

Equivalent high-resolution identification of neuronal cell types with single-nucleus and single-cell RNA-sequencing

Trygve E. Bakken¹, Rebecca D. Hodge¹, Jeremy M. Miller¹, Zizhen Yao¹, Thuc N. Nguyen¹, Brian Aevermann², Eliza Barkan¹, Darren Bertagnolli¹, Tamara Casper¹, Nick Dee¹, Emma Garren¹, Jeff Goldy¹, Lucas T. Gray¹, Matthew Kroll¹, Roger S. Lasken², Kanan Lathia¹, Sheana Parry¹, Christine Rimorin¹, Richard H. Scheuermann², Nicholas J. Schork², Soraya I. Shehata¹, Michael Tieu¹, John W. Phillips¹, Amy Bernard¹, Kimberly A. Smith¹, Hongkui Zeng¹, Ed S. Lein¹, and Bosiljka Tasic¹

¹Allen Institute for Brain Science, Seattle, WA, USA

²J. Craig Venter Institute, La Jolla, CA, USA

January 20, 2018

1 Abstract

2 Transcriptional profiling of complex tissues by RNA-sequencing of single nuclei presents some advantages over whole cell
3 analysis. It enables unbiased cellular coverage, lack of cell isolation-based transcriptional effects, and application to archived
4 frozen specimens. Using a well-matched pair of single-nucleus RNA-seq (snRNA-seq) and single-cell RNA-seq (scRNA-seq)
5 SMART-Seq v4 datasets from mouse visual cortex, we demonstrate that similarly high-resolution clustering of closely related
6 neuronal types can be achieved with both methods if intronic sequences are included in nuclear RNA-seq analysis. More
7 transcripts are detected in individual whole cells (~11,000 genes) than nuclei (~7,000 genes), but the majority of genes have
8 similar detection across cells and nuclei. We estimate that the nuclear proportion of total cellular mRNA varies from 20% to
9 over 50% for large and small pyramidal neurons, respectively. Together, these results illustrate the high information content of
10 nuclear RNA for characterization of cellular diversity in brain tissues.

11 Introduction

12 Understanding neural circuits requires characterization of their cellular components. Cell types in mam-
13 malian brain have been defined based on shared morphological, electrophysiological and, more recently,
14 molecular properties (Poulin et al. 2016; Zeng and Sanes 2017; Bernard, Sorensen, and Lein 2009). scRNA-
15 seq has emerged as a high-throughput method for quantification of the majority of transcripts in thousands
16 of cells. scRNA-seq data have revealed diverse cell types in many mouse brain regions, including neocor-
17 tex (Tasic et al. 2016; Tasic et al. 2017; Zeisel et al. 2015), hypothalamus (Campbell et al. 2017), and
18 retina (Shekhar et al. 2016; Macosko et al. 2015).

19 However, scRNA-seq profiling does not provide an unbiased survey of neural cell types. Some cell types are
20 more vulnerable to the tissue dissociation process and are underrepresented in the final data set. For exam-
21 ple, in mouse neocortex, fast-spiking parvalbumin-positive interneurons and deep-projecting glutamatergic
22 neurons in layer 5b are observed in lower proportions than expected and need to be selectively enriched
23 using Cre-driver lines (Tasic et al. 2017) for sufficient sampling. In adult human neocortex, neurons largely

24 do not survive dissociation thereby causing over-representation of non-neuronal cells in single cell suspen-
25 sions (Darmanis et al. 2015). In contrast to whole cells, nuclei are more resistant to mechanical assaults and
26 can be isolated from frozen tissue (Krishnaswami et al. 2016; Lacar et al. 2016). Single nuclei have been
27 shown to provide sufficient gene expression information to define relatively broad cell classes in adult human
28 brain (Lake et al. 2016; Lake et al. 2017) and mouse hippocampus (Habib et al. 2016).

29 Previous studies have not addressed if the nucleus contains sufficient diversity and number of transcripts to
30 enable discrimination of closely related cell types at a resolution comparable to whole cells. A recent study
31 compared clustering results for single nuclei and whole cells isolated from mouse somatosensory cortex (Lake
32 et al. 2017), but it only showed similar ability to distinguish two very different cell classes: superficial- and
33 deep-layer excitatory neurons.

34 In this study, we compared 463 matched nuclei and whole cells from layer 5 of mouse primary visual cortex
35 (VISp) to investigate differences in single nucleus and single cell transcriptomes. We selected this brain
36 region because it contains a known variety of distinguishable yet highly similar cell types that would reveal
37 the cell-type detection limit of RNA-seq data obtained from single cells or nuclei (Tasic et al. 2016). We used
38 the same primary cell source and processed cells and nuclei with the same transcriptomic profiling method to
39 directly compare the resolution limit of cell type detection from well-matched sets of single cells and nuclei.
40 Furthermore, we compared the nuclear fraction of gene transcripts among cell types and identified functional
41 classes of transcripts that are enriched in the cytoplasm and nucleus.

42 Results

43 RNA-seq profiling of single nuclei and single cells

44 We isolated 487 NeuN-positive single nuclei from layer 5 of mouse VISp using fluorescence activated cell
45 sorting (FACS). Anti-NeuN staining was performed to enrich for neurons. In parallel, we isolated 12,866
46 tdT-positive single cells by FACS from all layers of mouse VISp and a variety of Cre-driver lines, as part of
47 a larger study on cortical cell type diversity (Tasic et al. 2017). For both single nuclei and cells, poly(A)-
48 transcripts were reverse transcribed and amplified with SMART-Seq v4, cDNA was tagged by Nextera
49 XT, and resulting libraries were sequenced to an average depth of 2.5 million reads (Figure 1A). RNA-seq
50 reads were mapped to the mouse genome using the STAR aligner (Dobin et al. 2013). Gene expression
51 was quantified as the sum of intronic and exonic reads per gene and was normalized as counts per million
52 (CPM) and log₂-transformed. For each nucleus and cell, the probabilities of gene detection dropouts were
53 estimated as a function of average expression level based on empirical noise models (Kharchenko, Silberstein,
54 and Scadden 2014).

55 463 out of 487 single nuclei (95%) passed quality control metrics, and each nucleus was matched to the
56 most similar nucleus and cell based on the maximum correlated expression of all genes, weighted for gene
57 dropouts. Nuclei had similarly high pairwise correlations to cells as to other nuclei suggesting that cells and
58 nuclei were well matched (Figure 1B). As expected, matched cells were derived almost exclusively from layer
59 5 and adjacent layers 4 and 6 (Figure S1B), and from Cre-driver lines that labeled cells in layer 5 (Figure 1C
60 and Figure S1A,C). The small minority of matched cells isolated from superficial layers were GABAergic
61 interneurons that have been detected in many layers (Tasic et al. 2017).

62 Comparison of nuclear and whole cell transcriptomes

63 scRNA-seq profiles nuclear and cytoplasmic transcripts, whereas snRNA-seq profiles nuclear transcripts.
64 Therefore, we expect that RNA-seq reads will differ between nuclei and cells. In nuclei, more than 50%
65 of reads that aligned to the mouse genome did not map to known spliced transcripts but to non-exonic
66 regions within gene boundaries. They were therefore annotated as intronic reads (Figure 2A). In contrast,
67 the majority of cells had less than 30% intronic reads with a minority of cells having closer to 50% intronic
68 reads, similar to nuclei. Median gene detection based on exonic reads was lower for nuclei (~5,000 genes) than
69 for cells (~9,500). Including both intronic and exonic reads increased gene detection for nuclei (~7,000) and
70 cells (~11,000), demonstrating that intronic reads provided additional information not captured by exons.
71 Whole brain control RNA displayed a read mapping distribution similar to cells, which is consistent with
72 dissociated single cells capturing the majority of transcripts in the whole cell.

73 Transcript dropouts likely result from both technical and biological variability, and both effects are more
74 pronounced in nuclei than in cells. When transcript dropouts were adjusted based on empirical noise models,
75 correlations between pairs of nuclei and pairs of cells increased, although cell-cell similarities remained sig-
76 nificantly higher (Figure 2B). A majority of expressed genes (21,279; 63%) showed similar detection (<10%
77 difference) in nuclei and cells, whereas 7,217 genes (21%) were detected in at least 25% more cells than
78 nuclei (Figure 2C and Table S1). 8,614 genes have significantly higher expression in cells than nuclei (>1.5
79 fold expression; FDR < 0.05) and many are involved in house-keeping functions such as mRNA processing
80 and translation (Figure 2D). Genetic markers of neuronal activity, such as immediate early genes *Fos*, *Egr1*,
81 and *Arc* also displayed up to 10-fold increased expression in cells, potentially a byproduct of tissue dissocia-
82 tion (Lacar et al. 2016). 159 genes have significantly higher expression in nuclei (Figure 2D and Table S2),
83 and they appear relevant to neuronal identity as they include connectivity and signaling genes (Figure S2A
84 and Table S4). Based on the sum of intronic and exonic reads, these 159 nucleus-enriched genes are on
85 average more than 10-fold longer than cell-enriched genes (Figure S2B), as recently reported for single nuclei
86 in mouse somatosensory cortex (Lake et al. 2017). When only exonic reads were used to quantify expression
87 in nuclei and cells, a different set of 146 genes were significantly enriched in nuclei (Table S3) and were only
88 slightly longer than cell-enriched genes. These genes were not associated with neuron-specific functions and
89 were significantly enriched for genes that participate in pre-mRNA splicing.

90 Intronic reads are required for high-resolution cell type identification from snRNA- 91 seq

92 Next, we applied an iterative clustering procedure (see Methods and Figure S3) to identify clusters of single
93 nuclei and cells that share gene expression profiles. To assess cluster robustness, we repeated clustering
94 100 times using 80% random subsets of nuclei and cells and calculated the proportion of clustering runs in
95 which each pair of samples clustered together. Co-clustering matrices were reordered using Ward's hierar-
96 chical clustering and represented as heatmaps with coherent clusters ordered as squares along the diagonal
97 (Figure 3A,B).

98 Clustering includes two steps – differentially expressed (DE) gene selection and distance measurement – that
99 are particularly sensitive to expression quantification. We repeated clustering using intronic and exonic reads
100 or only exonic reads for these steps, and ordered co-clustering matrices to match the results using all reads
101 for both steps. When using introns and exons, we found 11 distinct clusters of nuclei and cells, and clusters
102 had similar cohesion (average within cluster co-clustering) and separation (average co-clustering difference
103 with the closest cluster) (Figure 3C). Including intronic reads for either clustering step increased the number
104 of clusters detected for nuclei but not cells. Therefore, accounting for intronic reads in snRNA-seq was
105 critical to enable high-resolution cluster detection equivalent to that observed with scRNA-seq.

106 Equivalent cell types identified with nuclei and cells

107 We used hierarchical clustering of median gene expression values in each cluster to determine the relationships
108 between clusters. We find that cluster relationships represented as dendrograms are remarkably similar for
109 nuclei and cells (Figure 4A). We compared the 11 clusters identified with single nuclei and cells to reported
110 cell types in mouse VISp (Tasic et al. 2016). Each nucleus and cell cluster could be linked to a reported cell
111 type (Figure S4A) and to each other (Figure 4B) based on correlated expression of marker genes. Many genes
112 contributed to high expression correlations ($r > 0.85$) for all cluster pairs (Figure S4B). Conserved marker
113 gene expression confirmed that the same 11 cell types were identified with nuclei and cells (Figure 4C).
114 These cell types included nine excitatory neuron types from layers 4-6 and two inhibitory interneuron types.
115 Matched cluster proportions were mostly consistent, except two closely related layer 5a subtypes were under-
116 (L5a Batf3) or over-represented (L5a Hsd11b1) among cells (Figure S4C). This demonstrated that the initial
117 matching of cells to nuclei was relatively unbiased.

118 We hypothesized that most intronic reads were mapped to nuclear transcripts, so quantifying gene expression
119 in cells using only introns would approximate nuclear expression. This was supported by higher correlations
120 of average expression across all nuclei and cells using only intronic reads as compared to only exonic reads
121 (Figure S4D). Thus, a dendrogram based on the median expression (quantified using only intronic reads) of
122 nuclei and cell clusters paired all matching cell types, except for two closely related layer 5b subtypes (Fig-
123 ure 4D). Therefore, intronic reads can help facilitate comparisons between data sets derived from snRNA-seq
124 and scRNA-seq although small expression differences remain. A dendrogram based on exonic reads grouped
125 clusters first by sample type (nuclei and cells) and then by broad cell class (inhibitory and excitatory neu-
126 rons). Samples grouped by sample type likely due to differences in cytoplasmic transcripts that were profiled
127 in cells but not nuclei. A dendrogram based on intronic reads did not show this grouping because most
128 cytoplasmic transcripts are spliced so were quantified by exonic but not intronic reads.

129 While we detected the same cell types using nuclei and cells, we expected that gene expression captured with
130 cells included additional information from cytoplasmic transcripts. We compared the separation of matched
131 pairs of clusters based on co-clustering and found that all nuclei and cell clusters were similarly distinct,
132 except using single cells significantly increased the separation of two pairs of similar types: L4 Arf5 from
133 L5a Hsd11b1 and L5b Cdh13 from L5b Tph2 (Figure 4E). Next, we compared how well genes marked cell
134 types by calculating the degree of binary expression. Cell marker scores were, on average, 15% higher than
135 nucleus scores due to fewer expression dropouts in cells (Figure 4F), and this was consistent with mildly
136 improved cluster separation.

