

1 **Title**

2

3 Comprehensive single cell transcriptional profiling of a multicellular organism by combinatorial
4 indexing

5

6 **Authors**

7

8 Junyue Cao^{1,2}, Jonathan S. Packer^{1,*}, Vijay Ramani^{1,*}, Darren A. Cusanovich^{1,*}, Chau Huynh¹,
9 Riza Daza¹, Xiaojie Qiu^{1,2}, Choli Lee¹, Scott N. Furlan^{3,4,5}, Frank J. Steemers⁶, Andrew Adey^{7,8},
10 Robert H. Waterston^{1,#}, Cole Trapnell^{1,#}, Jay Shendure^{1,9,#}

11

12 ¹Department of Genome Sciences, University of Washington, Seattle, WA, USA

13 ²Molecular and Cellular Biology Program, University of Washington, Seattle, WA, USA

14 ³Ben Towne Center for Childhood Cancer Research, Seattle Children's Research Institute, Seattle, WA, USA

15 ⁴Department of Pediatrics, University of Washington, Seattle, WA, USA

16 ⁵Fred Hutchinson Cancer Research Center, Seattle, WA, USA

17 ⁶Illumina Inc., Advanced Research Group, San Diego, CA, USA

18 ⁷Department of Molecular & Medical Genetics, Oregon Health & Science University, Portland, OR, USA

19 ⁸Knight Cardiovascular Institute, Portland, OR, USA

20 ⁹Howard Hughes Medical Institute, Seattle, WA, USA

21

22 *These authors contributed equally to this work

23

24 #Correspondence to colettrap@uw.edu (CT), watersto@uw.edu (RHW) & shendure@uw.edu (JS)

25

26 **Abstract**

27
28 Conventional methods for profiling the molecular content of biological samples fail to resolve
29 heterogeneity that is present at the level of single cells. In the past few years, single cell RNA
30 sequencing has emerged as a powerful strategy for overcoming this challenge. However, its
31 adoption has been limited by a paucity of methods that are at once simple to implement and cost
32 effective to scale massively. Here, we describe a combinatorial indexing strategy to profile the
33 transcriptomes of large numbers of single cells or single nuclei without requiring the physical
34 isolation of each cell (Single cell Combinatorial Indexing RNA-seq or sci-RNA-seq). We show
35 that sci-RNA-seq can be used to efficiently profile the transcriptomes of tens-of-thousands of
36 single cells per experiment, and demonstrate that we can stratify cell types from these data. Key
37 advantages of sci-RNA-seq over contemporary alternatives such as droplet-based single cell
38 RNA-seq include sublinear cost scaling, a reliance on widely available reagents and equipment,
39 the ability to concurrently process many samples within a single workflow, compatibility with
40 methanol fixation of cells, cell capture based on DNA content rather than cell size, and the
41 flexibility to profile either cells or nuclei. As a demonstration of sci-RNA-seq, we profile the
42 transcriptomes of 42,035 single cells from *C. elegans* at the L2 stage, effectively 50-fold
43 “shotgun cellular coverage” of the somatic cell composition of this organism at this stage. We
44 identify 27 distinct cell types, including rare cell types such as the two distal tip cells of the
45 developing gonad, estimate consensus expression profiles and define cell-type specific and
46 selective genes. Given that *C. elegans* is the only organism with a fully mapped cellular lineage,
47 these data represent a rich resource for future methods aimed at defining cell types and states.
48 They will advance our understanding of developmental biology, and constitute a major step
49 towards a comprehensive, single-cell molecular atlas of a whole animal.
50

51 Introduction

52
53 Individual cells are the natural unit of form and function in biological systems. However,
54 conventional methods for profiling the molecular content of biological samples usually mask
55 cellular heterogeneity that is likely present even in ostensibly homogenous tissues. The
56 averaging of signals from large numbers of single cells sharply limits what we are able to learn
57 from the resulting data (1). For example, differences between samples cannot easily be attributed
58 to differences within cells of the same type vs. differences in cell type composition.

59
60 In the past few years, profiling the transcriptome of individual cells (*i.e.* single cell RNA-seq)
61 has emerged as a powerful strategy for resolving heterogeneity in biological samples. The
62 expression levels of mRNA species are readily linked to cellular function, and therefore can be
63 used to classify cell types in heterogeneous samples (2–10) as well as to order cell states in
64 dynamic systems (11). Although methods for single cell RNA-seq have proliferated, they
65 universally rely on the isolation of individual cells within physical compartments, whether by
66 pipetting (12–14), sorting (2, 8, 13, 15–17), microfluidics-based deposition to microwells (18), or
67 by dilution to emulsion-based droplets (5, 19, 20). As a consequence, the cost of preparing single
68 cell RNA-seq libraries with these methods scales linearly with the numbers of cells processed.
69 Although droplet-based methods are generally more cost-effective than well-based methods, they
70 are nonetheless limited by linear cost-scaling, the need for specialized instrumentation, and an
71 incompatibility with profiling nuclei or fixed cells (21, 22). Furthermore, droplet-based systems
72 capture cells based on cell size, which may bias analyses of heterogeneous tissues.

73
74 We set out overcome these limitations with a new method for single cell RNA-seq based on the
75 concept of combinatorial indexing. Combinatorial indexing uses split-pool barcoding to uniquely
76 label a large number of single molecules or single cells, but without requiring the physical
77 isolation of each molecule or cell. We previously used combinatorial indexing of single
78 molecules (high molecular weight genomic DNA fragments) for both haplotype-resolved
79 genome sequencing and *de novo* genome assembly (23, 24). More recently, we and others
80 demonstrated single cell combinatorial indexing (“sci”) to efficiently profile chromatin
81 accessibility (sci-ATAC-seq) (25), genome sequences (sci-DNA-seq) (26), and genome-wide
82 chromosome conformation (sci-Hi-C) (27) in large numbers of single nuclei.

83
84 Here we describe sci-RNA-seq, a straightforward method for profiling the transcriptomes of
85 large numbers of single cells or nuclei per experiment, using off-the-shelf reagents and
86 conventional instrumentation. For single cells, the protocol relies on methanol fixation, which
87 can stabilize and preserve RNA in dissociated cells for weeks (28), thereby minimizing
88 perturbations to cell state before or during processing. In this proof-of-concept, we apply sci-
89 RNA-seq to profile the transcriptomes of ~16,000 mammalian cells in a single experiment, and
90 show that we can separate synthetic mixtures of inter- or intraspecies cell types. We then apply
91 sci-RNA-seq to *Caenorhabditis elegans* worms at the L2 stage, sequencing the transcriptomes of
92 42,035 single cells. This includes 37,734 somatic cells, effectively 50x coverage (*i.e.*
93 oversampling) of this organism’s entire cellular content (762 cells at the L2 stage). From these
94 data, we identify 27 distinct cell types, estimate their consensus expression profiles and define
95 cell-type specific and selective genes, including for some fine-grained cell types that are present
96 in only one or two cells per animal.

97 **Results**

98

99 Overview of method

100

101 In its current form, sci-RNA-seq relies on the following steps (**Fig. 1a**): (1) Cells are fixed and
102 permeabilized with methanol (alternatively, cells are lysed and nuclei recovered), and then split
103 across 96- or 384-well plates. (2) To introduce a first molecular index to the mRNA of cells
104 within each well, we perform *in situ* reverse transcription (RT) with a barcode-bearing, well-
105 specific polyT primer with unique molecular identifiers (UMI). (3) All cells are pooled and
106 redistributed by fluorescence activated cell sorting (FACS) to 96- or 384-well plates in limiting
107 numbers (*e.g.* 10-100 per well). (4) We then perform second strand synthesis, transposition with
108 Tn5 transposase (Nextera), lysis, and PCR amplification. The PCR primers target the barcoded
109 polyT primer on one end, and the Tn5 adaptor insertion on the other end, such that resulting PCR
110 amplicons preferentially capture the 3' ends of transcripts. Critically, these primers introduce a
111 second barcode, specific to each well of the PCR plate. (5) Amplicons are pooled and subjected
112 to massively parallel sequencing, resulting in 3'-tag digital gene expression profiles, with each
113 read associated with two barcodes corresponding to the first and second rounds of cellular
114 indexing (**Fig. 1b**).

115

116 Because the overwhelming majority of cells pass through a unique combination of wells, each
117 cell is effectively “compartmentalized” by the unique combination of barcodes that tag its
118 transcripts. As we previously described, the rate of “collisions”, *i.e.* two or more cells receiving
119 the same combination of barcodes, can be tuned by adjusting how many cells are distributed to
120 the second set of wells (25). The sub-linear cost-scaling of combinatorial indexing follows from
121 the fact that the number of possible barcode combinations is the product of the number of
122 barcodes used at each stage. Consequently, increasing the number of barcodes used in the two
123 rounds of indexing leads to an increased capacity for the number of cells that can be profiled and
124 a lower effective cost per cell (**Fig. S1**). Additional rounds of molecular indexing can potentially
125 offer even greater complexity and lower costs. Pertinent to experiments described below, we
126 note that multiple samples (*e.g.* different cell populations, tissues, individuals, time-points,
127 perturbations, replicates, etc.) can be concurrently processed within a single sci-RNA-seq
128 experiment, simply by using different subsets of wells for each sample during the first round of
129 indexing.

130

131 Application of sci-RNA-seq to mammalian cells

132

133 To evaluate the scalability of this strategy, we performed a 384 well x 384 well sci-RNA-seq
134 experiment. During the first round of indexing, half of the 384 wells contained pure populations
135 of either human (HEK293T or HeLa S3) or mouse (NIH/3T3) cells, whereas half of the wells
136 contained a mixture of human and mouse cells. After barcoded RT, cells were pooled and then
137 sorted to a new 384 well plate for further processing, including the second round of barcoding
138 and deep sequencing of pooled PCR amplicons. From this single experiment, we recovered
139 15,997 single cell transcriptomes.

140

141 At a read depth corresponding to ~25,000 reads per cell (~65% duplication rate), we recovered
142 (on average) 11,024 UMIs per HEK293T cell, 7,832 UMIs per HeLa S3 cell, and 5,260 UMIs

143 per NIH/3T3 cell. Cells originating in wells containing an interspecies mixture were readily
144 assignable as human or mouse, with a rate of collision 16% (6.3% expected) (**Fig. 1c**). Species
145 mixing experiments provide a means of estimating ‘impurities’ that might result from mRNA
146 leakage between permeabilized cells. In this initial experiment, the rate of human reads
147 appearing in mouse-assigned cells was 6.2%, and the rate of mouse reads appearing in human
148 cells was 2.9%.

149
150 To reduce mRNA leakage as well as to increase the number of unique transcripts profiled per
151 cell, we extensively optimized the protocol, most importantly the choice of reverse transcriptase
152 and the washing steps. In an experiment that is representative of the culmination of these
153 optimizations, we performed 96 well x 96 well sci-RNA-seq on five mammalian cell
154 populations, split across 96 wells during the first round of barcoding: HEK293T cells (8 wells);
155 HeLa S3 cells (8 wells); an intraspecies mixture of HEK293T and HeLa S3 cells (32 wells); and
156 interspecies mixtures of HEK293T and NIH/3T3 cells (24 wells) or nuclei (24 wells). After
157 barcoded reverse transcription across these 96 wells, cells were pooled, sorted to a new 96 well
158 plate for a second round of barcoding by PCR, and then all amplicons pooled.

159
160 We deeply sequenced this library to a read depth corresponding to ~250,000 reads per cell (~88%
161 duplication rate). We grouped reads that shared the same first and second round barcodes
162 (inferring that each such group originated from the same single cell), which yielded 744 single
163 cell and 175 single nucleus transcriptomes. Transcript tags were aligned to a combined human
164 and mouse reference genome with STAR (29). Transcriptomes originating in ‘human only’ wells
165 overwhelmingly mapped to the human genome (99%), including 51,311 UMIs per cell from
166 HEK293T-only wells and 31,276 UMIs per cell from HeLa S3-only wells (on average), much
167 higher transcript counts per cell than our original experiment. 81% of reads mapped to the
168 expected strand of genic regions (47% exonic, 34% intronic), and the remainder to the
169 unexpected strand of genic regions (9%) or to intergenic regions (10%). These proportions are
170 similar to other studies (17). Whereas exonic reads show the expected enrichment at the 3’ ends
171 of gene bodies, intronic reads do not, and may be the result of poly(dT) priming from poly(dA)
172 tracts in heterogeneous nuclear RNA (**Fig. S2**). However, because like exonic reads they
173 overwhelmingly derived from the expected strand of genic regions, we retained them in
174 subsequent analyses.

