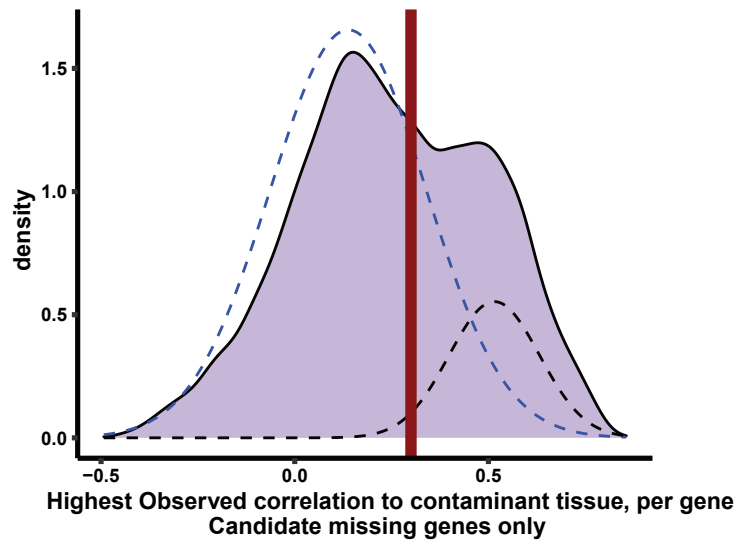
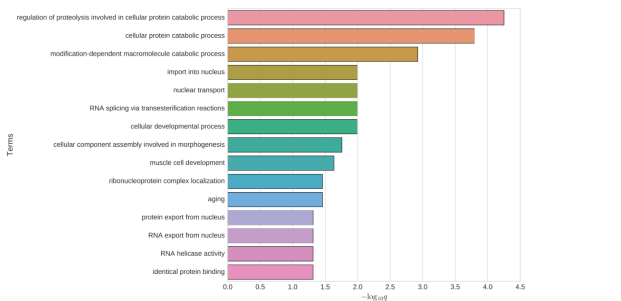
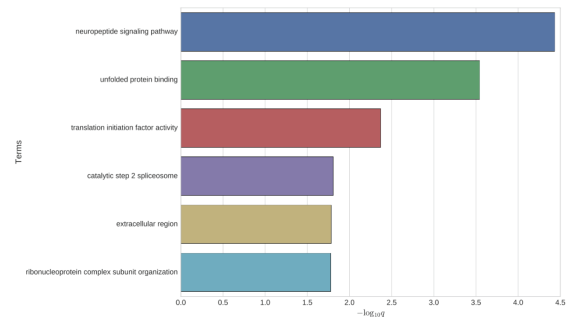
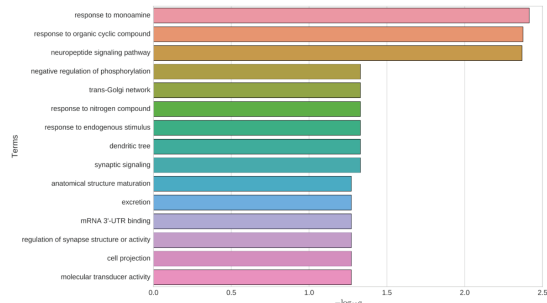
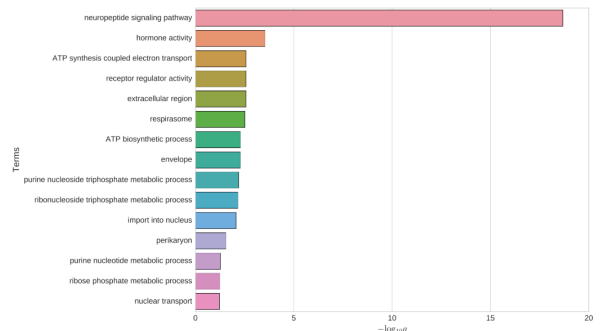