137 Nuclear content varies among cell types and for different transcripts

138 We estimated the nuclear proportion of mRNA for each cell type in two ways. Transcripts in the cytoplasm
139 are spliced so intronic reads should be restricted to the nucleus. First, we estimated the nuclear proportion
140 by calculating the ratio of the percentage of intronic reads in cells to the percentage of intronic reads in nuclei
141 (Figure 5A). Second, we estimated nuclear proportions by selecting three genes (*Malat1*, *Meg3*, and *Snhg11*)
142 with the highest expression in nuclei (Figure S4D) and calculating the ratio of the average expression in cells
143 versus nuclei (Figure 5B and Figure S5A). Both methods predicted that L4 Arf5 and L5a Hsd11b1 had a
144 significantly larger proportion of transcripts located in the nucleus compared to other cell types (Figure 5C).

145 Based on the comparison of scRNA-seq and snRNA-seq data, we estimate that L4 types have high nuclear
146 to cell volume (~50%), whereas L5 types have lower nuclear to cell volume. To evaluate this finding, we
147 measured nucleus and soma sizes of different cell types *in situ*. These types were labeled by different Cre-
148 transgenes and a Cre-reporter. *Nr5a1*-Cre and *Scnn1a-Tg3*-Cre mice almost exclusively label two cell types

149 (L4 *Arf5* and L5a *Hsd11b1*), whereas *Rbp4*-Cre mice label all layer 5 cell types including L5a *Hsd11b1*
150 (Figure S5B and Table S5). We measured the nuclear and cell sizes *in situ*, and calculated the nuclear
151 proportion of each cell as the ratio of nuclear to soma volume (Figure S5C). We found that the average
152 nuclear proportion was significantly lower for layer 5 cells compared to layer 4 cells, as predicted based on
153 RNA-seq data (Figure 5D).

154 In addition, nuclear proportion estimates based on *in situ* size measurements were systematically higher than
155 predicted for layer 5 but not layer 4 neurons. This could be the result of under-estimating the soma volume
156 based on cross-sectional area measurements of these large non-spherical (pyramidal) neurons. Alternatively,
157 layer 5 neuronal nuclei may have a lower density of nuclear transcripts or there may be cell type-specific biases
158 in our RNA-seq based estimates. We then performed an unbiased survey of nuclear proportions across the
159 full depth of cortex to test whether layer 4 or layer 5 neurons were exceptional compared to neurons in other
160 layers. We found that layer 5 neurons tend to be larger and have proportionally smaller nuclei (Figure S5D)
161 than other cortical neurons, and this feature is also found in rat primary visual cortex (Sigl-Glöckner and
162 Brecht 2017).

163 Next, we determined the nuclear versus cytoplasmic distribution of transcripts for individual genes. The
164 nuclear proportion of 11,932 transcripts was estimated by the ratio of nuclear to whole cell expression mul-
165 tiplied by the overall nuclear fraction of each cell type and averaged across cell types (Table S6). Different
166 functional classes of genes had strikingly different nuclear proportions (Figure 5E). Many non-coding trans-
167 cripts were localized in the nucleus, but some were abundantly expressed in the cytoplasm, such as the long
168 non-coding RNA (lncRNA) *Tunax* that is highly enriched in the brain, is conserved across vertebrates, and
169 has been associated with striatal pathology in Huntington's disease (Lin et al. 2014). Most protein-coding
170 transcripts were expressed in both the nucleus and cytoplasm with a small number restricted to the nucleus,
171 including the Parkinson's risk gene *Park2*. We found that pseudogenes were almost exclusively cytoplasmic
172 and were highly enriched for house-keeping functions.

173 We compared our estimates of nuclear enrichment in cortex to mouse liver and pancreas based on data
174 from (Halpern et al. 2015) and found moderately high correlation ($r = 0.61$) between 4,373 mostly house-
175 keeping genes that were expressed in all three tissues. Moreover, the shape of the distributions of nuclear
176 transcript proportions was highly similar between tissues with slightly higher proportions estimated in this
177 study. These results suggest that the mechanisms regulating the spatial localization of these transcripts – for
178 example, rates of nuclear export and cytoplasmic degradation (Halpern et al. 2015) – are conserved across
179 cell types.

180 Surprisingly, non-coding genes and pseudogenes are better markers of cell types, on average, than protein-
181 coding genes (Figure 5F). lncRNAs are known to have more specific expression among diverse human cell
182 lines (Djebali et al. 2012), and we show that this is also true for neuronal types in the mouse cortex. Many
183 pseudogene transcripts, most of which are enriched in the cytoplasm, were selectively depleted in the two
184 cell types, L4 *Arf5* and L5a *Hsd11b1*. This is consistent with our previous analysis that showed that neurons
185 of these types have relatively less cytoplasm. We also find that nucleus-enriched transcripts are slightly
186 better cell-type markers than cytoplasm-enriched transcripts, although this is highly variable across genes
187 (Figure 5G).

188 Finally, we compared our estimates of nuclear localization of transcripts for three genes – *Calb1*, *Grik1*,
189 and *Pvalb* – to relative counts of transcripts in nuclei and cytoplasm using multiplex RNA fluorescence *in*
190 *situ* hybridization (mFISH). We found that the relative nuclear proportions estimated by scRNA-seq and
191 mFISH were consistent although the absolute levels were quite variable (Figure 5H). Both methods confirmed
192 that *Pvalb* transcripts were mostly excluded from the nucleus, and this explained why 2 out of 35 nuclei

193 in the Pvalb-positive interneuron type (Pvalb Wt1) had no detectable *Pvalb* expression, whereas all cells of
194 this cell type had robust *Pvalb* expression.

195 Discussion

196 Unlike scRNA-seq, snRNA-seq enables transcriptomic profiling of tissues that are refractory to whole cell
197 dissociation and of archived frozen specimens. snRNA-seq is also less susceptible to perturbations in gene
198 expression that occur during cell isolation, such as increased expression of immediate early genes that can
199 obscure transcriptional signatures of neuronal activity (Lacar et al. 2016). However, these advantages come
200 at the cost of profiling less mRNA, and until this study, it was unclear if the nucleus contained sufficient
201 number and diversity of transcripts to distinguish highly related cell types.

202 To directly address this question, we profiled a well-matched set of 463 nuclei and 463 cells from layer 5 of
203 mouse primary visual cortex and identified 11 matching neuronal types: 2 interneuron types and 9 similar
204 excitatory neuron types. Including intronic reads in gene expression quantification was necessary to achieve
205 high-resolution cell type identification from single nuclei. Intronic reads substantially increased gene detection
206 to 7000 genes per nucleus. In addition, intronic reads were more frequently derived from long genes that
207 are known to have brain-specific expression (Gabel et al. 2015) and that help define neuronal connectivity
208 and signaling. Intronic reads may also reflect other cell-type specific features, such as retained introns or
209 alternative isoforms. For example, intron retention provides a mechanism for the nuclear storage and rapid
210 translation of long transcripts in response to neuronal activity (Mauger, Lemoine, and Scheiffele 2016).

211 We found that nuclei contain at least 20% of all cellular transcripts, and this percentage varies among cell
212 types. Two small pyramidal neuron types have large nuclei relative to cell size that contain more than half
213 of all transcripts. We detect 4000 more genes in single cells than single nuclei, but the majority of genes are
214 detected equally well in both. Cytoplasm-enriched transcripts are missed by profiling single nuclei but include
215 mostly house-keeping genes and pseudogenes, which are not related to neuronal identity. Nucleus-enriched
216 transcripts include protein-coding and non-coding genes that are more likely to be cell-type markers than
217 cytoplasmic transcripts. Overall, single cells do provide somewhat better detection of cell-type marker genes,
218 thereby resulting in slightly better cluster separation for two pairs of highly similar cell types. Therefore, as
219 more nuclei and cells are profiled, it is possible that finer discrimination of cell types may require single cell
220 profiling. However, the benefits of profiling single nuclei may outweigh potential loss in the finest cell type
221 resolution.

222 snRNA-seq is well suited for large-scale surveys of cellular diversity in various tissues and has the potential to
223 be less cell-type biased. For example, single cell profiling of adult human cortex isolated more interneurons
224 than excitatory neurons (Darmanis et al. 2015), whereas single nucleus profiling of the same tissue type
225 isolated 30% interneurons and 70% excitatory neurons (Lake et al. 2016), close to the proportions found *in*
226 *situ*. snRNA-seq also enables the use of stored frozen specimens to study cell types that will inform our
227 understanding of human diversity and disease. As large scale initiatives begin to characterize transcriptomic
228 cell types in the whole brain (Ecker et al. 2017) and whole organism (Regev et al. 2017), it is important to
229 understand the strengths and limitations of each mRNA profiling technique.

230 Materials and Methods

231 Tissue preparation

232 Tissue samples were obtained from adult (postnatal day (P) 53-59)) male and female transgenic mice carrying
233 a Cre transgene and a Cre-reporter transgene. Mice were anesthetized with 5% isoflurane and intracardially
234 perfused with either 25 or 50 ml of ice cold, oxygenated artificial cerebral spinal fluid (ACSF) at a flow
235 rate of 9 ml per minute until the liver appeared clear, or the full volume of perfusate had been flushed
236 through the vasculature. The ACSF solution consisted of 0.5mM CaCl₂, 25mM D-Glucose, 98mM HCl, 20mM
237 HEPES, 10mM MgSO₄, 1.25mM NaH₂PO₄, 3mM Myo-inositol, 12mM N-acetylcysteine, 96mM N-methyl-
238 D-glucamine, 2.5mM KCl, 25mM NaHCO₃, 5mM sodium L-Ascorbate, 3mM sodium pyruvate, 0.01mM
239 Taurine, and 2mM Thiourea. The brain was then rapidly dissected and mounted for coronal slice preparation
240 on the chuck of a Compresstome VF-300 vibrating microtome (Precisionary Instruments). Using a custom
241 designed photodocumentation configuration (Mako G125B PoE camera with custom integrated software),
242 a blockface image was acquired before each section was sliced at 250 μm intervals. The slice was then
243 hemisected along the midline, and both hemispheres were then transferred to chilled, oxygenated ACSF.

244 Each slice-hemisphere was transferred into a Sylgard-coated dissection dish containing 3 ml of chilled, oxy-
245 genated ACSF. Brightfield and fluorescent images between 4X and 20X were obtained of the intact tissue with
246 a Nikon Digital Sight DS-Fi1 or a Sentech STC-SC500POE camera mounted to a Nikon SMZ1500 dissecting
247 microscope. To guide anatomical targeting for dissection, boundaries were identified by trained anatomists,
248 comparing the blockface image and the slice image to a matched plane of the Allen Reference Atlas. In
249 general, three to five slices were sufficient to capture the targeted region of interest, allowing for expression
250 analysis along the anterior/posterior axis. The region of interest was then dissected and both brightfield and
251 fluorescent images of the dissections were acquired for secondary verification. The dissected regions were
252 transferred in ACSF to a microcentrifuge tube, and stored on ice. This process was repeated for all slices
253 containing the target region of interest, with each region of interest deposited into a new microcentrifuge
254 tube.

255 For whole cell dissociation, after all regions of interest were dissected, the ACSF was removed and 1 ml of
256 a 2 mg/ml pronase in ACSF solution was added. Tissue was digested at room temperature (approximately
257 22°C) for a duration that consisted of adding 15 minutes to the age of the mouse (in days; *i.e.*, P53 specimen
258 had a digestion time of 68 minutes). After digestion, the pronase solution was removed and replaced by
259 1 ml of ACSF supplemented with 1% Fetal Bovine Serum (FBS). The tissue was washed two more times
260 with the same solution and the sample was then triturated using fire-polished glass pipettes of decreasing
261 bore sizes (600, 300, and 150 μm). The cell suspension was incubated on ice in preparation for fluorescence-
262 activated cell sorting (FACS). FACS preparation involved adding 4'-6-diamidino-2-phenylindole (DAPI) at
263 a final concentration of 4 μg/ml to label dead (DAPI+) versus live (DAPI-) cells. The suspension was then
264 filtered through a fine-mesh cell strainer to remove cell aggregates. Cells were sorted by excluding DAPI
265 positive events and debris, and gating to include red fluorescent events (tdTomato-positive cells). Single
266 cells were collected into strip tubes containing 11.5μl of collection buffer (SMART-Seq v4 lysis buffer 0.83x,
267 Clontech #634894), RNase Inhibitor (0.17U/μl), and ERCCs (External RNA Controls Consortium, MIX1
268 at a final dilution of 1x10⁻⁸) (Baker et al. 2005; Risso et al. 2014). After sorting, strip tubes containing
269 single cells were centrifuged briefly and then stored at -80°C.