175
176 In this experiment, transcriptomes originating in wells containing an interspecies mixture of
177 human (HEK293T) and mouse (NIH/3T3) cells overwhelmingly mapped to the genome of one
178 species or the other (289 of 294 cells), with only 5 clear ‘collisions’ (3.4% collision rate; 4.3%
179 expected) (**Fig. 2a**). Excluding these collisions, we observed 24,454 UMIs per human cell and
180 17,665 UMIs per mouse cell (**Fig. 2b-c**), with an average of 1.9% and 3.3% of reads per cell
181 mapping to the incorrect species, respectively, a level of impurity that is lower than our original
182 experiment and comparable to the droplet-based Drop-Seq method (19).

183
184 We next sought to confirm that we can separate cell types with sci-RNA-seq, focusing on the
185 two human cell types (HEK293T and HeLa S3 cells) included in this experiment. We performed
186 t-stochastic neighbor embedding (t-SNE) on standardized expression values of single cells
187 derived from wells containing pure HEK293T cells, pure HeLa S3 cells, or mixed HEK293T and
188 HeLa S3 cells. The mixed cells readily separate into two clusters, with one corresponding to

189 HEK293T cells and the other to HeLa S3 cells (**Fig. 2d**). To further validate the result, we
190 identified single nucleotide variants (SNVs) that distinguish HEK293T and HeLa S3 cells and
191 found that our assignments are in agreement with the SNV-based assignments (**Fig. S3a**). Of
192 note, the intronic reads alone are sufficient to separate HEK293T and HeLa S3 cells (**Fig. S3b**).

193
194 To evaluate whether sci-RNA-seq is biased compared to bulk measurements, we compared *in*
195 *silico* aggregated transcriptomes from all 220 identified HEK293T cells against the output of a
196 related bulk RNA-seq workflow (Tn5-RNA-seq (30)) without methanol fixation. The resulting
197 estimates of gene expression were highly correlated (Pearson: 0.93; **Fig. 2e**).

198 199 Application of sci-RNA-seq to mammalian nuclei

200
201 Protocols for dissociating tissues to single cell suspensions are labor intensive, do not easily
202 standardize, potentially impact cell states, and potentially bias cell type composition. To address
203 this, several groups have developed single nucleus RNA-seq protocols (since single nuclei are
204 much more readily obtained than single cells), but these ‘one-nucleus-per-well’ approaches (8,
205 10, 17) do not efficiently scale.

206
207 We therefore evaluated sci-RNA-seq for profiling the transcriptomes of single nuclei extracted
208 from a synthetic mixture of mouse (NIH/3T3) and human (HEK293T) cells. From the above-
209 described experiment in which these nuclei were sorted to 24 of the 96 wells during the first
210 round of barcoding, we recovered 175 single nucleus transcriptome profiles. Transcript tags were
211 aligned to a combined human and mouse reference genome with STAR (29). As with single cells,
212 reads associated with single nuclei mapped overwhelmingly to the genome of one species or the
213 other (48 human nuclei; 124 mouse nuclei) with 3.4% collisions rate (4.3% expected) (**Fig. 3a**).

214
215 Excluding collisions and at a read depth corresponding to ~210,000 reads per nucleus (~88%
216 duplication rate), we observed 32,951 UMIs per human nucleus and 20,123 UMIs per mouse
217 nucleus (**Fig. 3b-c**), with an average of 2.2% and 1.9% of reads per cell mapping to the incorrect
218 species. 84% of reads mapped to the expected strand of genic regions (35% exonic, 49% intronic)
219 and 16% to intergenic regions or to the unexpected strand of genic regions. These proportions are
220 similar to previous descriptions of single nucleus RNA-seq (17), which together with the number
221 of transcript molecules recovered indicate that sci-RNA-seq can flexibly and scalably profile the
222 transcriptomes of either cells or nuclei. As a further check, we aggregated the transcriptomes of
223 identified NIH/3T3 nuclei ($n = 124$) and compared them against the aggregated transcriptomes of
224 identified NIH/3T3 cells ($n = 129$) and found the resulting estimates of gene expression to be
225 highly correlated (Pearson: 0.97; **Fig. 3d**).

226 227 Single cell RNA profiling of *C. elegans*

228
229 The roundworm *C. elegans* is the only multicellular organism for which all cells and cell types
230 are defined, as is its entire developmental lineage (31, 32). As the extent to which contemporary
231 single cell experimental and computational methods can comprehensively recover and
232 distinguish cell types remains a matter of debate, we applied sci-RNA-seq to whole *C. elegans*
233 larvae. Of note, the cells in *C. elegans* larvae are much smaller, more variably sized, and have

234 markedly lower mRNA content than mammalian cells, and therefore represent a much more
235 challenging test of sci-RNA-seq's technical robustness.

236
237 We pooled ~150,000 larvae synchronized at the L2 stage and dissociated them into single-cell
238 suspensions. We then performed *in situ* reverse transcription across 6 96-well plates (*i.e.* 576
239 first-round barcodes), with each well containing ~1,000 *C. elegans* cells and also 1,000 human
240 cells (HEK293T) as internal controls. After pooling cells from all plates together, we sorted the
241 mixture of *C. elegans* cells and HEK293T cells into 10 new 96-well plates for PCR barcoding
242 (*i.e.* 960 second-round barcodes), gating on DNA content to distinguish between *C. elegans* and
243 HEK293T cells. This sorting was carried out such that 96% of wells harbored only *C. elegans*
244 cells (140 each), and 4% of wells harbored a mix of *C. elegans* and HEK293T cells (140 *C.*
245 *elegans* cells and 10 HEK293T cells each).

246
247 This single experiment yielded 42,035 *C. elegans* single-cell transcriptomes (number of UMI
248 counts for protein-coding transcripts > 100). 93.7% of reads mapped to the expected strand of
249 genic regions (91.7% exonic, 2.0% intronic). At a sequencing depth of ~20,000 reads per cell
250 and a duplication rate of 79.9%, we identified a median of 575 UMIs mapping to protein coding
251 genes per *C. elegans* cell (mean of 1,121 UMIs per cell), likely reflective of the lower mRNA
252 content of *C. elegans* cells relative to mammalian cells (**Fig. S4a**). Importantly, control wells
253 containing both *C. elegans* and HEK293T cells demonstrated clear separation between the two
254 species (**Fig. S4b**).

255
256 Semi-supervised clustering analysis segregated the cells into 29 distinct groups with the largest
257 containing 13,205 cells (31.4%) and the smallest containing only 131 cells (0.3%) (**Fig. 4a**).
258 Somatic cell types comprised 37,734 cells. We identified genes that were expressed specifically
259 in a single cluster, and by comparing those genes to expression patterns reported in the literature
260 assigned the clusters to cell types (**Table S2**). Of the 29 clusters, 21 represented exactly one
261 literature-defined cell type, 5 contained >1 distinct cell types, and 3 contained cells of unclear
262 type. Neurons, which were present in 7 clusters from the global analysis, were independently
263 subclustered, revealing 10 major neuronal subtypes (**Fig. 4b**).

264
265 Overall, the global and neuron-specific clustering analyses allowed for the construction of
266 expression profiles of 27 well-defined cell types (**Fig 4e**; not counting “unclassified neurons”
267 and “other glia”). 11 of these cell types are represented by 1 to 6 cells in an individual L2 worm.
268 Examples of these fine-grained cell types include the distal tip cells of the developing gonad (2
269 cells), the excretory canal cell (1 single cell) and canal associated neurons (2 cells), and the
270 amphid/phasmid sheath cells (4 glial cells).

271
272 Each *C. elegans* animal contains exactly the same number of cells, which are generated
273 deterministically in a defined lineage (31, 32). Comparing the observed proportions of each cell
274 type to their known frequencies in L2 larvae showed that sci-RNA-Seq captured many cell types
275 at or near expected frequencies (**Fig. 4c**; 15 types had abundance \geq 50% of the expectation).
276 Body wall muscle was 2.4-fold more abundant in our data than in the animal, probably reflecting
277 its ease of dissociation. Only one abundant cell type, the intestinal cells, was absent. We
278 speculate that this was due to our gating strategy during FACS, which excluded cells based on

279 DAPI signal. Intestinal cells at the L2 stage are polyploid, containing 8 copies of each (haploid)
280 chromosome.

281
282 Previous analyses of single-cell transcriptomes have shown that cells can often be distinguished
283 with relatively light sequencing per cell on the basis of a small set of highly specific genes (33).
284 This raised the possibility that despite being able to detect many distinct cell types in the worm,
285 our molecular definition for each would be incomplete. However, we observed that half of all *C.*
286 *elegans* protein-coding genes were expressed in at least 80 cells in the full dataset, and 64% of
287 protein-coding genes were expressed in at least 20 cells. As many genes are expressed
288 specifically in embryos or adults, our gene expression measurements likely include most of the
289 protein-coding genes that are truly expressed in L2 *C. elegans*. The “whole-worm” expression
290 profile derived by aggregating all sci-RNA-Seq reads correlated well with published whole-
291 organism bulk RNA-seq (34) for L2 *C. elegans* (Fig. 4d, Spearman $\rho = 0.795$). Furthermore,
292 consensus expression profiles for each cell type segregated as expected in a hierarchical
293 clustering analysis (Fig. 4e). Thus, despite the fact that sci-RNA-Seq captures a minority of
294 transcripts in each cell, our ‘oversampling’ of the cellular composition of the organism (*i.e.*
295 37,734 somatic cells, effectively 50x coverage of the L2 worm’s 762 somatic cells), enables us to
296 construct faithful expression profiles for individual cell types.

297 298 Discussion

299
300 We describe a new method for single cell RNA-seq based on the principle of combinatorial
301 indexing of cells or nuclei. At the scale described here (*e.g.* 576 x 960 indexing), sci-RNA-seq
302 can be applied to profile the transcriptomes of tens-of-thousands of single cells per experiment.
303 Library preparation can be completed by a single individual in two days. The cost of library
304 preparation, currently \$0.03-\$0.20 per cell, is dominated by enzymes. The method relies on off-
305 the-shelf reagents and widely available instrumentation (*e.g.* FACS, sequencer).

306
307 sci-RNA-seq has several practical advantages over contemporary alternatives:

308
309 First, it is compatible (and indeed relies on) cell fixation, which can minimize perturbations to
310 cell state or RNA integrity before or during processing.

311
312 Second, sci-RNA-seq facilitates the concurrent processing of multiple samples (*e.g.*
313 corresponding to different cell populations, tissues, individuals, time-points, perturbations,
314 replicates, etc.) within a single experiment, simply by using different subsets of wells for each
315 sample during the first round of indexing. This may have the benefit of reducing batch effects,
316 relative to platforms requiring the serial processing of samples, an area of paramount concern for
317 the single cell RNA-seq field (35). Furthermore, given that the second barcode is introduced after
318 flow sorting, it is also possible to associate wells on the PCR plate with FACS-defined
319 subpopulations (*e.g.* corresponding to cell size, cell cycle, immunostaining, etc.), as we did while
320 sorting mixed HEK293T and *C. elegans* cells.

321
322 Third, as we show, sci-RNA-seq is readily compatible with the processing of nuclei. Single
323 nucleus RNA-seq is possible with “well-per-well” methods but is not currently supported on
324 droplet-based platforms. The ability to process nuclei, rather than cells, may be particularly

325 important for tissues for which unbiased cell disaggregation protocols are not well established
326 (which may be most tissues).

327

328 Fourth, sci-RNA-seq is highly scalable. Here we demonstrate up to 576 x 960 combinatorial
329 indexing, which enables the generation of $\sim 5 \times 10^4$ single cell transcriptomes. While beyond the
330 scope of this proof-of-concept, one can imagine simply using more barcoded RT and PCR
331 primers (*e.g.* 1,536 x 1,536 combinatorial indexing), which would enable processing of many
332 more cells with sub-linear scaling of the cost per cell. A complementary approach is simply to
333 introduce additional rounds of indexing, *e.g.* by using indexed Tn5 complexes during *in situ*
334 transposition (24, 25). With 384 x 384 x 384 combinatorial indexing, we can potentially uniquely
335 barcode the transcriptomes of >10 million cells within a single experiment.