Difference Density

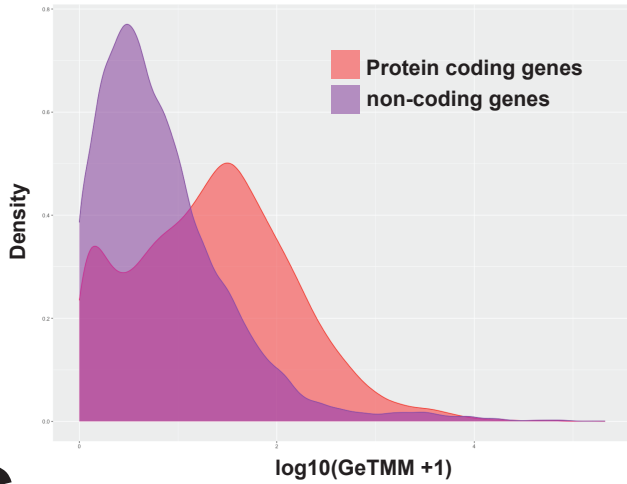


A**Full Population NNI S estimates**

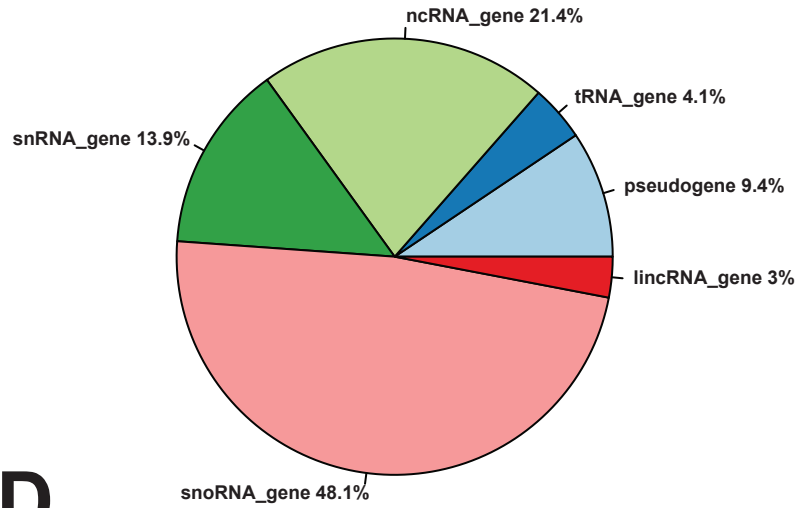
bioRxiv preprint doi: <https://doi.org/10.1101/2022.04.05.487209>; this version posted April 8, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

**B****30 cell bootstrap subsampled NNI S estimates****C****D****E****OLL****F****RIS****G****PVD****H****PVM**

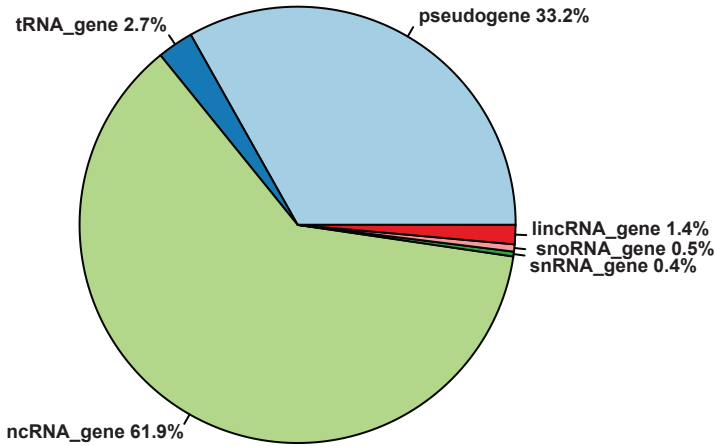
A Average GeTMM values for non-coding and mRNAs



B RNA classes in Pan-Neuronal noncoding RNAs



C RNA classes in Cell Type Specific noncoding RNAs



D Specific ncRNAs, grouped by modality