270 For nuclei isolation, dissected regions of interest were transferred to microcentrifuge tubes, snap frozen in a
271 slurry of dry ice and ethanol, and stored at -80°C until the time of use. To isolate nuclei, frozen tissues were
272 placed into a homogenization buffer that consisted of 10mM Tris pH 8.0, 250mM sucrose, 25mM KCl, 5mM
273 MgCl₂, 0.1% Triton-X 100, 0.5% RNasin Plus RNase inhibitor (Promega), 1X protease inhibitor (Promega),
274 and 0.1mM DTT. Tissues were placed into a 1ml dounce homogenizer (Wheaton) and homogenized using 10

275 strokes of the loose dounce pestle followed by 10 strokes of the tight pestle to liberate nuclei . Homogenate
276 was strained through a 30µm cell strainer (Miltenyi Biotech) and centrifuged at 900xg for 10 minutes to pellet
277 nuclei. Nuclei were then resuspended in staining buffer containing 1X PBS supplemented with 0.8% nuclease-
278 free BSA and 0.5% RNasin Plus RNase inhibitor. Mouse anti-NeuN antibody (EMD Millipore, MAB377,
279 Clone A60) was added to the nuclei at a final dilution of 1:1000 and nuclei suspensions were incubated at
280 4°C for 30 minutes. Nuclei suspensions were then centrifuged at 400xg for 5 minutes and resuspended in
281 clean staining buffer (1X PBS, 0.8% BSA, 0.5% RNasin Plus). Secondary antibody (goat anti-mouse IgG
282 (H+L), Alexa Fluor 594 conjugated, ThermoFisher Scientific) was applied to nuclei suspensions at a dilution
283 of 1:5000 for 30 minutes at 4°C. After incubation in secondary antibody, nuclei suspensions were centrifuged
284 at 400xg for 5 minutes and resuspended in clean staining buffer. Prior to FACS, DAPI was applied to nuclei
285 suspensions at a final concentration of 0.1µg/ml and nuclei suspensions were filtered through a 35µm nylon
286 mesh to remove aggregates. Single nuclei were captured by gating on DAPI-positive events, excluding debris
287 and doublets, and then gating on Alexa Fluor 594 (NeuN) signal. Strip tubes containing FACS isolated
288 single nuclei were then briefly centrifuged and frozen at -80°C.

289 RNA amplification and library preparation for RNA-seq

290 The SMART-Seq v4 Ultra Low Input RNA Kit for Sequencing (Clontech #634894) was used per the ma-
291 nufacturer's instructions for reverse transcription of single cell RNA and subsequent cDNA synthesis. Single
292 cells were stored in 8-strips at -80°C in 11.5 µl of collection buffer (SMART-Seq v4 lysis buffer at 0.83x,
293 RNase Inhibitor at 0.17 U/µl, and ERCC MIX1 at a final dilution of 1x10⁻⁸ dilution). Twelve to 24 8-well
294 strips were processed at a time (the equivalent of 1-2 96-well plates). At least 1 control strip was used per
295 amplification set, containing 2 wells without cells but including ERCCs, 2 wells without cells or ERCCs, and
296 either 4 wells of 10 pg of Mouse Whole Brain Total RNA (Zyagen, MR-201) or 2 wells of 10 pg of Mouse
297 Whole Brain Total RNA (Zyagen, MR-201) and 2 wells of 10 pg Control RNA provided in the Clontech
298 kit. Mouse whole cells were subjected to 18 PCR cycles after the reverse transcription step, whereas mouse
299 nuclei were subjected to 21 PCR cycles. AMPure XP Bead (Agencourt AMPure beads XP PCR, Beckman
300 Coulter A63881) purification was done using the Agilent Bravo NGS Option A instrument. A bead ratio of
301 1x was used (50 µl of AMPure XP beads to 50 µl cDNA PCR product with 1 µl of 10x lysis buffer added, as
302 per Clontech instructions), and purified cDNA was eluted in 17 µl elution buffer provided by Clontech. All
303 samples were quantitated using PicoGreen® on a Molecular Dynamics M2 SpectraMax instrument. A por-
304 tion of the samples, and all controls, were either run on the Agilent Bioanalyzer 2100 using High Sensitivity
305 DNA chips or the Advanced Analytics Fragment Analyzer (96) using the High Sensitivity NGS Fragment
306 Analysis Kit (1bp-6000bp) to qualify cDNA size distribution. An average of 7.3 ng of cDNA was synthesized
307 across all non-control samples. Purified cDNA was stored in 96-well plates at -20°C until library preparation.

308 Sequencing libraries were prepared using NexteraXT (Illumina, FC-131-1096) with NexteraXT Index Kit
309 V2 Set A (FC-131-2001). NexteraXT libraries were prepared at 0.5x volume, but otherwise followed the
310 manufacturer's instructions. An aliquot of each amplified cDNA sample was first normalized to 30 pg/µl
311 with Nuclease-Free Water (Ambion), then this normalized sample aliquot was used as input material into
312 the NexteraXT DNA Library Prep (for a total of 75pg input). AMPure XP bead purification was done using
313 the Agilent Bravo NGS Option A instrument. A bead ratio of 0.9x was used (22.5 ul of AMPure XP beads
314 to 25 ul library product, as per Illumina protocol), and all samples were eluted in 22 µl of Resuspension
315 Buffer (Illumina). All samples were run on either the Agilent Bioanalyzer 2100 using High Sensitivity DNA
316 chips or the Advanced Analytics Fragment Analyzer (96) using the High Sensitivity NGS Fragment Analysis
317 Kit (1bp-6000bp) to for sizing. All samples were quantitated using PicoGreen using a Molecular Dynamics
318 M2 SpectraMax instrument. Molarity was calculated for each sample using average size as reported by
319 Bioanalyzer or Fragment Analyzer and pg/µl concentration as determined by PicoGreen. Samples (5 µl
320 aliquot) were normalized to 2-10 nM with Nuclease-free Water (Ambion), then 2 µl from each sample within
321 one 96-index set was pooled to a total of 192 µl at 2-10 nM concentration. A portion of this library pool
322 was sent to an outside vendor for sequencing on an Illumina HS2500. All of the library pools were run using

323 Illumina High Output V4 chemistry. Covance Genomics Laboratory, a Seattle-based subsidiary of LabCorp
324 Group of Holdings, performed the RNA-Sequencing services. An average of 229 M reads were obtained per
325 pool, with an average of 2.0-3.1 M reads/cell across the entire data set.

326 RNA-Seq data processing

327 Raw read (fastq) files were aligned to the GRCm38 mouse genome sequence (Genome Reference Consortium,
328 2011) with the RefSeq transcriptome version GRCm38.p3 (current as of 1/15/2016) and updated by remov-
329 ing duplicate Entrez gene entries from the gtf reference file for STAR processing. For alignment, Illumina
330 sequencing adapters were clipped from the reads using the fastqMCF program (Aronesty 2011). After clip-
331 ping, the paired-end reads were mapped using Spliced Transcripts Alignment to a Reference (STAR) (Dobin
332 et al. 2013) using default settings. STAR uses and builds its own suffix array index which considerably
333 accelerates the alignment step while improving sensitivity and specificity, due to its identification of alterna-
334 tive splice junctions. Reads that did not map to the genome were then aligned to synthetic constructs (i.e.
335 ERCC) sequences and the *E.coli* genome (version ASM584v2). Quantification was performed using summer-
336 izeOverlaps from the R package GenomicAlignments (Lawrence et al. 2013). Read alignments to the genome
337 (exonic, intronic, and intergenic counts) were visualized as beeswarm plots using the R package *beeswarm*.

338 Expression levels were calculated as counts per million (CPM) of exonic plus intronic reads, and $\log_2(\text{CPM}$
339 $+ 1)$ transformed values were used for a subset of analyses as described below. Gene detection was calculated
340 as the number of genes expressed in each sample with $\text{CPM} > 0$. CPM values reflected absolute transcript
341 number and gene length, i.e. short and abundant transcripts may have the same apparent expression level
342 as long but rarer transcripts. Intron retention varied across genes so no reliable estimates of effective gene
343 lengths were available for expression normalization. Instead, absolute expression levels were estimated as
344 fragments per kilobase per million (FPKM) using only exonic reads so that annotated transcript lengths
345 could be used.

346 Selection of single nuclei and matched cells

347 463 of 487 (95%) of single nuclei isolated from layer 5 of mouse VISp passed quality control criteria: $>500,000$
348 genome-mapped reads, $>75\%$ reads aligned, and $>50\%$ unique reads. 12,866 single cells isolated from layers
349 1-6 of mouse VISp passed quality control criteria: $>200,000$ transcriptome mapped reads and >1000 genes
350 detected ($\text{CPM} > 0$).

351 Gene expression was more likely to drop out in samples with lower quality cDNA libraries and for low ex-
352 pressing genes. To estimate gene dropouts due to stochastic transcription or technical artifacts (Kharchenko,
353 Silberstein, and Scadden 2014), expression noise models were fit separately to single nuclei and cells using
354 the “knn.error.models” function of the R package *scde* (version 2.2.0) with default settings and eight nearest
355 neighbors. Noise models were used to calculate a dropout weight matrix that represented the likelihood of
356 expression dropouts based on average gene expression levels of similar nuclei or cells using mode-relative
357 weighting (“SCDE by Kharchenko Lab at Harvard DBMI”, n.d.). The probability of dropout for each
358 sample (s) and gene (g) was estimated based on two expression measurements: average expected expression
359 level of similar samples, $p(x_{\bar{g}})$, and observed expression levels, $p(x_{sg})$, using the “scde.failure.probability”
360 and “scde.posterior” functions. The dropout weighting was calculated as a combination of these probabili-
361 ties: $W_{sg} = 1 - \sqrt{p(x_{sg}) \cdot \sqrt{p(x_{sg}) \cdot p(x_{\bar{g}})}}$.

362 Dropout weighted Pearson correlations were calculated between all pairs of nuclei and cells using 42,003
363 genes expressed in at least one nucleus and one cell. The cell with the highest correlation to any nucleus

364 was selected as the best match, and this cell and nucleus were removed from further analysis. This process
365 was repeated until 463 best matching cells were selected, and the expression correlations were compared to
366 correlations of the best matching pairs of nuclei (Figure 1B). The Cre-lines and dissected cortical layers of
367 origin of the best matching cells were summarized as bar plots (Figure S1). Unweighted Pearson correlations
368 were also calculated between all pairs of nuclei and cells to test the effect of accounting for dropouts on
369 sample similarities (Figure 2B).

370 Differential expression analysis

371 Gene detection was estimated as the proportion of cells and nuclei expressing each gene ($CPM > 0$). In order
372 to estimate the expected variability of gene detection as a result of population sampling, cells were randomly
373 split into two sets of 231 and 232 cells and genes were grouped into 50 bins based on detection in the first
374 set of cells. For each bin of genes, the 97.5 percentile of detection was calculated for the second set of cells.
375 A 95% confidence interval of gene detection was constructed by reflecting these binned quantiles across
376 the line of unity. Data were summarized with a hexagonal binned scatter plot and a log-transformed color
377 scale using the R package *ggplot2* (Wickham 2009).

378 Differential expression between nuclei and cells was calculated with the R package *limma* (Ritchie et al.
379 2015) using default settings and $\log_2(CPM + 1)$ expression defined based on two sets of reads: introns plus
380 exons and only exons. Significantly differentially expressed were defined as having >1.5 -fold change and
381 a Benjamini-Hochberg corrected P-value < 0.05 . Gene expression distributions of nuclei or cells within a
382 cluster were visualized using violin plots, density plots rotated 90 degrees and reflected on the Y-axis.

383 Differences in alignment statistics and gene counts were calculated between cells, nuclei, and total RNA
384 controls (or just cells and nuclei) with analysis of variance using the “aov” function in R (Chambers, Freeny,
385 and Heiberger 1992). P-values for all comparisons were $P < 10^{-13}$.

386 Two sets of nucleus- and cell-enriched genes (introns plus exons and exons only) were tested for gene ontology
387 (GO) enrichment using the ToppGene Suite (Chen et al. 2009). Significantly enriched (Benjamini-Hochberg
388 false discovery rate < 0.05) GO terms were summarized as tree maps with box sizes proportional to $-\log_{10}(P$ -
389 values) using REVIGO (Supek et al. 2011) (Figure S2).

390 Clustering

391 Nuclei and cells were grouped into transcriptomic cell types using an iterative clustering procedure based
392 on community detection in a nearest neighbor graph as described in (Levine et al. 2015). Clustering was
393 performed using gene expression quantified with exonic reads only or intronic plus exonic reads for two key
394 clustering steps: selecting significantly variable genes and calculating pairwise similarities between nuclei.
395 Four combinations of expression quantification for nuclei and cells resulted in eight independent clustering
396 runs.

397 For each gene, $\log_2(CPM + 1)$ expression was centered and scaled across samples. Noise models were used to
398 select significantly variable genes (adjusted variance > 1.25). Dimensionality reduction was performed with
399 principal components analysis (PCA) on variable genes, and the covariance matrix was adjusted to account
400 for gene dropouts using the product of dropout weights across genes for each pair of samples. A maximum
401 of 20 principal components (PCs) were retained for which more variance was explained than the broken stick
402 null distribution, a conservative method of PC retention (Jackson 1993).

403 Nearest-neighbor distances between all samples were calculated using the “nn2” function of the R pack-
404 age *RANN*, and Jaccard similarity coefficients between nearest-neighbor sets were computed. Jaccard coeffi-
405 cients measured the proportion of nearest neighbors shared by each sample and were used as edge weights in
406 constructing an undirected graph of samples. Louvain community detection was used to cluster this graph
407 with 15 nearest neighbors. Considering more than 15 neighbors reduced the power to detect small clusters
408 due to the resolution limit of community detection (Fortunato and Barthelemy 2007). Considering fewer
409 than 15 neighbors increased over-splitting, as expected based on simulations by (Reichardt and Bornholdt
410 2006). Fewer nearest neighbors were used only when there were 15 or fewer samples total.