336

337 In this proof-of-concept, we apply sci-RNA-Seq to generate the first catalog of single cell
338 transcriptomes at the scale of a whole organism. In a single experiment, we generated $\sim 50X$
339 coverage of the somatic cellular composition of L2 *C. elegans*, detecting 27 cell types and
340 constructing consensus transcriptional profiles for each. While not all cell types are detected at
341 expected frequencies, 15 cell types had abundance $\geq 50\%$ of the expectation. Half of *C. elegans*
342 protein-coding genes were expressed in at least 80 cells in the full dataset, demonstrating that our
343 expression profile for this stage of *C. elegans* development is mostly comprehensive. Taken
344 together, these analyses show that sci-RNA-Seq constitutes a reliable platform for molecular
345 dissection of complex tissues and potentially whole organisms.

346

347 sci-RNA-seq further expands the repertoire of single cell molecular phenotypes that can be
348 resolved by combinatorial indexing (which now includes mRNA, chromatin accessibility,
349 genome sequence, and chromosome conformation). Looking forward, we anticipate that
350 additional forms of single cell profiling can be achieved with combinatorial indexing. Provided
351 that multiple aspects of cellular biology can be concurrently barcoded, combinatorial indexing
352 may also facilitate the scalable generation of ‘joint’ single cell molecular profiles (*e.g.* RNA-seq
353 and ATAC-seq from each of many single cells).

354

355 **Materials & Methods**

356

357 Cell Culture

358 All cells were cultured at 37°C with 5% CO₂, and were maintained in high glucose DMEM
359 (Gibco cat. no. 11965) supplemented with 10% FBS and 1X Pen/Strep (Gibco cat. no. 15140122;
360 100U/ml penicillin, 100µg/ml streptomycin). Cells were trypsinized with 0.25% trypsin-EDTA
361 (Gibco cat. no. 25200-056) and split 1:10 three times a week.

362

363 Generation of whole *C. elegans* cell suspensions

364 A *C. elegans* strain (RW12139 *stIs11435(unc-120::H1-Wcherry;unc-119(+));unc-119(tm4063)*)
365 carrying an integrated Punc-120::mCherry gene in a wild type background was used in all
366 experiments. A synchronized L2 population was obtained by two cycles of bleaching gravid
367 adults to isolate fertilized eggs allowing the eggs to hatch in the absence of food to generate a
368 population of starved L1 animals. Around 150,000 L1 larvae were plated on each 100 mm petri
369 plate seeded with NA22 bacteria and incubated at 24 °C for 15 hr to produce early L2 larvae.
370 Dissociated cells were recovered following a published protocol (36) with modification.
371 Specifically, L2 stage worms were collected by adding 10 ml sterile ddH₂O to each plate. The
372 collected L2s were pelleted by centrifugation at 1300 g for 1 min. The larval pellet was washed
373 five times with sterile ddH₂O to remove bacteria. The resulting pellet was transferred to a 1.6 ml
374 microcentrifuge tube. Around 40 µl of the final compact pellet was used for each cell
375 dissociation experiment. The worm pellet was treated with 250 µl of SDS-DTT solution (20 nM
376 HEPES pH8, 0.25% SDS, 200 mM DTT, 3% sucrose) for 4 min. Immediately after SDS-DTT
377 treatment, egg buffer was added to the SDS-DTT treated worms. Worms were pelleted at 500 g
378 for 1 min, then washed 5 times with Egg buffer (118 mM NaCl, 48 mM KCl, 3 mM CaCl₂, 3
379 mM MgCl₂, 5 mM HEPES (pH 7.2)). Pelleted SDS-DTT treated worms were digested with 200
380 µl of 15 mg/ml pronase (Sigma-Aldrich, St. Louis, MO) for 20 min. The treated worms were
381 broken up to release cells by pipetting up and down with 21G1 ¼ needle. When sufficient single
382 cells were observed the reaction was stopped by adding 900 µl L-15 medium containing 10%
383 fetal bovine serum. Cells were separated from worm debris by centrifuging the pronase-treated
384 worms at 150 g for 5 min at 4°C. The supernatant was transferred to 1.6 ml microcentrifuge tube
385 and centrifuged at 500 g for 5 min at 4°C. The cell pellet was washed twice with egg-buffer
386 containing 1% BSA.

387

388 Sample Processing

389 All cell lines were trypsinized, spun down at 300xg for 5 min (4°C). and washed once in 1X PBS.
390 *C. elegans* cells were dissociated as described above.

391

392 For sci-RNA-seq on whole cells, 5M cells were fixed in 5 mL ice-cold 100% methanol at -20 °C
393 for 10 min, washed twice with 1 ml ice-cold 1X PBS containing 1% Diethyl pyrocarbonate (0.1%
394 for *C. elegans* cells) (DEPC; Sigma-Aldrich), washed three times with 1 mL ice-cold PBS
395 containing 1% SUPERase In RNase Inhibitor (20 U/µL, Ambion) and 1% BSA (20 mg/ml,
396 NEB). Cells were resuspended in wash buffer at a final concentration of 5000 cells/ul. For all
397 washes, cells were pelleted through centrifugation at 300xg for 3 min, at 4°C.

398

399 For sci-RNA-seq on nuclei, 5M cells were combined and lysed using 1 mL ice-cold lysis buffer
400 (10 mM Tris-HCl, pH 7.4, 10 mM NaCl, 3 mM MgCl₂ and 0.1% IGEPAL CA-630 from (37)),

401 modified to also include 1% SUPERase In and 1% BSA). The isolated nuclei were then pelleted,
402 washed twice with 1 mL ice-cold 1X PBS containing 1% DEPC, twice with 500 μ L cold lysis
403 buffer, once with 500 μ L cold lysis buffer without IGEPAL CA-630, and then resuspended in
404 lysis buffer without IGEPAL CA-630 at a final concentration of 5000 nuclei/ μ L. For all washes,
405 nuclei were pelleted through centrifugation at 300xg for 3 min. at 4°C).

406
407 For cell-mixing experiments, trypsinized cells were counted and the appropriate number of cells
408 from each cell line were combined prior to fixation or lysis. Fixed cells or nuclei were then
409 distributed into 96- or 384-well plates (see **Table S1**). For each well, 1,000-10,000 cells or nuclei
410 (2 μ L) were mixed with 1 μ L of 25 μ M anchored oligo-dT primer (5'-
411 ACGACGCTCTTCCGATCTNNNNNNNN[10bp
412 index]TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTVN
413 -3', where "N" is any base and "V" is either "A", "C" or "G"; IDT) and 0.25 μ L 10 mM dNTP
414 mix (Thermo), denatured at 55°C for 5 min and immediately placed on ice. 1.75 μ L of first-
415 strand reaction mix, containing 1 μ L 5X Superscript IV First-Strand Buffer (Invitrogen), 0.25 μ L
416 100 mM DTT (Invitrogen), 0.25 μ L SuperScript IV reverse transcriptase (200 U/ μ L, Invitrogen),
417 0.25 μ L RNaseOUT Recombinant Ribonuclease Inhibitor (Invitrogen), was then added to each
418 well. Reverse transcription was carried out by incubating plates at 55°C for 10 min, and was
419 stopped by adding 5 μ L 2X stop solution (40 mM EDTA, 1 mM spermidine) to each well. All
420 cells (or nuclei) were then pooled, stained with 4',6-diamidino-2-phenylindole (DAPI, Invitrogen)
421 at a final concentration of 3 μ M, and sorted at varying numbers of cells/nuclei per well
422 (depending on experiment; see **Table S1**) into 5 μ L buffer EB using a FACSAria III cell sorter
423 (BD). 0.5 μ L mRNA Second Strand Synthesis buffer (NEB) and 0.25 μ L mRNA Second Strand
424 Synthesis enzyme (NEB) were then added to each well, and second strand synthesis was carried
425 out at 16°C for 150 min. The reaction was then terminated by incubation at 75°C for 20 min.

426
427 Tagmentation was carried out on double-stranded cDNA using the Nextera DNA Sample
428 Preparation kit (Illumina). Each well was mixed with 5 ng Human Genomic DNA (Promega), as
429 carrier to avoid over-tagmentation and reduce losses during purification, 5 μ L Nextera TD buffer
430 (Illumina) and 0.5 μ L TDE1 enzyme (Illumina), and then incubated at 55 °C for 5 min to carry
431 out tagmentation. Note that because the PCR primers used to amplify libraries are specific to the
432 RT products, tagmented carrier genomic DNA are not appreciably amplified or sequenced. The
433 reaction was then stopped by adding 12 μ L DNA binding buffer (Zymo) and incubating at room
434 temperature for 5 min. Each well was then purified using 36 μ L AMPure XP beads (Beckman
435 Coulter), eluted in 16 μ L of buffer EB (Qiagen), then transferred to a fresh multi-well plate.

436
437 For PCR reactions, each well was mixed with 2 μ L of 10 μ M P5 primer (5'-
438 AATGATACGGCGACCACCGAGATCTACAC[i5]ACACTCTTCCCTACACGACGCTCTT
439 CCGATCT-3'), 2 μ L of 10 μ M P7 primer (5'-
440 CAAGCAGAAGACGGCATAACGAGAT[i7]GTCTCGTGGGCTCGG-3'), and 20 μ L NEBNext
441 High-Fidelity 2X PCR Master Mix (NEB). Amplification was carried out using the following
442 program: 75°C for 3 min, 98°C for 30 sec, 18-22 cycles of (98°C for 10 sec, 66°C for 30 sec,
443 72°C for 1 min) and a final 72°C for 5 min. After PCR, samples were pooled and purified using
444 0.8 volumes of AMPure XP beads. Library concentrations were determined by Qubit (Invitrogen)
445 and the libraries were visualized by electrophoresis on a 6% TBE-PAGE gel. Libraries were

446 sequenced on the NextSeq 500 platform (Illumina) using a V2 75 cycle kit (Read 1: 18 cycles,
447 Read 2: 52 cycles, Index 1: 10 cycles, Index 2: 10 cycles).

448 Read alignments and construction of gene expression matrix

450 Base calls were converted to fastq format and demultiplexed using Illumina's bcl2fastq/
451 2.16.0.10 tolerating one mismatched base in barcodes (edit distance (ED) < 2). Demultiplexed
452 reads were then adaptor clipped using trim_galore/0.4.1 with default settings. Trimmed reads
453 were mapped to the human reference genome (hg19), mouse reference genome (mm10),
454 *C.elegans* reference genome (PRJNA13758) or a chimeric reference genome of hg19, mm10 and
455 PRJNA13758, using STAR/v 2.5.2b (38) with default settings and gene annotations (GENCODE
456 V19 for human; GENCODE VM11 for mouse, WormBase
457 PRJNA13758.WS253.canonical_gene set for *C.elegans*). Uniquely mapping reads were
458 extracted, and duplicates were removed using the unique molecular identifier (UMI) sequence
459 (ED < 2, including insertions and deletions), reverse transcription (RT) index, and read 2 end-
460 coordinate (*i.e.* reads with identical UMI, RT index, and tagmentation site were considered
461 duplicates). Finally, mapped reads were split into constituent cellular indices by further
462 demultiplexing reads using the RT index (ED < 2, including insertions and deletions). For
463 mixed-species experiment, the percentage of uniquely mapping reads for genomes of each
464 species was calculated. Cells with over 85% of UMIs assigned to one species were regarded as
465 species-specific cells, with the remaining cells classified as mixed cells. The collision rate was
466 calculated as twice the ratio of mixed cells (as we are blind to any collisions involving cells of
467 the same species). For gene body coverage analysis of exonic reads, the split human and mouse
468 single cell SAM files were concatenated and exonic reads were selected and analyzed using
469 RSEQC/2.6.1, using BED annotation files downloaded from the UCSC Golden Path. For read
470 position analysis for intronic reads, the split human and mouse single cell SAM files were
471 concatenated and intronic reads were selected; the fractional position of each intronic read along
472 the genomic distance between the TSS and transcript terminus was calculated, and these values
473 used to generate a density plot.

474
475 To generate digital expression matrices, we calculated the number of strand-specific UMIs
476 mapping to the exonic and intronic regions of each gene, for each cell; generally, fewer than 3%
477 of total UMIs strand-specifically mapped to multiple genes. For multi-mapped reads, reads were
478 assigned to the closest gene, except in cases where another intersected gene fell within 100 bp to
479 the end of the closest gene, in which case the read was discarded. For most analyses we included
480 both intronic and exonic UMIs in per-gene single-cell expression matrices.

481 t-SNE visualization of HEK293T cells and HeLa S3 cells

482 We visualized the clustering of sci-RNA-seq data from populations of pure HEK293T, pure
483 HeLa S3 and mixed HEK293T + HeLa S3 cells using t-Distributed Stochastic Neighbor
484 Embedding (tSNE). The top 3,000 genes with the highest variance in the digital gene expression
485 matrix for these cells were first given as input to Principal Components Analysis (PCA). The top
486 10 principal components were then used as the input to t-SNE, resulting in the two-dimensional
487 embedding of the data show in **Fig. 2D**. The process was repeated using only intronic reads (**Fig.**
488 **S3B**). For this analysis, the top 2,000 (instead of 3,000) highly variable genes were used as input
489 to PCA; all other parameters remained unchanged.

491

492 Genotyping of single HeLa cells by 3' tag sequences

493 HeLa S3 cell identity was verified on the basis of homozygous alleles not present in the hg19
494 assembly, using a callset derived from (39). Single-cell BAM files (with cellular indices encoded
495 in the “read_id” field) were concatenated, and then processed as follows using a python wrapper
496 of the samtools API (i.e. pysam). For each homozygous alternate SNV overlapping with a
497 GENCODE V19 defined gene ($n = 865,417$) in the HeLa S3 variant callset, we computed the
498 fraction of matching (i.e. HeLa S3 specific) alleles, and computed this value for all cells where at
499 least 1 read containing a polymorphic site. We then re-plotted in R the tSNE visualization shown
500 in **Fig. 2D**, now colored by the relative fraction of homozygous alternate alleles called for each
501 cell.

503 Analysis of *C. elegans* whole-organism sci-RNA-seq experiment

504 A digital gene expression matrix was constructed from the raw sequencing data as described
505 above. The dimensionality of this matrix was reduced first with PCA (40 components) and then
506 with t-SNE, giving a two-dimensional representation of the data. Similar to the approach in (40),
507 cells in this two-dimensional representation were clustered using the density peak algorithm (41)
508 as implemented in Monocle 2. Genes specific to each cluster were identified and compared to
509 microscopy-based expression profiles reported in the literature (**Table S2**), allowing the distinct
510 cell types represented in each cluster to be identified. Based on these results, we manually
511 merged two clusters that both corresponded to body wall muscle, and manually split two clusters
512 that included hypodermis, somatic gonad cells, and glia. Seven clusters exclusively contained
513 neurons. We identified neuronal subtypes applying PCA, t-SNE, and density peak clustering to
514 this subset of cells using the same approach as for the global cluster analysis.

516 Consensus expression profiles for each cell type were constructed by first dividing each column
517 in the gene-by-cell digital gene expression matrix by the cell's size factor and then for each cell
518 type, taking the mean of the normalized UMI counts for the subset of cells assigned to that cell
519 type. These mean normalized UMI counts were then re-scaled to transcripts per million. Cells
520 that were part of a cluster corresponding to a cell type but did not express any of the marker
521 genes used to define that cell type were excluded when generating these consensus expression
522 profiles.

524 Comparing sci-RNA-seq and bulk RNA-seq data for HEK293T cells

525 To compare aggregated sci-RNA-seq single cell transcriptomes with bulk RNA-seq, we
526 performed bulk RNA-seq using a modified protocol (30). In brief, 500 ng total RNA extracted
527 from three biological replicate HEK293T samples (extraction using RNeasy kit (Qiagen)) with
528 the RNeasy kit (Qiagen) were used for reverse transcription following the standard SuperScript
529 II protocol. 500 ng total RNA (in 9 μ L water) was mixed with 2 μ L 25 uM oligo-dT(VN) (5'-
530 ACGACGCTCTCCGATCTNNNNNNNN[10bp
531 index]TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTVN
532 -3', where “N” is any base and “V” is either “A”, “C” or “G”; IDT) and 1 μ L 10 mM dNTPs,
533 then incubated at 65°C for 5 min Following incubation, 8 μ L reaction mix (4 μ L 5X Superscript
534 II First-Strand Buffer, 2 μ L 100 mM DTT, 1 μ L SuperScript II reverse transcriptase, 1 μ L
535 RnaseOUT) was added. Reactions were incubated at 42°C for 50 min and terminated at 70°C for
536 15 min. For second strand synthesis, 2 μ L RT product was mixed with 6.5 μ L water, 1 μ L
537 mRNA Second Strand Synthesis buffer (NEB) and 0.25 μ L mRNA Second Strand Synthesis

538 enzyme (NEB). Second strand synthesis was carried out at 16°C for 150 min, followed by 75°C
539 for 20 min. Tagmentation was carried out by adding 10 µL Nextera TD buffer, 1 µL Nextera Tn5
540 enzyme and incubating at 55°C for 5 min. Tagmented cDNA was purified using a Clean &
541 ConcentratorTM-100 kit (Zymo) and eluted in 16 µL buffer EB. PCR, purification, and
542 quantification were then performed as detailed above.

543
544 For comparing single cell RNA-seq and bulk RNA-seq, single cell gene counts of exonic reads
545 and intronic reads were added for the same gene from sci-RNA-seq of pure HEK293T cells as
546 well as HEK293T cells identified from HEK293T and NIH/3T3 mixed cells. Counts for bulk
547 RNA-seq of HEK293T cells were extracted based on the RT barcode and aggregated separately,
548 again adding exonic and intronic read counts per gene. Transcript counts were converted to
549 transcripts per million (TPM) and then transformed to $\log(\text{TPM} + 1)$. Pearson correlation
550 coefficients were calculated between the aggregated sci-RNA-seq and bulk RNA-seq data using
551 R.

552 553 Cost estimation

554 For 576 x 960 sci-RNA-seq, reagent costs are largely enzyme-driven and include SuperScript IV
555 reverse transcriptase (\$934), second strand synthesis mix (\$750), Nextera Tn5 enzyme (\$5,000),
556 NEBnext master mix (\$1,150) and other reagents and plates (\$500). If we sort 60 cells per well
557 (assuming recovery rate is 100%) for 960 wells (5% collision rate), then the reagent cost of
558 library preparation per single cell is around \$0.14 (expected yield of around 55,000 cells).
559 However, it is worth noting that simply increasing the number of cells sorted per well also
560 decreases costs (*e.g.* sort 150 cells to each well would yield around 140,000 cells at a cost of
561 \$0.05 per cell), but also results in an increased collision rate (12%). Alternatively, by increasing
562 to 1,536 barcodes during the first (RT-based) round of indexing, we can sort up to 320 cells per
563 well at a 10% collision rate, thereby reducing the cost per cell to less than \$0.025 per cell.
564 Straightforward reductions in reaction volumes at all steps may also lead to further reductions in
565 costs, as would additional rounds of molecular indexing.

566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583

584 **References**

585

- 586 1. C. Trapnell, Defining cell types and states with single-cell genomics. *Genome Res.* **25**,
587 1491–1498 (2015).
- 588 2. D. Ramsköld *et al.*, Full-length mRNA-Seq from single-cell levels of RNA and individual
589 circulating tumor cells. *Nat. Biotechnol.* **30**, 777–782 (2012).
- 590 3. A. K. Shalek *et al.*, Single-cell transcriptomics reveals bimodality in expression and
591 splicing in immune cells. *Nature.* **498**, 236–240 (2013).
- 592 4. Q. F. Wills *et al.*, Single-cell gene expression analysis reveals genetic associations masked
593 in whole-tissue experiments. *Nat. Biotechnol.* **31**, 748–752 (2013).
- 594 5. G. X. Y. Zheng *et al.*, Massively parallel digital transcriptional profiling of single cells.
595 *bioRxiv*, 65912 (2016).
- 596 6. A. A. Pollen *et al.*, Low-coverage single-cell mRNA sequencing reveals cellular
597 heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat.*
598 *Biotechnol.* **32**, 1053–1058 (2014).
- 599 7. A. Zeisel *et al.*, Cell types in the mouse cortex and hippocampus revealed by single-cell
600 RNA-seq. *Science.* **347**, 1138–1142 (2015).
- 601 8. B. B. Lake *et al.*, Neuronal subtypes and diversity revealed by single-nucleus RNA
602 sequencing of the human brain. *Science.* **352**, 1586–1590 (2016).
- 603 9. I. Tirosh *et al.*, Dissecting the multicellular ecosystem of metastatic melanoma by single-
604 cell RNA-seq. *Science.* **352**, 189–196 (2016).
- 605 10. W. Zeng *et al.*, Single-nucleus RNA-seq of differentiating human myoblasts reveals the
606 extent of fate heterogeneity. *Nucleic Acids Res.* (2016), doi:10.1093/nar/gkw739.
- 607 11. C. Trapnell *et al.*, The dynamics and regulators of cell fate decisions are revealed by
608 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 609 12. F. Tang *et al.*, mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods.* **6**,
610 377–382 (2009).
- 611 13. S. Islam *et al.*, Characterization of the single-cell transcriptional landscape by highly
612 multiplex RNA-seq. *Genome Res.* **21**, 1160–1167 (2011).
- 613 14. T. Hashimshony, F. Wagner, N. Sher, I. Yanai, CEL-Seq: Single-Cell RNA-Seq by
614 Multiplexed Linear Amplification. *Cell Rep.* **2**, 666–673 (2012).
- 615 15. S. Picelli *et al.*, Smart-seq2 for sensitive full-length transcriptome profiling in single cells.
616 *Nat. Methods.* **10**, 1096–1098 (2013).

- 617 16. Y. Sasagawa *et al.*, Quartz-Seq: a highly reproducible and sensitive single-cell RNA
618 sequencing method, reveals non-genetic gene-expression heterogeneity. *Genome Biol.* **14**,
619 R31 (2013).
- 620 17. R. V. Grindberg *et al.*, RNA-sequencing from single nuclei. *Proc. Natl. Acad. Sci.* **110**,
621 19802–19807 (2013).
- 622 18. H. C. Fan, G. K. Fu, S. P. A. Fodor, Combinatorial labeling of single cells for gene
623 expression cytometry. *Science.* **347**, 1258367 (2015).
- 624 19. E. Z. Macosko *et al.*, Highly Parallel Genome-wide Expression Profiling of Individual Cells
625 Using Nanoliter Droplets. *Cell.* **161**, 1202–1214 (2015).
- 626 20. A. M. Klein *et al.*, Droplet Barcoding for Single-Cell Transcriptomics Applied to
627 Embryonic Stem Cells. *Cell.* **161**, 1187–1201 (2015).
- 628 21. A. A. Kolodziejczyk, J. K. Kim, V. Svensson, J. C. Marioni, S. A. Teichmann, The
629 Technology and Biology of Single-Cell RNA Sequencing. *Mol. Cell.* **58**, 610–620 (2015).
- 630 22. S. Liu, C. Trapnell, Single-cell transcriptome sequencing: recent advances and remaining
631 challenges. *F1000Research.* **5** (2016), doi:10.12688/f1000research.7223.1.
- 632 23. A. Adey *et al.*, In vitro, long-range sequence information for de novo genome assembly via
633 transposase contiguity. *Genome Res.* **24**, 2041–2049 (2014).
- 634 24. S. Amini *et al.*, Haplotype-resolved whole-genome sequencing by contiguity-preserving
635 transposition and combinatorial indexing. *Nat. Genet.* **46**, 1343–1349 (2014).
- 636 25. D. A. Cusanovich *et al.*, Multiplex single cell profiling of chromatin accessibility by
637 combinatorial cellular indexing. *Science.* **348**, 910–914 (2015).
- 638 26. S. A. Vitak *et al.*, Sequencing thousands of single-cell genomes with combinatorial
639 indexing. *Nat. Methods.* **advance online publication** (2017), doi:10.1038/nmeth.4154.
- 640 27. V. Ramani *et al.*, Massively multiplex single-cell Hi-C. *Nat. Methods.* **advance online**
641 **publication** (2017), doi:10.1038/nmeth.4155.
- 642 28. J. Alles *et al.*, Cell fixation and preservation for droplet-based single-cell transcriptomics.
643 *bioRxiv*, 99473 (2017).
- 644 29. V. Ntranos, G. M. Kamath, J. M. Zhang, L. Pachter, D. N. Tse, Fast and accurate single-cell
645 RNA-seq analysis by clustering of transcript-compatibility counts. *Genome Biol.* **17**, 112
646 (2016).
- 647 30. J. Gertz *et al.*, Transposase mediated construction of RNA-seq libraries. *Genome Res.* **22**,
648 134–141 (2012).