411 Clustering significance was tested by comparing the observed modularity to the expected modularity of an
412 Erdős-Rényi random graph with a matching number of nodes and average connection probability. Expected
413 modularity was calculated as the maximum estimated by two reported equations (Guimerà, Sales-Pardo, and
414 Amaral 2004; Reichardt and Bornholdt 2006). Samples were split into clusters only if the observed modularity
415 was greater than the expected modularity, and only clusters with distinct marker genes were retained. Marker
416 genes were defined for all cluster pairs using two criteria: 1) significant differential expression (Benjamini-
417 Hochberg false discovery rate < 0.05) using the R package *limma* and 2) either binary expression (CPM > 1
418 in >50% samples in one cluster and <10% in the second cluster) or >100-fold difference in expression. Pairs
419 of clusters were merged if either cluster lacked at least one marker gene.

420 Clustering was applied iteratively to each sub-cluster until the occurrence of one of four stop criteria: 1)
421 fewer than six samples (due to a minimum cluster size of three); 2) no significantly variable genes; 3) no
422 significantly variable PCs; 4) no significant clusters.

423 To assess the robustness of clusters, the iterative clustering procedure described above was repeated 100 times
424 for random sets of 80% of samples. A co-clustering matrix was generated that represented the proportion of
425 clustering iterations that each pair of samples were assigned to the same cluster. Average-linkage hierarchical
426 clustering was applied to this matrix followed by dynamic branch cutting using “cutreeHybrid” in the R
427 package *WGCNA* (Langfelder, Zhang, and Horvath 2007) with cut height ranging from 0.01 to 0.99 in steps
428 of 0.01. A cut height was selected that resulted in the median number of clusters detected across all 100
429 iterations. Cluster cohesion (average within cluster co-clustering) and separation (difference between within
430 cluster co-clustering and maximum between cluster co-clustering) was calculated for all clusters. Marker genes
431 were defined for all cluster pairs as described above, and clusters were merged if they had a co-clustering
432 separation <0.25 or either cluster lacked at least one marker gene.

433 Scoring marker genes based on cluster specificity

434 Many genes were expressed in the majority of nuclei or cells in a subset of clusters. A marker score (beta)
435 was defined for all genes to measure how binary expression was among clusters, independent of the number
436 of clusters labeled. First, the proportion (x_i) of samples in each cluster that expressed a gene above back-
437 ground level (CPM > 1) was calculated. Then, scores were defined as the squared differences in proportions
438 normalized by the sum of absolute differences plus a small constant (ϵ) to avoid division by zero. Scores
439 ranged from 0 to 1, and a perfectly binary marker had a score equal to 1.

$$\beta = \frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2}{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j| + \epsilon}.$$

440 Cluster dendrograms

441 Clusters were arranged by transcriptomic similarity based on hierarchical clustering. First, the average
442 expression level of the top 1200 marker genes (i.e. highest beta scores) was calculated for each cluster.
443 A correlation-based distance matrix ($D_{xy} = \frac{1-\rho(x,y)}{2}$) was calculated, and complete-linkage hierarchical
444 clustering was performed using the “hclust” R function with default parameters. The resulting dendrogram
445 branches were reordered to show inhibitory clusters followed by excitatory clusters, with larger clusters first,
446 while retaining the tree structure. Note that this measure of cluster similarity is complementary to the
447 co-clustering separation described above. For example, two clusters with similar gene expression patterns
448 but a few binary marker genes may be close on the tree but highly distinct based on co-clustering.

449 Matching clusters based on marker gene expression

450 Nuclei and cell clusters were independently compared to published mouse VISp cell types (Tasic et al.
451 2016). The proportion of nuclei or cells expressing each gene with CPM > 1 was calculated for all clusters.
452 Approximately 400 genes were markers in both data sets (beta score > 0.3) and were expressed in the
453 majority of samples of between one and five clusters. Markers expressed in more than five clusters were
454 excluded to increase the specificity of cluster matching. Weighted correlations were calculated between all
455 pairs of clusters across these genes and weighted by beta scores to increase the influence of more informative
456 genes. Heatmaps were generated to visualize all cluster correlations. All nuclei and cell clusters had reciprocal
457 best matching clusters from Tasic et al. and were labeled based on these reported cluster names.

458 Next, nuclei and cell clusters were directly compared using the above analysis. All 11 clusters had reciprocal
459 best matches that were consistent with cluster labels assigned based on similarity to published types. The
460 most highly conserved marker genes of matching clusters were identified by selecting genes expressed in a
461 single cluster (>50% of samples with CPM > 1) and with the highest minimum beta score between nuclei
462 and cell clusters. Two additional marker genes were identified that discriminated two closely related clusters.
463 Violin plots of marker gene expression were constructed with each gene on an independent, linear scale.

464 Nuclei and cell clusters were also compared by calculating average cluster expression based only on intronic
465 or exonic reads and calculating a correlation-based distance using the top 1200 marker genes as described
466 above. Hierarchical clustering was applied to all clusters quantified using the two sets of reads. In addition,
467 the average $\log_2(\text{CPM} + 1)$ expression across all nuclei and cells was calculated using intronic or exonic
468 reads.

469 Cluster separation was calculated for individual nuclei and cells as the average within cluster co-clustering
470 of each sample minus the maximum average between cluster co-clustering. Separations for matched pairs of
471 clusters were visualized with box plots and compared using a Student’s *t*-test, and significance was tested
472 after Bonferroni correction for multiple testing. Finally, a linear model was fit to beta marker scores for
473 genes that were expressed in at least one but not all cell and nuclear clusters, and the intercept was set to
474 zero.

475 Estimating proportions of nuclear transcripts

476 The nuclear proportion of transcripts was estimated in two ways. First, all intronic reads were assumed to
477 be from transcripts localized to the nucleus so that the proportion of intronic reads measured in cells should
478 decrease linearly with the nuclear proportion of the cell as nuclear reads are diluted with cytoplasmic reads.
479 For each cell type, the nuclear proportion was estimated as the proportion of intronic reads in cells divided

480 by the proportion of intronic reads in matched nuclei. Second, the nuclear proportion was estimated as the
481 average ratio of cell to nuclear expression (CPM) using only exonic reads of three highly expressed nuclear
482 genes (*Snhg11*, *Malat1*, and *Meg3*). The standard deviation of nuclear proportion estimates were calculated
483 based on standard error propagation of variation in intronic read proportions and expression levels. Nuclear
484 proportion estimates were compared with linear regression, and the estimate based on relative expression
485 levels was used for further analysis.

486 The nuclear proportion of transcripts for all genes was estimated for each cell type as the ratio of average
487 expression (CPM) in nuclei versus matched cells multiplied by the nuclear proportion of all transcripts.
488 Estimated proportions greater than 1 were set equal to 1 for each cell type, and a weighted average proportion
489 was calculated for each gene with weights equal to the average $\log_2(\text{CPM} + 1)$ expression in each cell type.
490 11,932 genes were expressed in at least one nuclear or cell cluster (>50% samples expressed with CPM >
491 1) and were annotated as one of three gene types – protein-coding, protein non-coding, or pseudogene –
492 using gene metadata from NCBI ([ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/GENE_INFO/Mammalia/Mus_](ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/GENE_INFO/Mammalia/Mus_musculus.gene_info.gz)
493 [musculus.gene_info.gz](ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/GENE_INFO/Mammalia/Mus_musculus.gene_info.gz); downloaded 10/12/2017). For each type, histograms of gene counts with different
494 nuclear proportions were generated. Next, beta marker score distributions were visualized as violin plots,
495 and differences across gene types were compared with a Kruskal-Wallis rank sum test followed by Wilcoxon
496 signed rank unpaired tests. Finally, genes were grouped into 10 bins of estimated nuclear proportions, from
497 high cytoplasmic enrichment to high nuclear enrichment, and beta marker score distributions were visualized
498 as box plots. A linear regression was fit to marker scores versus nuclear proportion.

499 Nuclear transcript proportions were compared to nuclear proportions estimated for mouse liver and pan-
500 creatic beta cells based on data from (Halpern et al. 2015). Ratios of normalized nuclear and cytoplasmic
501 transcript counts were calculated in four tissue replicates. Average ratios were calculated for genes with at
502 least one count in either fraction in at least one tissue. Nuclear proportion estimates for all genes with data
503 from both data sets ($n = 4373$) were compared with Pearson correlation, a linear model with intercept set
504 equal to zero, and histograms with a bin width of 0.02.

505

506 **Colorimetric *in situ* hybridization**

507 *In situ* hybridization data for mouse cortex was from the Allen Mouse Brain Atlas (Lein et al. 2007). All data
508 is publicly accessible through www.brain-map.org. Data was generated using a semiautomated technology
509 platform as described in (Lein et al. 2007). Mouse ISH data shown is from primary visual cortex (VISp) in
510 the Paxinos Atlas (Paxinos and others 2013).

511 **Multiplex fluorescence RNA *in situ* hybridization and quantification of nuclear** 512 **versus cytoplasmic transcripts**

513 The RNAscope multiplex fluorescent kit was used according to the manufacturer's instructions for fresh
514 frozen tissue sections (Advanced Cell Diagnostics), with the exception that 16 μm tissue sections were fixed
515 with 4% PFA at 4°C for 60 minutes and the protease treatment step was shortened to 15 minutes at room
516 temperature. Probes used to identify nuclear and cytoplasmic enriched transcripts were designed antisense
517 to the following mouse genes: *Calb1*, *Grik1*, and *Pvalb*. Following hybridization and amplification, stained
518 sections were imaged using a 60X oil immersion lens on a Nikon TiE epifluorescence microscope.

519 To determine if spots fell within the nucleus or cytoplasm, a boundary was drawn around the nucleus to
520 delineate its border using measurement tools within Nikon Elements software. To delineate the cytoplasmic

521 boundary of each cell, a circle with a diameter of 15 μ m was drawn and centered over the cell (Fig. 5). RNA
522 spots in each channel were quantified manually using counting tools available in the Nikon Elements software.
523 Spots that fell fully within the interior boundary of the nucleus were classified as nuclear transcripts. Spots
524 that fell outside of the nucleus but within the circle that defined the cytoplasmic boundary were classified
525 as cytoplasmic transcripts. Additionally, if spots intersected the exterior boundary of the nucleus they were
526 classified as cytoplasmic transcripts. To prevent double counting of spots and ambiguities in assigning spots
527 to particular cells, labeled cells whose boundaries intersected at any point along the circumference of the
528 circle delineating their cytoplasmic boundary were excluded from the analysis. A linear regression was fit to
529 nuclear versus soma probe counts, and the slope was used to estimate the nuclear proportion.

530 *In situ* quantification of nucleus and soma size

531 Coronal brain slices from *Nr5a1-Cre;Ai14*, *Scnn1a-Tg3-Cre;Ai14*, and *Rbp4-Cre-KL100;Ai14* mice were stain-
532 ed with anti-dsRed (Clontech #632496) to enhance tdTomato signal in red channel and DAPI to label nuclei.
533 Maximum intensity projections from six confocal stacks of 1- μ m intervals were processed for analysis. Initial
534 segmentation was performed by CellProfiler (Lamprecht, Sabatini, and Carpenter 2007) to identify nuclei
535 from the DAPI signal and soma from the tdTomato signal. Segmentation results were manually verified and
536 any mis-segmented nuclei or somata were removed or re-segmented if appropriate. Area measurement of
537 segmented nuclei and somata was performed in CellProfiler in Layer 4 from *Nr5a1-Cre;Ai14* and *Scnn1a-*
538 *Tg3-Cre;Ai14* mice, and in Layer 5 from *Rbp4-Cre-KL100;Ai14* mice. A linear regression was fit to nuclear
539 versus soma area to highlight the differences between Cre-lines.

540 For measurements of nucleus and soma size agnostic to Cre driver, we used 16 μ m-tissue sections from P56
541 mouse brain. To label nuclei, DAPI was applied to the tissue sections at a final concentration of 1mg/ml.
542 To label cell somata, tissue sections were stained with Neurotrace 500/525 fluorescent Nissl stain (Ther-
543 moFisher Scientific) at a dilution of 1:100 in 1X PBS for 5 minutes, followed by brief washing in 1X PBS.
544 Sections were coverslipped with Fluoromount-G (Southern Biotech) and visualized on a Nikon TiE epiflu-
545 orescence microscope using a 40x oil objective. Soma and nuclei area measurements were taken by tracing
546 the boundaries of the Nissl-stained soma or DAPI-stained nucleus, respectively, using cell measurement tools
547 available in the Nikon TiE microscope software. All cells with a complete nucleus clearly present within the
548 section were measured, except that we excluded glial cells which had very small nuclei and scant cytoplasm.
549 Measurements were taken within a 40x field of view across an entire cortical column encompassing layers
550 1-6, and the laminar position of each cell (measured as depth from the pial surface) was tracked along with
551 the nucleus and soma area measurements for each cell.

552 For each cell in the experiments above, the nuclear proportion was estimated as the ratio of nucleus and soma
553 area raised to the $3/2$ power. This transformation was required to convert area to volume measurements
554 and assumed that the 3-dimensional geometries of soma and nuclei were reflected by their cross-sectional
555 profiles. This is true for approximately symmetrical shapes such as most nuclei and some somata, but will
556 lead to under- or over-estimates of nuclear proportions for asymmetrical cells. Therefore, the estimated
557 nuclear proportion of any individual cell may be inaccurate, but the average nuclear proportion for many
558 cells should be relatively unbiased.