- 649 31. J. E. Sulston, E. Schierenberg, J. G. White, J. N. Thomson, The embryonic cell lineage of
650 the nematode *Caenorhabditis elegans*. *Dev. Biol.* **100**, 64–119 (1983).
- 651 32. J. E. Sulston, H. R. Horvitz, Post-embryonic cell lineages of the nematode, *Caenorhabditis*
652 *elegans*. *Dev. Biol.* **56**, 110–156 (1977).
- 653 33. G. Heimberg, R. Bhatnagar, H. El-Samad, M. Thomson, Low Dimensionality in Gene
654 Expression Data Enables the Accurate Extraction of Transcriptional Programs from
655 Shallow Sequencing. *Cell Syst.* **2**, 239–250 (2016).
- 656 34. M. E. Boeck *et al.*, *Genome Res.*, in press, doi:10.1101/gr.202663.115.
- 657 35. P.-Y. Tung *et al.*, Batch effects and the effective design of single-cell gene expression
658 studies. *bioRxiv*, 62919 (2016).
- 659 36. S. Zhang, D. Banerjee, J. R. Kuhn, Isolation and Culture of Larval Cells from *C. elegans*.
660 *PLOS ONE*. **6**, e19505 (2011).
- 661 37. J. D. Buenrostro, P. G. Giresi, L. C. Zaba, H. Y. Chang, W. J. Greenleaf, Transposition of
662 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-
663 binding proteins and nucleosome position. *Nat. Methods*. **10**, 1213–1218 (2013).
- 664 38. A. Dobin *et al.*, STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. **29**, 15–21
665 (2013).
- 666 39. A. Adey *et al.*, The haplotype-resolved genome and epigenome of the aneuploid HeLa
667 cancer cell line. *Nature*. **500**, 207–211 (2013).
- 668 40. N. Habib *et al.*, Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn
669 neurons. *Science*. **353**, 925–928 (2016).
- 670 41. A. Rodriguez, A. Laio, Clustering by fast search and find of density peaks. *Science*. **344**,
671 1492–1496 (2014).
- 672 42. D. Moerman, Sarcomere assembly in *C. elegans* muscle. *WormBook* (2006),
673 doi:10.1895/wormbook.1.81.1.
- 674 43. A. A. Beg, E. M. Jorgensen, EXP-1 is an excitatory GABA-gated cation channel. *Nat.*
675 *Neurosci.* **6**, 1145–1152 (2003).
- 676 44. L. Tilleman *et al.*, An N-Myristoylated Globin with a Redox-Sensing Function That
677 Regulates the Defecation Cycle in *Caenorhabditis elegans*. *PLOS ONE*. **7**, e48768 (2012).
- 678 45. J. P. Ardizzi, H. F. Epstein, Immunochemical localization of myosin heavy chain isoforms
679 and paramyosin in developmentally and structurally diverse muscle cell types of the
680 nematode *Caenorhabditis elegans*. *J. Cell Biol.* **105**, 2763–2770 (1987).

- 681 46. M. Köppen *et al.*, Cooperative regulation of AJM-1 controls junctional integrity in
682 *Caenorhabditis elegans* epithelia. *Nat. Cell Biol.* **3**, 983–991 (2001).
- 683 47. L. McMahon, R. Legouis, J.-L. Vonesch, M. Labouesse, Assembly of *C. elegans* apical
684 junctions involves positioning and compaction by LET-413 and protein aggregation by the
685 MAGUK protein DLG-1. *J. Cell Sci.* **114**, 2265–2277 (2001).
- 686 48. C. McKeown, V. Praitis, J. Austin, *sma-1* encodes a betaH-spectrin homolog required for
687 *Caenorhabditis elegans* morphogenesis. *Dev. Camb. Engl.* **125**, 2087–2098 (1998).
- 688 49. A. Hrus *et al.*, *C. elegans* Agrin Is Expressed in Pharynx, IL1 Neurons and Distal Tip Cells
689 and Does Not Genetically Interact with Genes Involved in Synaptogenesis or Muscle
690 Function. *PLOS ONE*. **2**, e731 (2007).
- 691 50. V. Ghai, R. B. Smit, J. Gaudet, Transcriptional regulation of HLH-6-independent and
692 subtype-specific genes expressed in the *Caenorhabditis elegans* pharyngeal glands. *Mech.*
693 *Dev.* **129**, 284–297 (2012).
- 694 51. C. Thacker, J. A. Sheps, A. M. Rose, *Caenorhabditis elegans* *dpy-5* is a cuticle procollagen
695 processed by a proprotein convertase. *Cell. Mol. Life Sci. CMLS.* **63**, 1193–1204 (2006).
- 696 52. L. Hao, R. Johnsen, G. Lauter, D. Baillie, T. R. Bürglin, Comprehensive analysis of gene
697 expression patterns of hedgehog-related genes. *BMC Genomics.* **7**, 280 (2006).
- 698 53. J. B. Rand, Acetylcholine. *WormBook Online Rev. C Elegans Biol.*, 1–21 (2007).
- 699 54. E. Efimenko *et al.*, *Caenorhabditis elegans* DYF-2, an Orthologue of Human WDR19, Is a
700 Component of the Intraflagellar Transport Machinery in Sensory Cilia. *Mol. Biol. Cell.* **17**,
701 4801–4811 (2006).
- 702 55. E. M. Jorgensen, GABA. *WormBook Online Rev. C Elegans Biol.*, 1–13 (2005).
- 703 56. Y. Zhang *et al.*, Identification of genes expressed in *C. elegans* touch receptor neurons.
704 *Nature.* **418**, 331–335 (2002).
- 705 57. P. J. Brockie, D. M. Madsen, Y. Zheng, J. Mellem, A. V. Maricq, Differential expression of
706 glutamate receptor subunits in the nervous system of *Caenorhabditis elegans* and their
707 regulation by the homeodomain protein UNC-42. *J. Neurosci. Off. J. Soc. Neurosci.* **21**,
708 1510–1522 (2001).
- 709 58. T. Janssen *et al.*, Discovery of a Cholecystokinin-Gastrin-Like Signaling System in
710 Nematodes. *Endocrinology.* **149**, 2826–2839 (2008).
- 711 59. M. Treinin, M. Chalfie, A mutated acetylcholine receptor subunit causes neuronal
712 degeneration in *C. elegans*. *Neuron.* **14**, 871–877 (1995).

- 713 60. M. Treinin, B. Gillo, L. Liebman, M. Chalfie, Two functionally dependent acetylcholine
714 subunits are encoded in a single *Caenorhabditis elegans* operon. *Proc. Natl. Acad. Sci. U. S.*
715 *A.* **95**, 15492–15495 (1998).
- 716 61. M. J. Alkema, M. Hunter-Ensor, N. Ringstad, H. R. Horvitz, Tyramine Functions
717 independently of octopamine in the *Caenorhabditis elegans* nervous system. *Neuron.* **46**,
718 247–260 (2005).
- 719 62. L. S. Nelson, M. L. Rosoff, C. Li, Disruption of a neuropeptide gene, *flp-1*, causes multiple
720 behavioral defects in *Caenorhabditis elegans*. *Science.* **281**, 1686–1690 (1998).
- 721 63. H. C. Korswagen, A. M. van der Linden, R. H. Plasterk, G protein hyperactivation of the
722 *Caenorhabditis elegans* adenylyl cyclase *SGS-1* induces neuronal degeneration. *EMBO J.*
723 **17**, 5059–5065 (1998).
- 724 64. D. Combes, Y. Fedon, J.-P. Toutant, M. Arpagaus, Multiple *ace* genes encoding
725 acetylcholinesterases of *Caenorhabditis elegans* have distinct tissue expression. *Eur. J.*
726 *Neurosci.* **18**, 497–512 (2003).
- 727 65. G. Jafari *et al.*, Genetics of Extracellular Matrix Remodeling During Organ Growth Using
728 the *Caenorhabditis elegans* Pharynx Model. *Genetics.* **186**, 969–982 (2010).
- 729 66. M. Harterink *et al.*, Neuroblast migration along the anteroposterior axis of *C. elegans* is
730 controlled by opposing gradients of Wnts and a secreted Frizzled-related protein. *Dev.*
731 *Camb. Engl.* **138**, 2915–2924 (2011).
- 732 67. R. J. Hobson *et al.*, *SER-7*, a *Caenorhabditis elegans* 5-HT7-like Receptor, Is Essential for
733 the 5-HT Stimulation of Pharyngeal Pumping and Egg Laying. *Genetics.* **172**, 159–169
734 (2006).
- 735 68. M. Furuya, H. Qadota, A. D. Chisholm, A. Sugimoto, The *C. elegans* eyes absent ortholog
736 *EYA-1* is required for tissue differentiation and plays partially redundant roles with *PAX-6*.
737 *Dev. Biol.* **286**, 452–463 (2005).
- 738 69. A. Oishi *et al.*, *FLR-2*, the glycoprotein hormone alpha subunit, is involved in the neural
739 control of intestinal functions in *Caenorhabditis elegans*. *Genes Cells.* **14**, 1141–1154
740 (2009).
- 741 70. L. Emtage, G. Gu, E. Hartwig, M. Chalfie, Extracellular proteins organize the
742 mechanosensory channel complex in *C. elegans* touch receptor neurons. *Neuron.* **44**, 795–
743 807 (2004).
- 744 71. X. Wang *et al.*, The *C. elegans* L1CAM homologue *LAD-2* functions as a coreceptor in
745 *MAB-20/Sema2*-mediated axon guidance. *J. Cell Biol.* **180**, 233–246 (2008).
- 746 72. S. Yu, L. Avery, E. Baude, D. L. Garbers, Guanylyl cyclase expression in specific sensory
747 neurons: A new family of chemosensory receptors. *Proc. Natl. Acad. Sci.* **94**, 3384–3387
748 (1997).

- 749 73. J. M. Gray *et al.*, Oxygen sensation and social feeding mediated by a *C. elegans* guanylate
750 cyclase homologue. *Nature*. **430**, 317–322 (2004).
- 751 74. K. Kim, C. Li, Expression and regulation of an FMRFamide-related neuropeptide gene
752 family in *Caenorhabditis elegans*. *J. Comp. Neurol.* **475**, 540–550 (2004).
- 753 75. P. W. McDonald *et al.*, Vigorous motor activity in *Caenorhabditis elegans* requires efficient
754 clearance of dopamine mediated by synaptic localization of the dopamine transporter DAT-
755 1. *J. Neurosci. Off. J. Soc. Neurosci.* **27**, 14216–14227 (2007).
- 756 76. R. Lints, S. W. Emmons, Patterning of dopaminergic neurotransmitter identity among
757 *Caenorhabditis elegans* ray sensory neurons by a TGFbeta family signaling pathway and a
758 Hox gene. *Dev. Camb. Engl.* **126**, 5819–5831 (1999).
- 759 77. T. C. Jacob, J. M. Kaplan, The EGL-21 carboxypeptidase E facilitates acetylcholine release
760 at *Caenorhabditis elegans* neuromuscular junctions. *J. Neurosci. Off. J. Soc. Neurosci.* **23**,
761 2122–2130 (2003).
- 762 78. D. Sieburth *et al.*, Systematic analysis of genes required for synapse structure and function.
763 *Nature*. **436**, 510–517 (2005).
- 764 79. T. Cai, M. W. Krause, W. F. Odenwald, R. Toyama, A. L. Notkins, The IA-2 gene family:
765 homologs in *Caenorhabditis elegans*, *Drosophila* and zebrafish. *Diabetologia*. **44**, 81–88
766 (2001).
- 767 80. D. D. Ikeda *et al.*, CASY-1, an ortholog of calsynenins/alcadeins, is essential for learning
768 in *Caenorhabditis elegans*. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 5260–5265 (2008).
- 769 81. H. M. Zhou, I. Brust-Mascher, J. M. Scholey, Direct visualization of the movement of the
770 monomeric axonal transport motor UNC-104 along neuronal processes in living
771 *Caenorhabditis elegans*. *J. Neurosci. Off. J. Soc. Neurosci.* **21**, 3749–3755 (2001).
- 772 82. G. P. Mullen *et al.*, UNC-41/stonin functions with AP2 to recycle synaptic vesicles in
773 *Caenorhabditis elegans*. *PloS One*. **7**, e40095 (2012).
- 774 83. T. Bacaj, M. Tevlin, Y. Lu, S. Shaham, Glia are essential for sensory organ function in *C.*
775 *elegans*. *Science*. **322**, 744–747 (2008).
- 776 84. A. Yoshida *et al.*, A glial K(+)/Cl(-) cotransporter modifies temperature-evoked dynamics
777 in *Caenorhabditis elegans* sensory neurons. *Genes Brain Behav.* **15**, 429–440 (2016).
- 778 85. R. Y. Yu, C. Q. Nguyen, D. H. Hall, K. L. Chow, Expression of ram-5 in the structural cell
779 is required for sensory ray morphogenesis in *Caenorhabditis elegans* male tail. *EMBO J.* **19**,
780 3542–3555 (2000).
- 781 86. E. A. Perens, S. Shaham, *C. elegans* daf-6 encodes a patched-related protein required for
782 lumen formation. *Dev. Cell.* **8**, 893–906 (2005).