559 Code availability

560 Data and code to reproduce all figures are publicly available from GitHub at [https://github.com/AllenInstitute/](https://github.com/AllenInstitute/NucCellTypes)
561 [NucCellTypes](https://github.com/AllenInstitute/NucCellTypes).

562 **Competing interests**

563 The authors declare no competing interests.

564 **Acknowledgements**

565 The authors thank the Allen Institute for Brain Science founders, P. G. Allen and J. Allen, for their vision,
566 encouragement, and support.

567 **Figures**

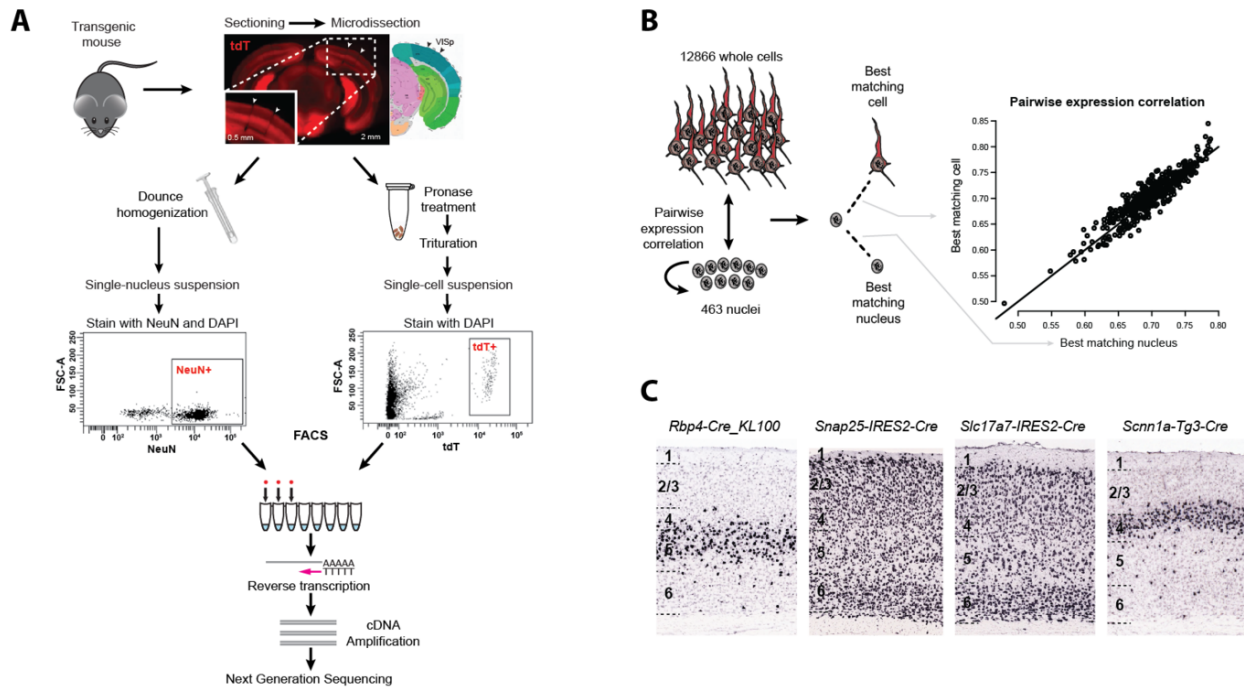


Figure 1: Identification of an expression-matched set of single nuclei and whole cells from mouse primary visual cortex (VISp). **(A)** Whole brains were dissected from transgenic mice, coronal slices were sectioned, and individual layers of VISp were microdissected. Nuclei were dissociated from layer 5, stained with DAPI and against the neuronal marker NeuN. Single NeuN-positive nuclei were isolated by fluorescence-activated cell sorting (FACS). In parallel, whole cells were dissociated from all layers, and single td-Tomato reporter-positive cells were isolated. Single nucleus and cell mRNA were reverse transcribed, amplified, and sequenced to measure transcriptome-wide expression levels. **(B)** Left: 463 nuclei from layer 5 and 12,866 whole cells from all layers passed quality control metrics, and the expression correlation was calculated between each nucleus and all other nuclei and cells. Expression similarity can vary based on sample quality, so nuclei were compared to each other to provide a baseline expected similarity. For each nucleus, the best matching nucleus and cell were selected based on maximal correlation. Right: Cells and nuclei displayed comparable expression similarities to all nuclei, with 95% of correlations between 0.63 and 0.78. This suggested that nuclei and cells were well matched. **(C)** Chromogenic RNA *In situ* hybridization (ISH) images of all VISp layers from four mouse Cre-lines from which the best matching cells were most commonly derived. As expected, all Cre-lines label cells in layer 5 and adjacent layers.

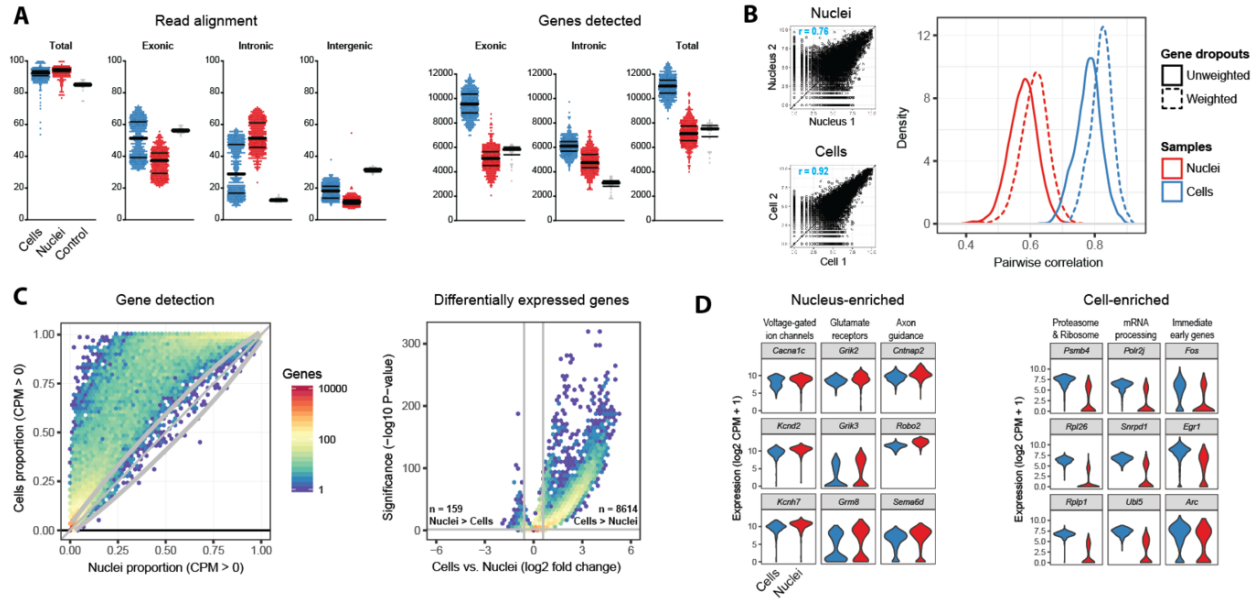


Figure 2: Comparison of nuclear and whole cell transcriptomes. **(A)** Left: Percentage of RNA-seq reads mapping to genomic regions for cells, nuclei, and whole brain control RNA. Bars indicate median and 25th and 75th quantiles. Note that among cells exonic and intronic read alignment is bimodal. Right: Gene detection (counts per million, CPM > 0) based on reads mapping to exons, introns, or both introns and exons. **(B)** Left: The most similar pair of cells have more highly correlated gene expression ($r = 0.92$) than the most similar pair of nuclei ($r = 0.76$), due to fewer gene dropouts. Right: Cells have consistently more similar expression to each other than nuclei, even after correcting for gene dropouts based on an expression noise model. **(C)** Left: Binned scatter plot showing all genes are detected (CPM > 0) with equal or greater reliability in cells than nuclei. Grey lines show the variation in detection that is expected by chance (95% confidence interval). Right: Binned scatter plot showing 0.4% of genes are significantly more highly expressed (fold change > 1.5, adjusted P-value < 0.05) in nuclei, and 20.5% of genes are more highly expressed in cells. The log-transformed color scale indicates the number of genes in each bin. **(D)** Nuclear enriched genes are highly enriched for genes involved in neuronal connectivity, synaptic transmission, and intrinsic firing properties. Cell enriched genes are predominantly related to mRNA processing and protein translation and degradation. In addition, immediate early gene expression is increased up to 10-fold in cells, despite comparable isolation protocols for cells and nuclei.

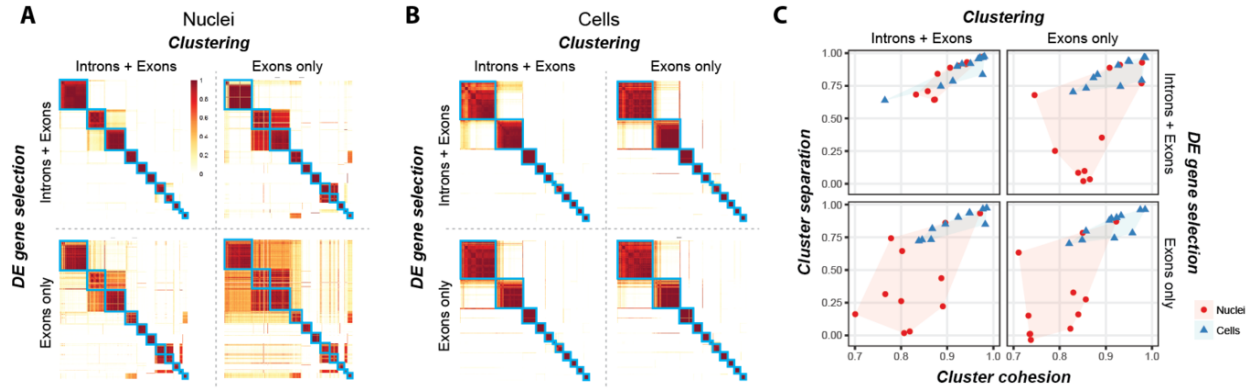


Figure 3: Single nuclei provide comparable clustering resolution to cells with inclusion of intronic reads. **(A)** Co-clustering heatmaps show the proportion of 100 clustering iterations that each pair of nuclei were assigned to the same cluster. Clustering was performed using gene expression quantified with exonic reads or intronic plus exonic reads for two key clustering steps: selecting significantly differentially expressed (DE) genes and calculating pairwise similarities between nuclei. Co-clustering heatmaps were generated for each combination of gene expression values, and blue boxes highlight 11 clusters of nuclei that consistently co-clustered using introns and exons (upper left heatmap) and were overlaid on the remaining heatmaps. The row and column order of nuclei is the same for all heatmaps. **(B)** Co-clustering heatmaps were generated for cells as described for nuclei in **(A)**, and blue boxes highlight 11 clusters of cells. **(C)** Cluster cohesion (average within cluster co-clustering) and separation (difference between within cluster co-clustering and maximum between cluster co-clustering) are plotted for nuclei and cells and all combinations of reads. Including introns in gene expression quantification dramatically increases cohesion and separation of nuclei but not cell clusters.