- 783 87. D. Levitan, I. Greenwald, LIN-12 protein expression and localization during vulval
784 development in *C. elegans*. *Dev. Camb. Engl.* **125**, 3101–3109 (1998).
- 785 88. T. A. Starich, D. H. Hall, D. Greenstein, Two classes of gap junction channels mediate
786 soma-germline interactions essential for germline proliferation and gametogenesis in
787 *Caenorhabditis elegans*. *Genetics*. **198**, 1127–1153 (2014).
- 788 89. B. D. Ackley *et al.*, The basement membrane components nidogen and type XVIII collagen
789 regulate organization of neuromuscular junctions in *Caenorhabditis elegans*. *J. Neurosci.*
790 *Off. J. Soc. Neurosci.* **23**, 3577–3587 (2003).
- 791 90. R. P. Johnson, S. H. Kang, J. M. Kramer, *C. elegans* dystroglycan DGN-1 functions in
792 epithelia and neurons, but not muscle, and independently of dystrophin. *Dev. Camb. Engl.*
793 **133**, 1911–1921 (2006).
- 794 91. E. J. Cram, H. Shang, J. E. Schwarzbauer, A systematic RNA interference screen reveals a
795 cell migration gene network in *C. elegans*. *J. Cell Sci.* **119**, 4811–4818 (2006).
- 796 92. S. H. Kang, J. M. Kramer, Nidogen is nonessential and not required for normal type IV
797 collagen localization in *Caenorhabditis elegans*. *Mol. Biol. Cell.* **11**, 3911–3923 (2000).
- 798 93. H. Komatsu *et al.*, OSM-11 facilitates LIN-12 Notch signaling during *Caenorhabditis*
799 *elegans* vulval development. *PLoS Biol.* **6**, e196 (2008).
- 800 94. C. W. Whitfield, C. Bénard, T. Barnes, S. Hekimi, S. K. Kim, Basolateral localization of
801 the *Caenorhabditis elegans* epidermal growth factor receptor in epithelial cells by the PDZ
802 protein LIN-10. *Mol. Biol. Cell.* **10**, 2087–2100 (1999).
- 803 95. S. A. Kostas, A. Fire, The T-box factor MLS-1 acts as a molecular switch during
804 specification of nonstriated muscle in *C. elegans*. *Genes Dev.* **16**, 257–269 (2002).
- 805 96. I. Kawasaki *et al.*, PGL-1, a predicted RNA-binding component of germ granules, is
806 essential for fertility in *C. elegans*. *Cell.* **94**, 635–645 (1998).
- 807 97. R. E. Navarro, E. Y. Shim, Y. Kohara, A. Singson, T. K. Blackwell, *cgh-1*, a conserved
808 predicted RNA helicase required for gametogenesis and protection from physiological
809 germline apoptosis in *C. elegans*. *Dev. Camb. Engl.* **128**, 3221–3232 (2001).
- 810 98. A. R. Jones, R. Francis, T. Schedl, GLD-1, a cytoplasmic protein essential for oocyte
811 differentiation, shows stage- and sex-specific expression during *Caenorhabditis elegans*
812 germline development. *Dev. Biol.* **180**, 165–183 (1996).
- 813 99. M. Hanazawa *et al.*, The *Caenorhabditis elegans* eukaryotic initiation factor 5A homologue,
814 IFF-1, is required for germ cell proliferation, gametogenesis and localization of the P-
815 granule component PGL-1. *Mech. Dev.* **121**, 213–224 (2004).
- 816 100. G. Aspöck, H. Kagoshima, G. Niklaus, T. R. Bürglin, *Caenorhabditis elegans* has scores of
817 hedgehog-related genes: sequence and expression analysis. *Genome Res.* **9**, 909–923 (1999).

- 818 101. A. Hahn-Windgassen, M. R. Van Gilst, The Caenorhabditis elegans HNF4alpha Homolog,
819 NHR-31, mediates excretory tube growth and function through coordinate regulation of the
820 vacuolar ATPase. *PLoS Genet.* **5**, e1000553 (2009).
- 821 102. A. Karabinos, H. Schmidt, J. Harborth, R. Schnabel, K. Weber, Essential roles for four
822 cytoplasmic intermediate filament proteins in Caenorhabditis elegans development. *Proc.*
823 *Natl. Acad. Sci. U. S. A.* **98**, 7863–7868 (2001).
- 824 103. A. Patton *et al.*, Endocytosis function of a ligand-gated ion channel homolog in
825 Caenorhabditis elegans. *Curr. Biol. CB.* **15**, 1045–1050 (2005).
- 826 104. P. M. Loria, J. Hodgkin, O. Hobert, A conserved postsynaptic transmembrane protein
827 affecting neuromuscular signaling in Caenorhabditis elegans. *J. Neurosci. Off. J. Soc.*
828 *Neurosci.* **24**, 2191–2201 (2004).

829

830 **Acknowledgements**

831

832 We thank members of the Shendure, Trapnell and Waterston labs for helpful discussions; D.
833 Prunkard, and L. Gitari in the Pathology Flow Cytometry Core Facility for their exceptional
834 assistance in flow sorting; J. Rose, D. Maly, L. VandenBosch, and T. Reh lab for sharing the
835 NIH/3T3 cell line. HeLa S3 cells were used as part of this study. Henrietta Lacks, and the HeLa
836 cell line that was established from her tumor cells in 1951, have made significant contributions to
837 scientific progress and advances in human health. We are grateful to Henrietta Lacks, now
838 deceased, and to her surviving family members for their contributions to biomedical research.
839 This work was funded by grants from the NIH (DP1HG007811 and R01HG006283 to JS;
840 U41HG007355 and R01GM072675 to RHW) and the W. M. Keck Foundation (to CT and JS).
841 DAC was supported in part by T32HL007828 from the National Heart, Lung, and Blood Institute.
842 JS is an Investigator of the Howard Hughes Medical Institute.

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861 **Figure Legends**

862

863 **Figure 1: sci-RNAseq enables massively multiplexed single cell transcriptome profiling.** a.)
864 Schematic of the sci-RNA-seq workflow. Methanol-fixed cells or unfixed nuclei are split to one
865 or more 96-well or 384-well plates for reverse transcription with different barcodes (first round
866 of barcoding) in each well. Cells from different wells are pooled together and flow-sorted into
867 one or more 96-well or 384-well plates for second strand synthesis, tagmentation and PCR with
868 well-specific barcode combinations (second round of barcoding). The resulting PCR amplicons
869 are pooled and deep sequenced to generate single cell 3' digital gene expression profiles. b.) sci-
870 RNA-seq library amplicons include Illumina adapters, PCR indices (i5 and i7), a reverse
871 transcription barcode and UMI, in addition to the cDNA fragment to be sequenced. Read 1
872 covers the reverse transcription barcode (10 bp) and unique molecular identifier (UMI, 8 bp).
873 Read 2 covers the cDNA fragment. The combination of the PCR indices and the reverse
874 transcription barcode define a cellular index. c.) Scatter plot of unique human and mouse UMI
875 counts from a 384 x 384 sci-RNA-seq experiment. This 384-well experiment included multiple
876 different mixtures of cells (see Methods), but only cells originating from a well containing mixed
877 human (HEK293T or HeLa S3) and mouse (NIH/3T3) cells during the first round of barcoding
878 are plotted here. Inferred mouse cells are colored in blue; inferred human cells are colored in red,
879 and "collisions" are colored in grey.

880

881 **Figure 2: Representative results from an optimized protocol for sci-RNA-seq.** a.) Scatter
882 plot of unique human and mouse cell UMI counts from a 96-well sci-RNA-seq experiment. This
883 96-well experiment included multiple different mixtures of cells (see Methods), but only cells
884 originating from a mixture of human (HEK293T) and mouse (NIH/3T3) are plotted here.
885 Inferred mouse cells are colored in blue; inferred human cells are colored in red, and "collisions"
886 are colored in grey. b.-c.) Boxplots showing the number of UMIs (b) and genes (c) detected per
887 cell in interspecies mixing experiments. d.) tSNE plot of human cell line (HEK293T and HeLa
888 S3) mixtures and pure populations. Cells originating in wells containing pure HEK293T (red),
889 pure HeLa S3 (yellow) or a mixture of the two (pink) were all clustered together with tSNE. e.)
890 Correlation between gene expression measurements from aggregated sci-RNA-seq data vs. bulk
891 RNA-seq data, together with a linear regression line (red) and $y=x$ line (black).

892

893 **Figure 3: sci-RNA-seq is compatible with isolated nuclei as starting material.** a.) Scatter plot
894 of unique human and mouse nuclei UMI counts from a 96-well sci-RNA-seq experiment. This
895 96-well experiment included multiple different mixtures of cells (see Methods), but only cells
896 originating from a mixture of human (HEK293T) and mouse (NIH/3T3) nuclei are plotted here.
897 Inferred murine cells are colored in blue; inferred human cells are colored in red, and "collisions"
898 are colored in grey. b.-c.) Boxplots showing the number of UMIs (b) and genes (c) detected per
899 cell in nuclear sci-RNA-seq experiments. d.) Correlation between gene expression measurements
900 in aggregated sci-RNA-seq profiles of NIH/3T3 cells vs. NIH/3T3 nuclei, together with a linear
901 regression line (red) and $y=x$ line (black).

902

903 **Figure 4: A single sciRNA-seq experiment provides a near-comprehensive view of the**
904 **single cell transcriptomes comprising the *C. elegans* larva.** a) t-SNE visualization of the high-
905 level cell types identified. b) t-SNE visualization of neuronal subtypes. Dimensionality reduction
906 and clustering was applied to cells in neuronal clusters shown in (a). c) Bar plot showing the

907 proportion of somatic cells profiled with sci-RNA-seq that could be identified as belonging to
908 each cell type compared to the proportion of cells from that type present in an L2 *C. elegans*
909 individual. d) Scatter plot showing the log-scaled transcripts per million (TPM) of genes in the
910 aggregation of all sci-RNA-seq reads (x axis) or in bulk RNA-seq (y axis; geometric mean of 3
911 experiments), e) Heatmap showing the relative expression of genes in consensus transcriptomes
912 for each cell type estimated by sci-RNA-seq. Genes are included if they have a size-factor-
913 normalized mean expression of ≥ 0.05 in at least one cell type (8,318 genes in total). The raw
914 expression data (UMI count matrix) is log-transformed, column centered and scaled (using the R
915 function scale), and the resulting values are clamped to the interval [-2, 2].

916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934

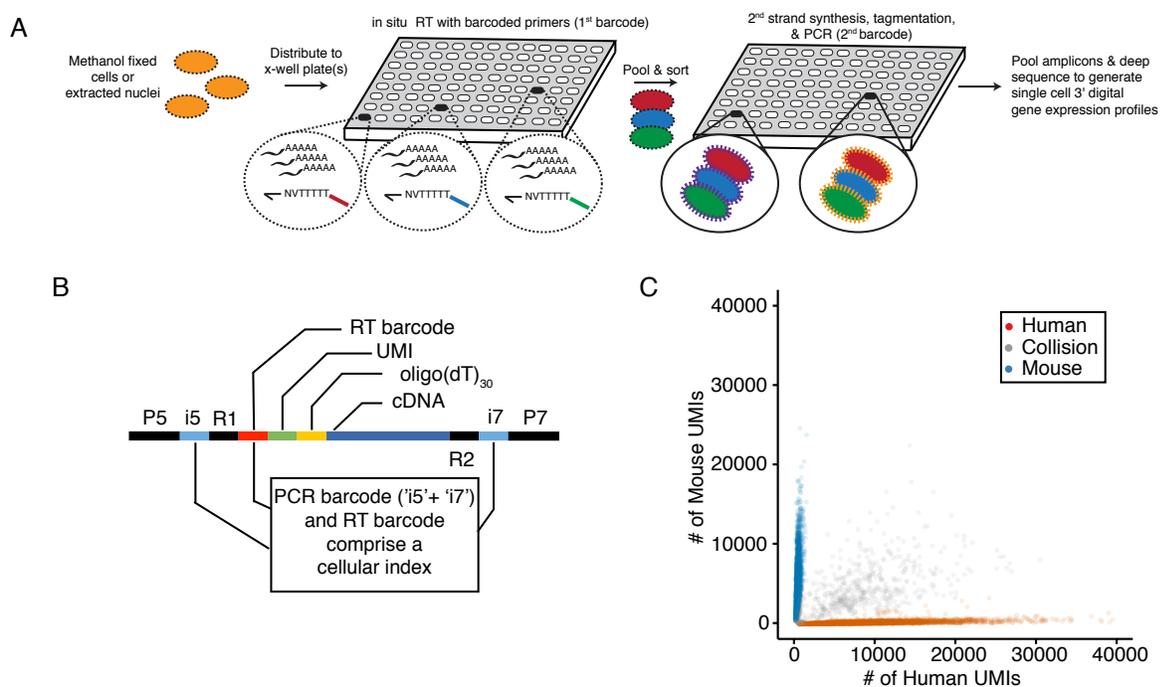


Figure 1

935

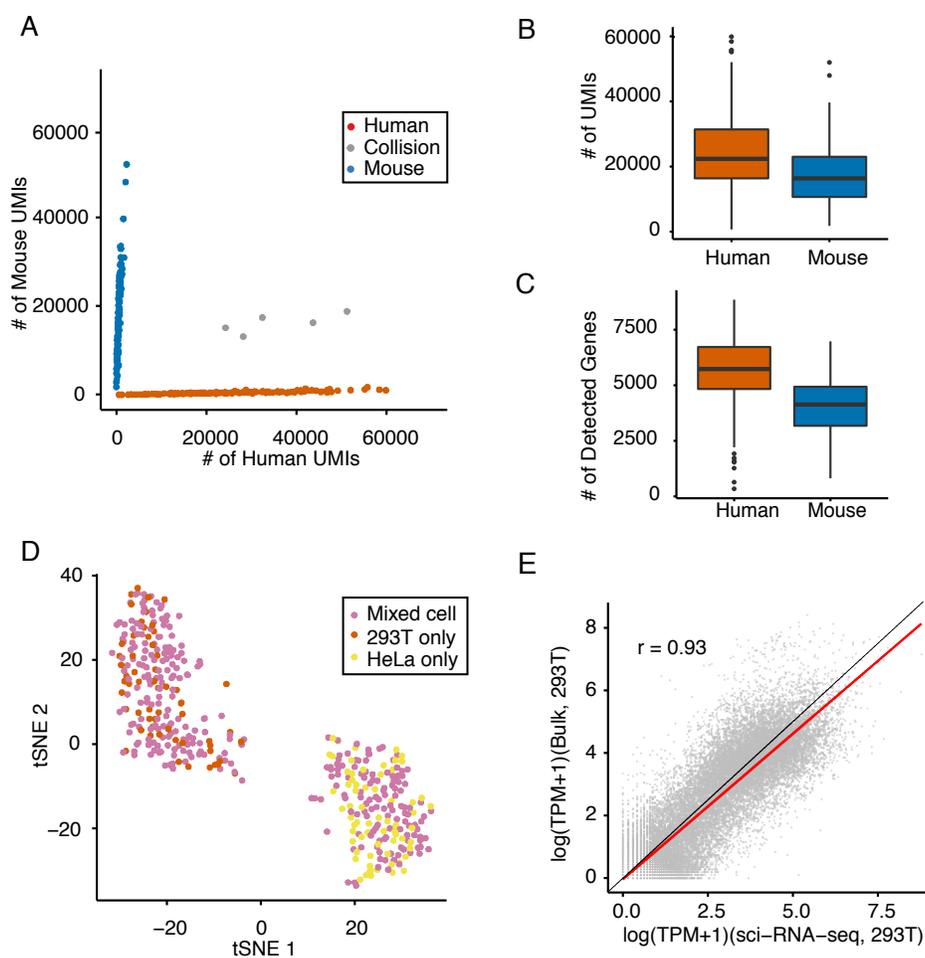


Figure 2

936

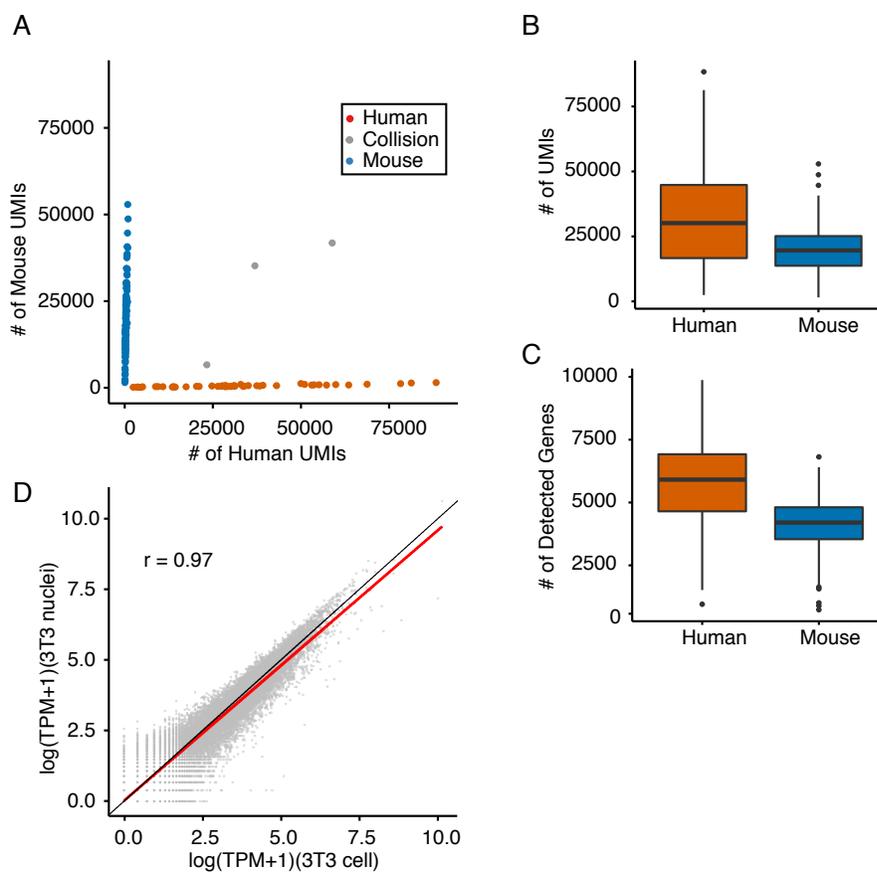


Figure 3

937
938

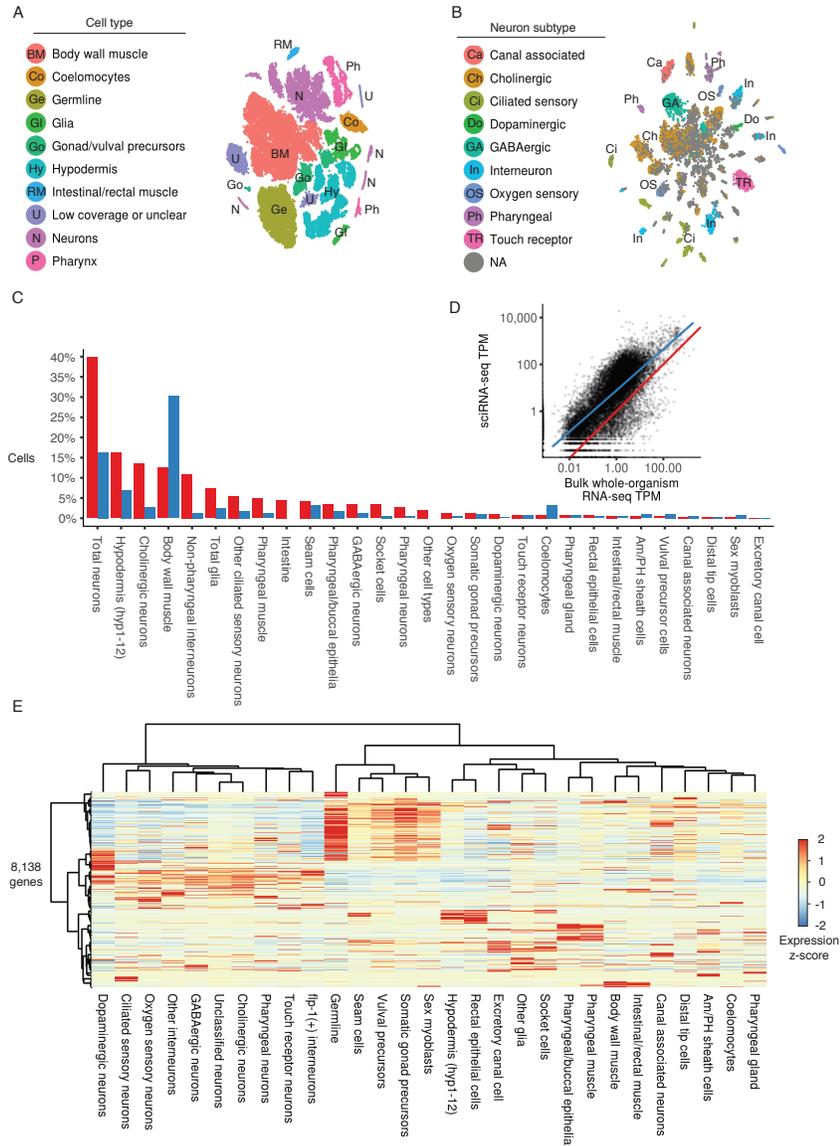


Figure 4

939

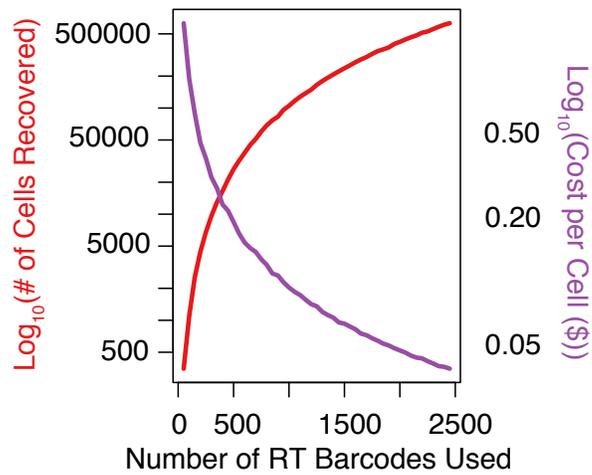


Figure S1 | Combinatorial indexing with increasing numbers of reverse-transcription barcodes enables sublinear scaling of cost per cell. Plot showing how detection capacity (i.e. the number of cells detected in a sci-RNA-seq experiment, red) and cost per cell (blue) vary as a function of the number of cellular indices used, assuming a collision rate of 5%.

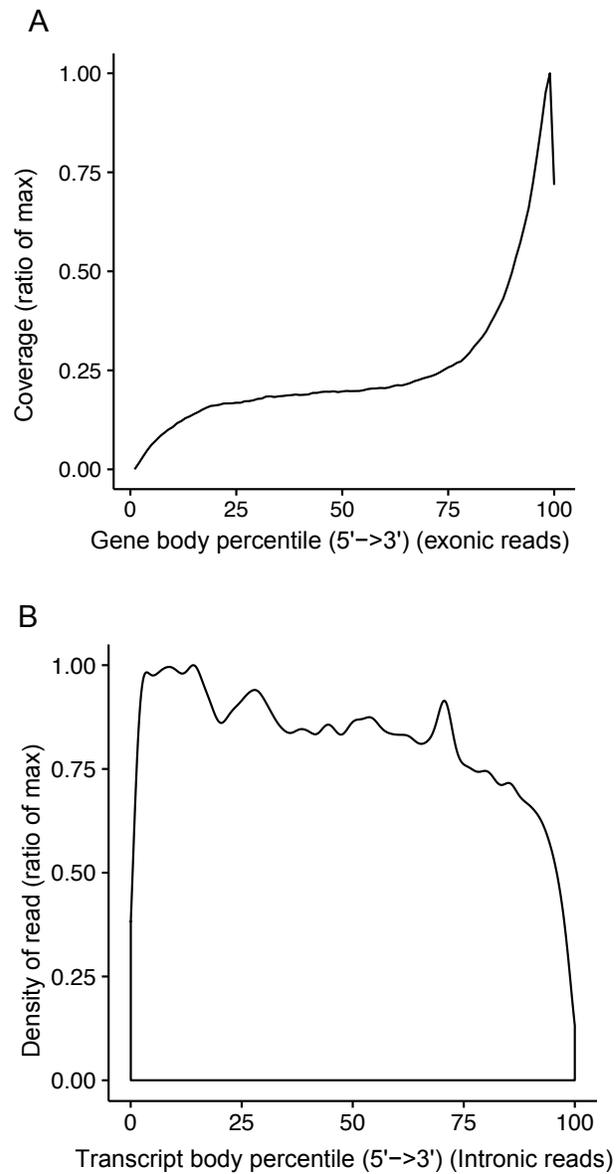


Figure S2 | Metagene plot for unique intronic and exonic sci-RNA-seq reads. A.)