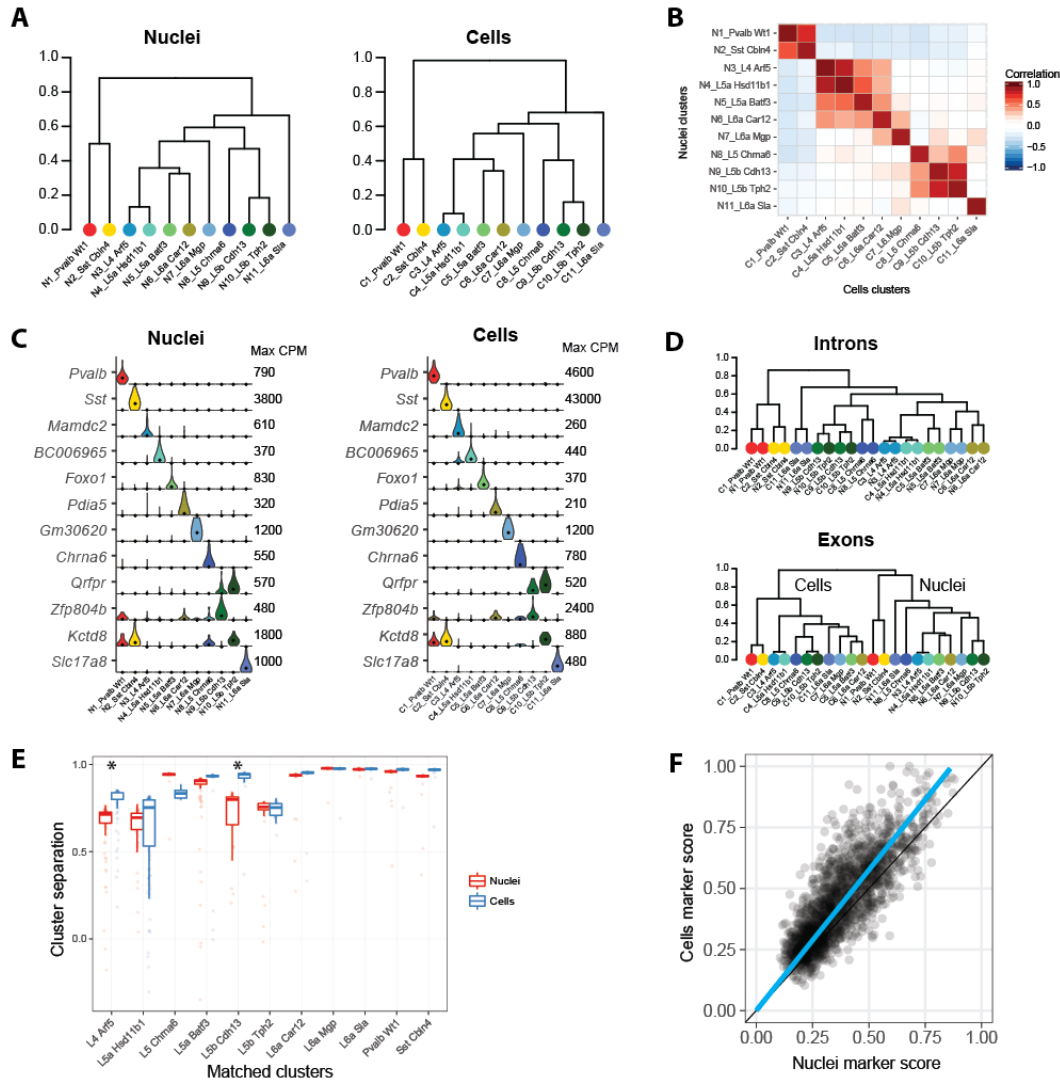


Figure 4: Equivalent neuronal cell types identified with nuclei and cells. **(A)** Cluster dendrograms for nuclei and cells based on hierarchical clustering of average expression of the top 1200 cluster marker genes. 11 clusters are labeled based on dendrogram leaf order and the closest matching mouse VISp cell type described in (Tasic et al. 2016) based on correlated marker gene expression (see Figure S4). **(B)** Pairwise correlations between nuclear and cell clusters using average cluster expression of the top 490 shared marker genes. **(C)** Violin plots of cell type specific marker genes expressed in matching nuclear and cell clusters. Plots are on a linear scale, max CPM indicates the maximum expression of each gene, and black dots indicate median expression. **(D)** Hierarchical clustering of nuclear and cell clusters using the top 1200 marker genes with expression quantified by intronic or exonic reads. Intronic reads group nine matching nuclear and cell clusters together at the leaves, while two closely related deep layer 5 excitatory neuron types group by sample type. In contrast, exonic reads completely segregate clusters by sample type. **(E)** Box plots of cluster separations for all samples in matched nuclear and cell clusters. Clusters are equally well separated for all but two cell types, L4 Arf5 and L5b Cdh13, that are moderately but significantly (Wilcoxon signed rank unpaired tests; Bonferroni corrected P-value < 0.05) more distinct with cells than nuclei. **(F)** Cell type marker genes are consistently detected in nuclei and cells, although marker scores (see Methods) were on average 15% higher for cells.

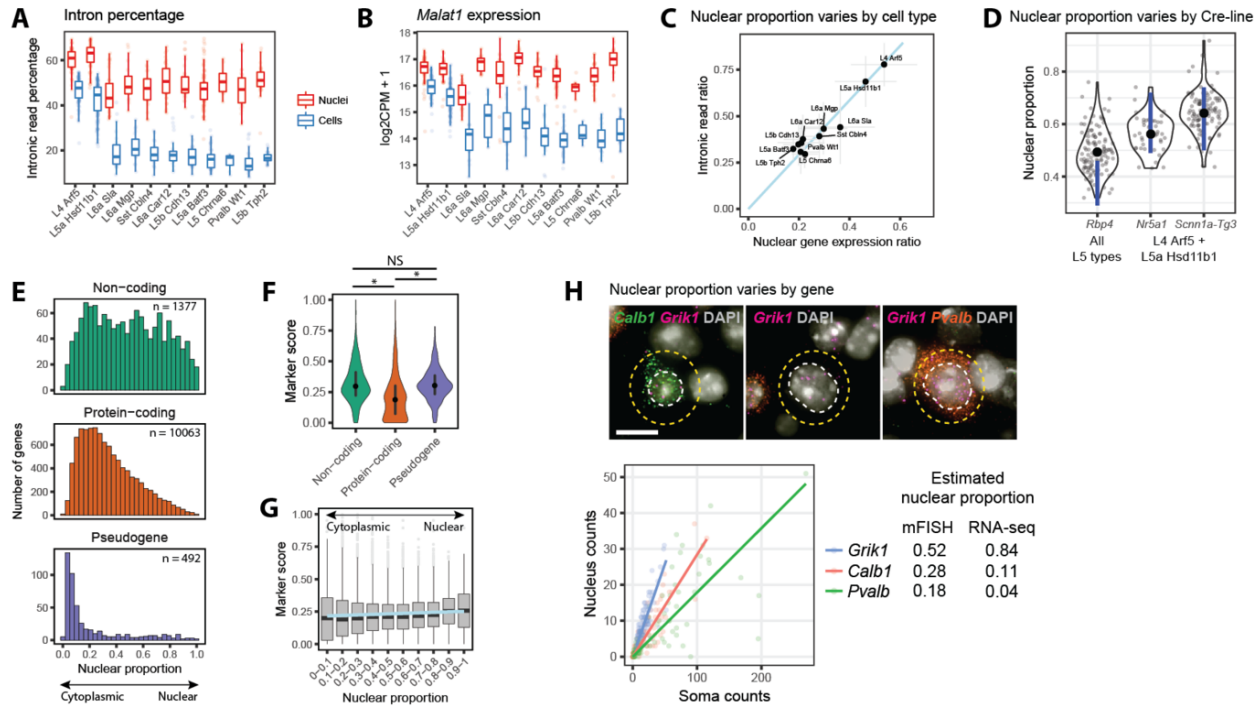


Figure 5: Nuclear transcript content varies among cell types and genes. **(A)** Box plots showing median (bars), 25th and 75th quantiles (boxes), and range (whiskers) of percentages of reads mapping to introns for matched nuclei and cell clusters. **(B)** Box plots of log₂-transformed expression of the nuclear non-coding RNA, *Malat1*, in matched nuclei and cell clusters. **(C)** The nuclear fraction of transcripts in cell types was estimated with two methods: the ratio of intronic read percentages in cells compared to nuclei; and the average ratio of expression in cells compared to nuclei of three highly expressed genes (*Snhg11*, *Meg3*, and *Malat1*) that are localized to the nucleus. The relative ranking of nuclear fractions was consistent (Spearman rank correlation = 0.84), although estimates based on the intronic read ratio were consistently 50% higher. **(D)** Estimated nuclear proportion (ratio of nucleus and soma volume) of neurons labeled by three mouse Cre-lines in Layers 4 and 5 (see Supplementary Figure S5D). Single neuron measurements (grey points) were summarized as violin plots, and average nuclear proportions (black points) were compared to the range of estimated proportions (blue lines) based on intronic read ratios and nuclear gene expression. **(E)** Histograms of nuclear fraction estimates for 11,932 genes expressed (CPM > 1) in at least one nuclear or cell cluster and grouped by type of gene. **(F)** Violin plots of marker score distributions with median and inter-quartile intervals. Non-coding genes and pseudogenes are on average better markers of cell types than protein-coding genes. Kruskal–Wallis rank sum test, post hoc Wilcoxon signed rank unpaired tests: *P < 1 × 10⁻⁵⁰ (Bonferroni-corrected), NS, not significant. **(G)** Box plots of cell type marker scores for genes grouped by estimated nuclear enrichment. Nucleus-enriched genes have significantly higher marker scores (linear regression; P = 2.3 × 10⁻⁸). **(H)** Validation of the estimated nuclear proportion of transcripts for *Calb1*, *Grik1*, and *Pvalb* using multiplex fluorescent *in situ* hybridization (mFISH). Top: For each gene, transcripts were labeled with fluorescent probes and counted in the nucleus (white) and soma (yellow). Bottom: Probe counts in the nucleus and soma across all cells with linear regression fits to estimate nuclear transcript proportions for each gene. Estimated proportions based on mFISH and RNA-seq data are summarized on the right.

568 Supplemental Figures

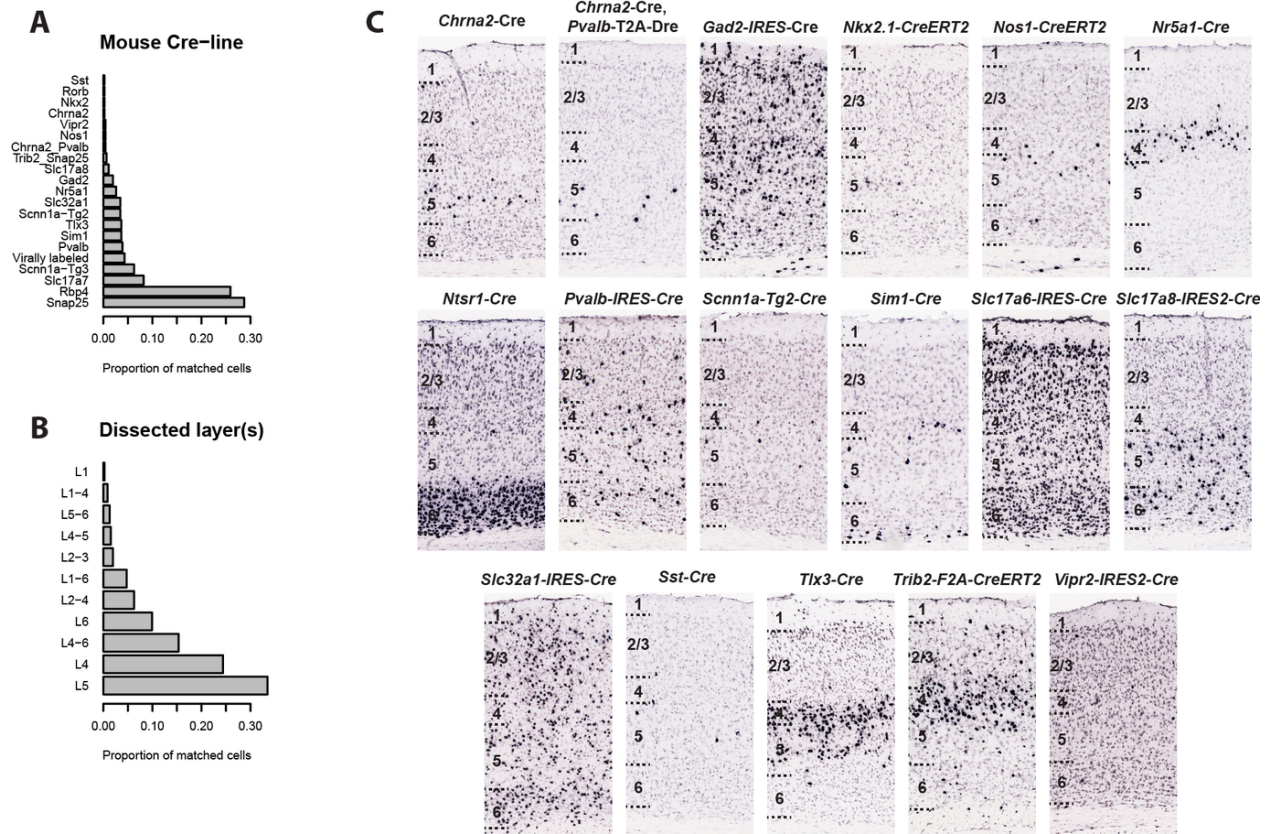


Figure S1: [Figure 1 - supplemental] Properties of 463 cells matched to nuclei. **(A)** Proportion of matched cells isolated from transgenic mouse lines that label different subsets of cortical neurons. Note that a small number of “virally labeled” cells (<5%) were FAC sorted from wild-type mice based on retrograde labeling by viral injections into various cortical and subcortical structures. **(B)** Proportion of matched cells dissected from one or more adjacent layers of cortex. **(C)** ISH images from additional mouse Cre-lines from which the best matching cells were most commonly derived. ISH images show all cortical layers within VISp.