Sci-RNA-seq on fixed cells demonstrates 3'-biased exonic coverage along gene body (intronic region excluded). B.) Density plot for the intronic reads number mapping to different percentiles of transcript body (intronic region included). y-axis is scaled to the ratio of max.

941
942
943

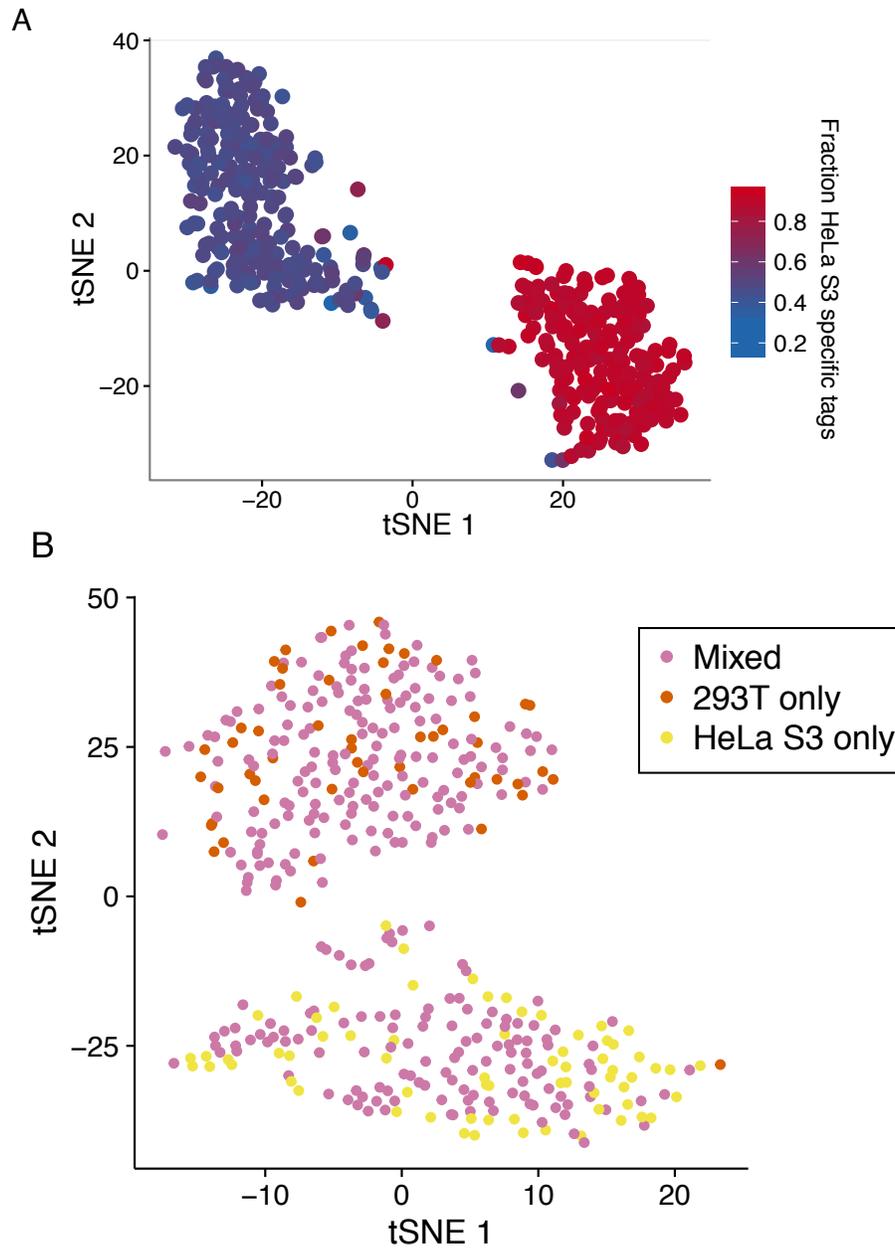


Figure S3 | Quality control for sci-RNA-seq experiments using synthetic mixtures of HeLa S3 and HEK293T cell lines. A.) tSNE plot (as in Figure 2d), with cells colored by fraction of reads harboring HeLa S3 specific SNVs relative to hg19 assembly. B.) tSNE using digital gene expression matrices constructed from intronic reads only. Cells are colored by the population of cells from which they derived, with pure HEK293T in red, pure HeLa S3 in yellow, and mixed HEK293T+HeLa S3 in pink

944
945
946

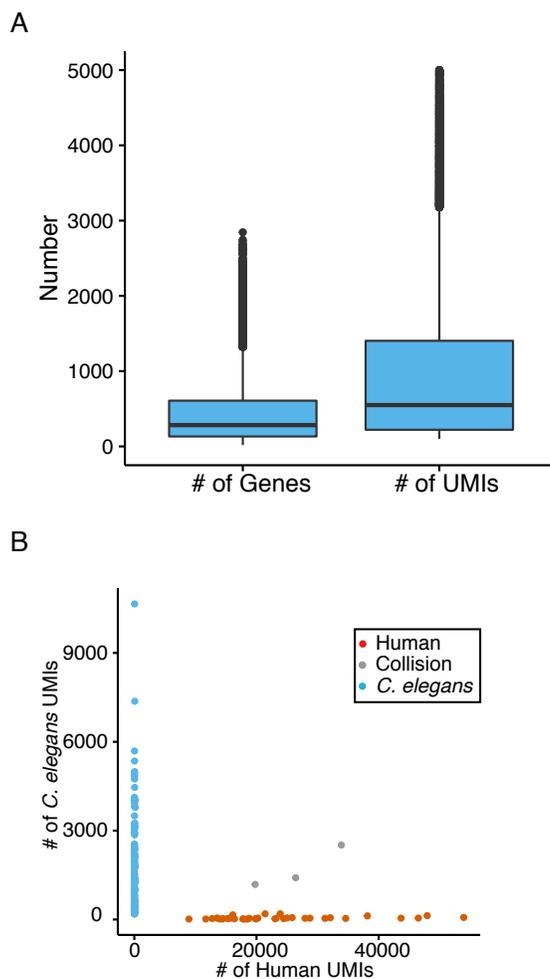


Figure S4 | Quality control metrics for *C. elegans* sci-RNA-seq experiments. A.) Distribution of number of protein coding genes and UMI counts (mapping to protein coding genes) detected per *C. elegans* cell. B.) Scatter plot of unique human and *C. elegans* cell UMI counts from a sci-RNA-seq experiment performed on mixture of HEK293T (human) and *C. elegans* cells.

947
948
949
950

951 **Table S1. Summary of experiments**

952

Experiment ID	Technique version	# of first round barcodes	Cell populations barcoded (# of wells during first round of barcoding)	# of second round barcodes	cells sorted per well
1	Initial protocol	384	Pure HEK293T cells (32) Pure HeLa S3 cells (32) Pure NIH/3T3 cells (32) HEK293T & NIH/3T3 cells (96) HeLa S3 & NIH/3T3 cells (96) HEK293T & HeLa S3 cells (96)	384	50
2	Optimized protocol	96	Pure HEK293T cells (8) Pure HeLa S3 cells (8) HEK293T & NIH/3T3 cells (24) HEK293T & HeLa S3 cells (24) HEK293T & NIH/3T3 nuclei (24)	96	12
3	Optimized protocol	576	<i>C. elegans</i> & HEK293T cells (576)	960	96% of wells: 140 <i>C. elegans</i> cells 4% of wells: 140 <i>C. elegans</i> cells + 10 HEK293T cells

953

954

955 **Table S2: Marker genes for cell types identified in sci-RNA-seq *C. elegans* data.** A cell is
 956 assigned to a cell type (i.e. used when constructing a consensus expression profile for that cell
 957 type) if it is in a t-SNE cluster enriched for expression of marker genes listed for that cell type
 958 and the individual cell expresses at least one of those marker genes (or ≥ 2 from a set of marker
 959 genes where listed).

Cell type	# cells assigned	# genes expressed (>2% of cells)	Marker genes	References
Body wall muscle	11,389	4,321	≥ 2 of <i>lev-11</i> , <i>myo-3</i> , <i>ttn-1</i> , <i>unc-54</i>	(42)
Intestinal/rectal muscle	165	4,561	<i>exp-1</i> , <i>glb-26</i>	(43, 44)
Pharyngeal muscle	451	3,332	<i>myo-2</i> , <i>myo-1</i> , <i>myo-5</i> , <i>tnt-4</i>	(45)
Pharyngeal/buccal epithelia	678	2,781	<i>ajm-1</i> , <i>dlg-1</i> , <i>sma-1</i> , <i>agr-1</i> + absence of <i>myo-1/2/5</i> and <i>tnt-4</i>	(46–49)
Pharyngeal gland	238	3,472	<i>phat-4</i> , <i>phat-2</i> , <i>phat-6</i>	(50)
Hypodermis (hyp1-12)	2,646	4,428	<i>dpy-4</i> , <i>dpy-5</i> , <i>dpy-13</i>	(51)
Seam cells	1,231	4,842	<i>grd-13</i> , <i>grd-10</i> , <i>grd-3</i> and absence of <i>dpy-4/5/13</i>	(52)
Cholinergic neurons	1,038	2,265	<i>unc-17</i> , <i>cho-1</i> , <i>cha-1</i> , <i>acr-18</i> , <i>acr-15</i>	(53)
Ciliated sensory neurons (excluding dopaminergic neurons)	587	2,308	<i>R102.2</i> , <i>dyf-2</i> , <i>che-3</i> , <i>nphp-4</i>	(54)
GABAergic neurons	435	1,810	<i>unc-25</i>	(55)
Touch receptor neurons	321	1,703	<i>mec-17</i> , <i>mec-7</i>	(56)
Other interneurons (excluding flp-1(+) interneurons)	263	2,499	<i>glr-3</i> (RIA), <i>nlp-12</i> (DVA), Both <i>des-2</i> and <i>deg-3</i> (PVC/PVD), or <i>tbh-1</i> (RIC)	(57–61)
flp-1(+) interneurons	211	1,484	<i>flp-1</i>	(62)
Canal associated neurons	187	2,241	≥ 2 of <i>acy-2</i> , <i>ace-3</i> , <i>mig-6</i> , <i>cwn-1</i>	(63–66)
Pharyngeal neurons	173	1,975	<i>ser-7</i> , <i>eya-1</i> , <i>flr-2</i> , <i>pha-4</i>	(67–69)
Oxygen sensory neurons	173	2,043	Both <i>mec-1</i> , <i>lad-2</i>	(70–74)

			(ALN/PLN/SDQ), or <i>gcy-32, gcy-36,</i> <i>gcy-37</i> (AQR/PQR/URX), or <i>flp-17</i> (BAG)	
Dopaminergic neurons	70	1,821	<i>dat-1, cat-2</i>	(75, 76)
Unclassified neurons	2690	1,895	<i>egl-21, sbt-1, ida-1,</i> <i>casy-1, unc-104, unc-41</i>	(77–82)
Amphid/phasmid sheath cells	338	4,143	<i>vap-1, fig-1</i>	(83)
Socket cells	236	3,192	<i>grl-1, grl-2, grd-15</i>	(52)
Unclassified glia	394	2,983	<i>kcc-3, ram-5, daf-6</i>	(84–86)
Somatic gonad precursors (excluding distal tip cells)	376	5,827	≥ 2 of <i>lin-12,</i> <i>inx-9, cle-1, dgn-1</i>	(87–90)
Distal tip cells	128	5,030	Both <i>mig-6, nid-1</i>	(91, 92)
Vulval precursor cells	377	5,069	≥ 2 of <i>lin-12, osm-11,</i> <i>let-23</i>	(87, 93, 94)
Sex myoblasts	311	4,466	<i>egl-15</i>	(95)
Germline	4301	5,086	≥ 2 of <i>pgl-1, cgh-1,</i> <i>gld-1, iff-1</i>	(96–99)
Rectal epithelial cells	168	4,490	<i>grd-1, grd-12</i>	(52, 100)
Excretory canal cell	37	5,057	Both of: <i>nhr-31, ifa-4</i>	(101, 102)
Coelomocytes	1232	1,991	<i>cup-4, lgc-26, unc-122</i>	(103, 104)
Unassigned	11192	NA	NA	NA

960