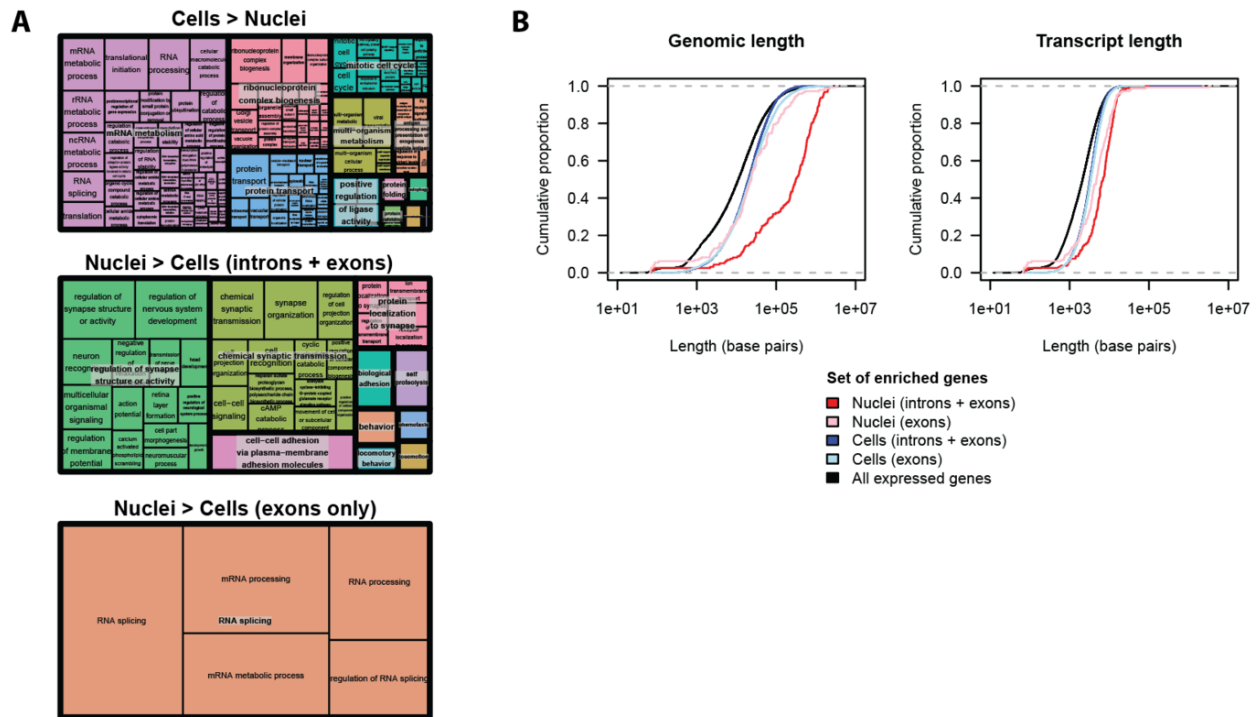


Figure S2: [Figure 2 - supplemental] Nuclear enrichment of transcripts related to neuron function can be explained by nuclear intron retention of long genes. **(A)** REVIGO (Supek et al. 2011) summaries of gene ontology (GO) enrichment of genes enriched in cells or nuclei. Including introns dramatically changes the functional categories of nuclear but not cell enriched genes. **(B)** Cumulative distribution of genomic and transcript lengths for genes enriched in nuclei and cells (fold change > 1.5) based on expression of exons or introns plus exons. Using introns plus exons, the median genomic length of nuclear enriched genes is 16-fold longer than cell enriched genes. Using exons only, there is no significant difference in genomic lengths (Kolmogorov-Smirnov test P-value = 0.27).

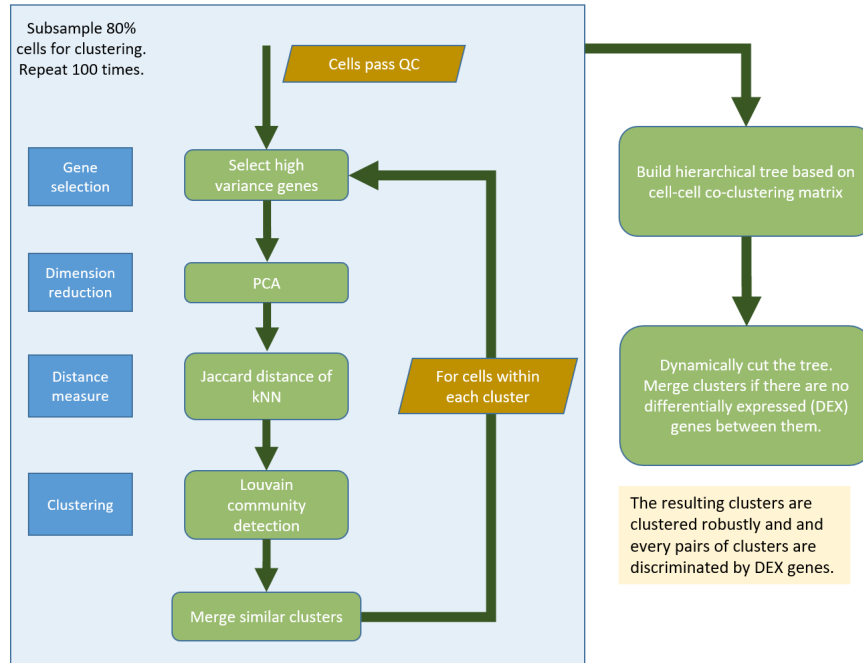


Figure S3: [Figure 3 - supplemental] Overview of single nucleus RNA-seq clustering pipeline. See methods for a detailed description of clustering steps.

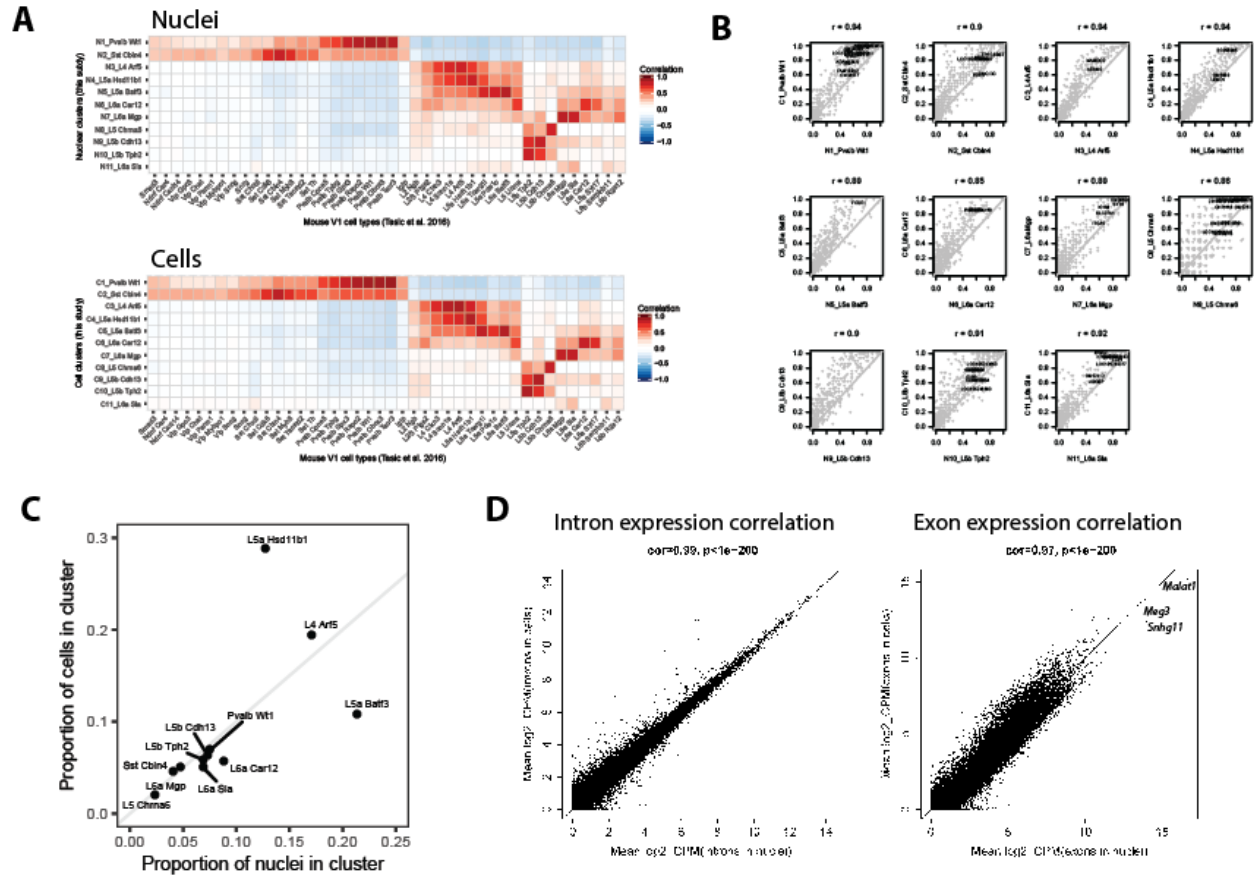


Figure S4: [Figure 4 - supplemental] Nuclear and cell clusters are well matched based on marker gene expression. **(A)** Pairwise correlations between previously reported mouse VISp cell type clusters (Tasic et al. 2016) and nuclear and cell clusters using average cluster expression of the top shared marker genes. Heatmaps show remarkably similar correlation patterns, supporting the existence of a well matched set of nuclear and cell clusters. Nuclear and cell clusters were annotated based on the reciprocal best matching published cluster name and mapped to two interneuron types and five of eight layer 5 excitatory neuron types. **(B)** Comparisons of the proportion of nuclei or cells expressing marker genes (CPM > 1) for matched pairs of clusters. Correlations are reported at the top of each scatter plot, and cell type specific markers are labeled. As expected based on Figure 2C, gene detection is consistently higher in cells than nuclei. **(C)** Matched clusters have similar proportions of nuclei and cells (except for two closely related cell types, L5a Hsd11b1 and L5a Batf3), which supports the accuracy of the initial correlation based mapping of single nuclei to cells. **(D)** Average gene expression quantified based on intronic reads is more highly correlated between cells and nuclei than expression quantified based on exonic reads, particularly for highly expressed genes. *Malat1*, *Meg3*, and *Snhg11* are the three highest expressing genes in nuclei and have consistently lower expression in cells, as expected based on their reported nuclear localization.

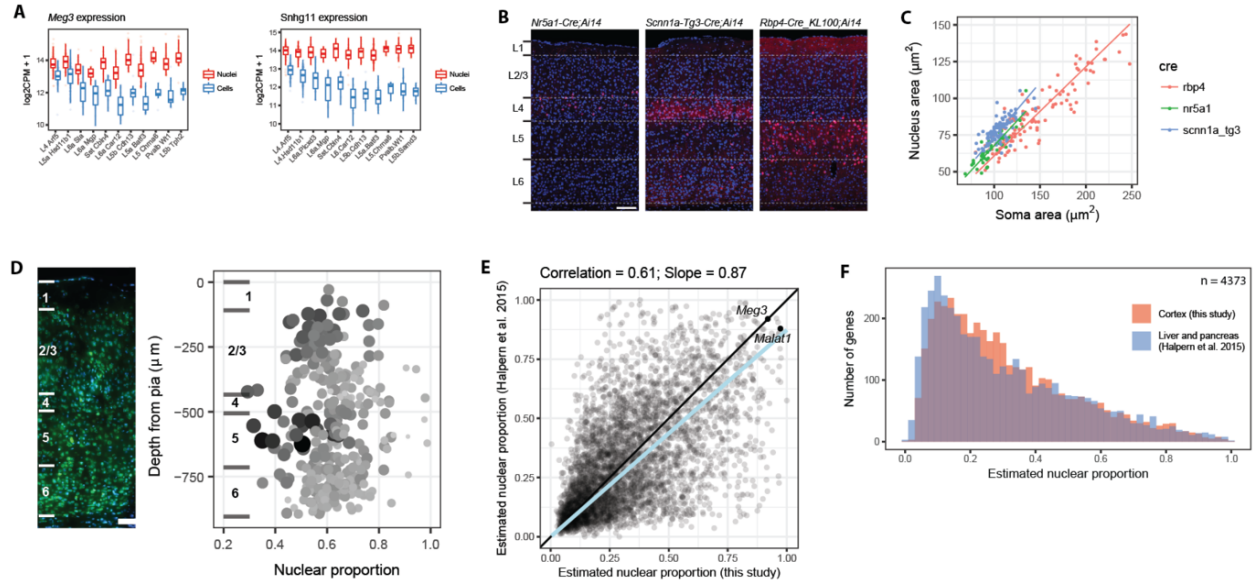


Figure S5: [Figure 5 - supplemental] Nuclear proportion estimates are supported by multiple genes and consistent with previously reported values. **(A)** Box plots of \log_2 -transformed expression of two nuclear transcripts, *Meg3* and the small nucleolar RNA *Snhg11*, in matched nuclear and cell clusters. **(B)** Representative sections of VISp from three Cre-driver mouse lines with layer boundaries, nuclei labeled with DAPI (blue), and subsets of neurons labeled with tdTomato (red). Scale bar is $100 \mu\text{m}$. **(C)** Nucleus and soma area measurements from three Cre-lines, and linear regressions to estimate nuclear proportions. **(D)** Left: Section of VISp from wild type mouse labeled with DAPI and Neurotrace 500 fluorescent Nissl stain with layer boundaries indicated by white lines. Scale bar is $100 \mu\text{m}$. Right: Nuclear proportion was quantified based on nucleus and soma area measurements and plotted as a function of cortical depth. Size and darkness of points are proportional to soma area. **(E)** Average nuclear proportions of 4,373 genes (mostly house-keeping) also expressed in mouse pancreatic beta-cells and liver cells (Halpern et al. 2015) are moderately correlated with and approximately 13% less than estimated proportions in this study. **(F)** The distributions of nuclear proportions are highly similar with slightly higher reported cytoplasmic enrichment for reported genes. Note that the matched set of genes includes 99% protein-coding genes so the distributions more closely resemble those genes in Figure 5D.

569 Supplemental Tables

570 Table S1 [Figure 2 - supplemental]. Average gene expression and detection in matched nuclei and cells.

571 Table S2 [Figure 2 - supplemental]. Differentially expressed genes in cells versus nuclei using intronic plus
572 exonic reads.

573 Table S3 [Figure 2 - supplemental]. Differentially expressed genes in cells versus nuclei using only exonic
574 reads.

575 Table S4 [Figure 2 - supplemental]. Gene ontology (GO) enrichment of differentially expressed genes in cells
576 and nuclei.

577 Table S5 [Figure 4 - supplemental]. Cre-driver line composition of cell clusters.

578 Table S6 [Figure 5 - supplemental]. Gene properties including the number of clusters with any expression,
579 maximum cluster expression, cell type marker score, and estimated nuclear proportion of transcripts.

References

- 580
- 581 Poulin, Jean-Francois, Bosiljka Tasic, Jens Hjerling-Leffler, Jeffrey M Trimarchi, and Rajeshwar Awatramani.
582 2016. “Disentangling Neural Cell Diversity Using Single-Cell Transcriptomics”. *Nature Neuroscience* 19 (9).
583 Springer Nature: 1131–41. doi:10.1038/nn.4366.
- 584 Zeng, Hongkui, and Joshua R. Sanes. 2017. “Neuronal Cell-Type Classification: Challenges Oppor-
585 tunities and the Path Forward”. *Nature Reviews Neuroscience* 18 (9). Springer Nature: 530–46.
586 doi:10.1038/nrn.2017.85.
- 587 Bernard, Amy, Staci A Sorensen, and Ed S Lein. 2009. “Shifting the Paradigm: New Approaches for
588 Characterizing and Classifying Neurons”. *Current Opinion in Neurobiology* 19 (5). Elsevier BV: 530–36.
589 doi:10.1016/j.conb.2009.09.010.
- 590 Tasic, B, V Menon, TN Nguyen, TK Kim, T Jarsky, Z Yao, B Levi, et al. 2016. “Adult Mouse Cortical Cell
591 Taxonomy Revealed by Single Cell Transcriptomics”. *Nat Neurosci* 19: 335–46.
- 592 Tasic, Bosiljka, Zizhen Yao, Kimberly A Smith, Lucas T Graybuck, and Thuc Nghi Nguyen. 2017. “Shared
593 and Distinct Transcriptomic Cell Types across Neocortical Areas”. *BioRxiv*. doi:dx.doi.org/10.1101/229542.
- 594 Zeisel, A, AB Muñoz-Manchado, S Codeluppi, P Lönnerberg, Manno G La, A Juréus, S Marques, et al.
595 2015. “Cell Types in the Mouse Cortex and Hippocampus Revealed by Single-Cell RNA-Seq”. *Science* 347:
596 1138–42.
- 597 Campbell, JN, EZ Macosko, H Fenselau, TH Pers, A Lyubetskaya, D Tenen, M Goldman, et al. 2017. “A
598 Molecular Census of Arcuate Hypothalamus and Median Eminence Cell Types”. *Nat Neurosci* 20: 484–96.
- 599 Shekhar, K, SW Lapan, IE Whitney, NM Tran, EZ Macosko, M Kowalczyk, X Adiconis, et al. 2016. “Com-
600 prehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics”. *Cell* 166: 1308–23.e30.
- 601 Macosko, EZ, A Basu, R Satija, J Nemes, K Shekhar, M Goldman, I Tirosh, et al. 2015. “Highly Parallel
602 Genome-Wide Expression Profiling of Individual Cells Using Nanoliter Droplets”. *Cell* 161: 1202–14.
- 603 Darmanis, S, SA Sloan, Y Zhang, M Enge, C Caneda, LM Shuer, Gephart MG Hayden, BA Barres, and
604 SR Quake. 2015. “A Survey of Human Brain Transcriptome Diversity at the Single Cell Level.”. *Proc Natl*
605 *Acad Sci U S A* 112: 7285–90.
- 606 Krishnaswami, Suguna Rani, Rashel V Grindberg, Mark Novotny, Pratap Venepally, Benjamin Lacar, Kunal
607 Bhutani, Sara B Linker, et al. 2016. “Using Single Nuclei for RNA-Seq to Capture the Transcriptome of
608 Postmortem Neurons”. *Nature Protocols* 11 (3). Springer Nature: 499–524. doi:10.1038/nprot.2016.015.
- 609 Lacar, Benjamin, Sara B. Linker, Baptiste N. Jaeger, Suguna Krishnaswami, Jerika Barron, Martijn Kelder,
610 Sarah Parylak, et al. 2016. “Nuclear RNA-Seq of Single Neurons Reveals Molecular Signatures of Activa-
611 tion”. *Nature Communications* 7 (April). Springer Nature: 11022. doi:10.1038/ncomms11022.
- 612 Lake, BB, R Ai, GE Kaeser, NS Salathia, YC Yung, R Liu, A Wildberg, et al. 2016. “Neuronal Subtypes
613 and Diversity Revealed by Single-Nucleus RNA Sequencing of the Human Brain”. *Science* 352: 1586–90.
- 614 Lake, BB, S Chen, BC Sos, J Fan, GE Kaeser, YC Yung, TE Duong, et al. 2017. “Integrative Single-Cell
615 Analysis of Transcriptional and Epigenetic States in the Human Adult Brain”. *Nat Biotechnol*.
- 616 Habib, N, Y Li, M Heidenreich, L Swiech, I Avraham-Davidi, JJ Trombetta, C Hession, F Zhang, and A
617 Regev. 2016. “Div-Seq: Single-Nucleus RNA-Seq Reveals Dynamics of Rare Adult Newborn Neurons.”.
618 *Science* 353: 925–28.
- 619 Lake, BB, S Codeluppi, YC Yung, D Gao, J Chun, PV Kharchenko, S Linnarsson, and K Zhang. 2017. “A
620 Comparative Strategy for Single-Nucleus and Single-Cell Transcriptomes Confirms Accuracy in Predicted
621 Cell-Type Expression from Nuclear RNA”. *Sci Rep* 7: 6031.

- 622 Dobin, A, CA Davis, F Schlesinger, J Drenkow, C Zaleski, S Jha, P Batut, M Chaisson, and TR Gingeras.
623 2013. “STAR: Ultrafast Universal RNA-Seq Aligner”. *Bioinformatics* 29: 15–21.
- 624 Kharchenko, PV, L Silberstein, and DT Scadden. 2014. “Bayesian Approach to Single-Cell Differential
625 Expression Analysis”. *Nat Methods* 11: 740–42.
- 626 Sigl-Glöckner, J, and M Brecht. 2017. “Polyploidy and the Cellular and Areal Diversity of Rat Cortical
627 Layer 5 Pyramidal Neurons”. *Cell Rep* 20: 2575–83.
- 628 Lin, N, KY Chang, Z Li, K Gates, ZA Rana, J Dang, D Zhang, et al. 2014. “An Evolutionarily Conserved
629 Long Noncoding RNA TUNA Controls Pluripotency and Neural Lineage Commitment”. *Mol Cell* 53:
630 1005–19.
- 631 Halpern, Keren B, I Caspi, D Lemze, M Levy, S Landen, E Elinav, I Ulitsky, and S Itzkovitz. 2015. “Nuclear
632 Retention of mRNA in Mammalian Tissues”. *Cell Rep* 13: 2653–62.
- 633 Djebali, S, CA Davis, A Merkel, A Dobin, T Lassmann, A Mortazavi, A Tanzer, et al. 2012. “Landscape of
634 Transcription in Human Cells”. *Nature* 489: 101–8.
- 635 Gabel, HW, B Kinde, H Stroud, CS Gilbert, DA Harmin, NR Kastan, M Hemberg, DH Ebert, and ME
636 Greenberg. 2015. “Disruption of DNA-Methylation-Dependent Long Gene Repression in Rett Syndrome.”.
637 *Nature* 522: 89–93.
- 638 Mauger, O, F Lemoine, and P Scheiffele. 2016. “Targeted Intron Retention and Excision for Rapid Gene
639 Regulation in Response to Neuronal Activity.”. *Neuron* 92: 1266–78.
- 640 Ecker, JR, DH Geschwind, AR Kriegstein, J Ngai, P Osten, D Polioudakis, A Regev, N Sestan, IR Wick-
641 ersham, and H Zeng. 2017. “The BRAIN Initiative Cell Census Consortium: Lessons Learned toward
642 Generating a Comprehensive Brain Cell Atlas.”. *Neuron* 96: 542–57.
- 643 Regev, Aviv, Sarah A Teichmann, Eric S Lander, Ido Amit, Christophe Benoist, Ewan Birney, Bernd
644 Bodenmiller, et al. 2017. “Science Forum: The Human Cell Atlas”. *ELife* 6 (December). eLife Sciences
645 Organisation Ltd. doi:10.7554/elife.27041.
- 646 Baker, SC, SR Bauer, RP Beyer, JD Brenton, B Bromley, J Burrill, H Causton, et al. 2005. “The External
647 RNA Controls Consortium: a Progress Report.”. *Nat Methods* 2: 731–34.
- 648 Risso, D, J Ngai, TP Speed, and S Dudoit. 2014. “Normalization of RNA-Seq Data Using Factor Analysis
649 of Control Genes or Samples.”. *Nat Biotechnol* 32: 896–902.
- 650 Aronesty, Erik. 2011. “Ea-Utils : Command-Line Tools for Processing Biological Sequencing Data;
651 <https://Github.com/ExpressionAnalysis/Ea-Utils>”.
- 652 Lawrence, M, W Huber, H Pagès, P Aboyoun, M Carlson, R Gentleman, MT Morgan, and VJ Carey. 2013.
653 “Software for Computing and Annotating Genomic Ranges”. *PLoS Comput Biol* 9: e1003118.
- 654 n.d. <http://hms-dbmi.github.io/scde/diffexp.html>. <http://hms-dbmi.github.io/scde/diffexp.html>.
- 655 [html](http://hms-dbmi.github.io/scde/diffexp.html).
- 656 Wickham, Hadley. 2009. *Ggplot2: Elegant Graphics for Data Analysis*. 2nd ed. Springer Publishing
657 Company, Incorporated.
- 658 Ritchie, Matthew E, Belinda Phipson, Di Wu, Yifang Hu, Charity W Law, Wei Shi, and Gordon K Smyth.
659 2015. “Limma Powers Differential Expression Analyses for RNA-Sequencing and Microarray Studies”. *Nu-
660 cleic Acids Research* 43 (7). Oxford University Press: e47–e47.
- 661 Chambers, J. M., A. Freeny, and R. M. Heiberger. 1992. *Analysis of Variance; Designed Experiments*.
662 Edited by J. M. Chambers and T. J. Hastie. *Statistical Models in S*. Wadsworth & Brooks/Cole.

- 663 Chen, J., E. E. Bardes, B. J. Aronow, and A. G. Jegga. 2009. “ToppGene Suite for Gene List Enrichment
664 Analysis and Candidate Gene Prioritization”. *Nucleic Acids Research* 37 (Web Server). Oxford University
665 Press (OUP): W305–W311. doi:10.1093/nar/gkp427.
- 666 Supek, Fran, Matko Bošnjak, Nives Škunca, and Tomislav Šmuc. 2011. “REVIGO Summarizes and Visual-
667 izes Long Lists of Gene Ontology Terms”. Edited by Cynthia Gibas. *PLoS ONE* 6 (7). Public Library of
668 Science (PLoS): e21800. doi:10.1371/journal.pone.0021800.
- 669 Levine, JH, EF Simonds, SC Bendall, KL Davis, el-AD Amir, MD Tadmor, O Litvin, et al. 2015. “Data-
670 Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells That Correlate with Prognosis”. *Cell*
671 162: 184–97.
- 672 Jackson, Donald A. 1993. “Stopping Rules in Principal Components Analysis: A Comparison of Heuristical
673 and Statistical Approaches”. *Ecology* 74 (8): 2204–14. doi:10.2307/1939574.
- 674 Fortunato, Santo, and M. Barthelemy. 2007. “Resolution Limit in Community Detection”. *Proceedings of*
675 *the National Academy of Sciences* 104 (1): 36–41. doi:10.1073/pnas.0605965104.
- 676 Reichardt, Jörg, and Stefan Bornholdt. 2006. “Statistical Mechanics of Community Detection”. *Physical*
677 *Review E* 74 (1). APS: 016110.
- 678 Guimerà, Roger, Marta Sales-Pardo, and Luís A Nunes Amaral. 2004. “Modularity from
679 Fluctuations in Random Graphs and Complex Networks”. *Physical Review E* 70 (2): 25101.
680 doi:10.1103/PhysRevE.70.025101.
- 681 Langfelder, Peter, Bin Zhang, and Steve Horvath. 2007. “Defining Clusters from a Hierarchical Cluster
682 Tree: the Dynamic Tree Cut Package for R”. *Bioinformatics* 24 (5). Oxford University Press: 719–20.
- 683 Lein, ES, MJ Hawrylycz, N Ao, M Ayres, A Bensinger, A Bernard, AF Boe, et al. 2007. “Genome-Wide
684 Atlas of Gene Expression in the Adult Mouse Brain”. *Nature* 445: 168–76.
- 685 Paxinos, George, and others. 2013. *Paxinos and Franklin’s the Mouse Brain in Stereotaxic Coordinates*.
686 Academic Press.
- 687 Lamprecht, MR, DM Sabatini, and AE Carpenter. 2007. “CellProfiler: Free, Versatile Software for Auto-
688 mated Biological Image Analysis”. *Biotechniques* 42: 71–75.
- 689 Supek, F, M Bošnjak, N Škunca, and T Šmuc. 2011. “REVIGO Summarizes and Visualizes Long Lists of
690 Gene Ontology Terms”. *PLoS One* 6: e21800.