# Mapping the transcriptional diversity of genetically and anatomically defined cell populations in the mouse brain.

**Ken Sugino**[1]*, **Erin Clark**[2,3], **Anton Schulmann**[1,3], **Yasuyuki Shima**[2], **Lihua Wang**[1], **David L. Hunt**[1], **Bryan M. Hooks**[1], **Dimitri Tränkner**[1], **Jayaram Chandrashekar**[1], **Serge Picard**[1], **Andrew Lemire**[1], **Nelson Spruston**[1], **Adam Hantman**[1], **Sacha B. Nelson**[2]*

**\*For correspondence:**
suginok@janelia.hhmi.org (KS);
hantmana@janelia.hhmi.org (AH);
nelson@brandeis.edu (SN)

[1]Janelia Research Campus; [2]Brandeis Universiy; [3]equal contribution

## Abstract

Understanding the principles governing neuronal diversity is a fundamental goal for neuroscience. Here we provide an anatomical and transcriptomic database of nearly 200 genetically identified cell populations. By separately analyzing the robustness and pattern of expression differences across these cell populations, we identify two gene classes contributing distinctly to neuronal diversity. Short homeobox transcription factors distinguish neuronal populations combinatorially, and exhibit extremely low transcriptional noise, enabling highly robust expression differences. Long neuronal effector genes, such as channels and cell adhesion molecules, contribute disproportionately to neuronal diversity, based on their patterns rather than robustness of expression differences. By linking transcriptional identity to genetic strains and anatomical atlases we provide an extensive resource for further investigation of mouse neuronal cell types.

## Introduction

The extraordinary diversity of vertebrate neurons has been appreciated since the proposal of the neuron doctrine (*Ramon y Cajal, 1894*). Classically, this diversity was characterized by neuronal morphology, physiology, and circuit connectivity, but increasingly, defined genetically through driver and reporter strains (*Gong et al., 2003*; *Madisen et al., 2009*; *Taniguchi et al., 2011*; *Shima et al., 2016*) or genomically by their genome-wide expression profiles. The first genome-wide studies of mammalian neuronal diversity employed *in situ* hybridization (*Lein et al., 2006*) or microarrays (*Sugino et al., 2005*; *Doyle et al., 2008*), while more recent studies have utilized advances in single cell (SC) RNA-seq (*Zeisel et al., 2015*, *2018*; *Tasic et al., 2016*, *2018*; *Paul et al., 2017*). In theory, SC RNA-seq can be applied in an unbiased fashion to discover all cell types that comprise a tissue, but manipulation of these cell types to better understand their biological composition and function often require the use of genetic tools such as mouse driver strains. Differences in techniques for cell isolation, library preparation or clustering have not yet led to a consensus view of the number or identity of the neuronal cell types comprising most parts of the mouse nervous system. Furthermore, the relationship between cell populations defined transcriptionally and those that can be specified genetically and anatomically using existing strains has received far less attention (though see *Tasic et al. 2018*).

Here we attempt to strengthen the link between genomically and genetically defined cell types in the mouse brain by performing RNA-seq on a large set of genetically identified and fluorescently labeled neurons from micro-dissected brain regions. In total, we profiled 179 sorted neuronal populations and 15 nonneuronal populations. Because each sample of sorted cells may contain more than one "atomic" cell type, we refer to these as genetically- and anatomically-identified cell populations (GACPs). To assess homogeneity, we quantitatively compared our sorted cell populations to publicly available single cell datasets, which revealed a comparable level of homogeneity, but a much lower level of noise in the sorted population profiles.

Although neuronal diversity has long been recognized, the question of how this diversity arises has not been addressed sufficiently in a genomic context (*Arendt et al., 2016*; *Muotri and Gage, 2006*). We identify two different sets of genes that distinguish GACPs based on the robustness or pattern of their expression differences. The most robust expression differences are those of homeobox transcription factors. These genes also have the lowest transcriptional noise suggesting differential chromatin regulation. Chromatin accessibility measurements reveal that the promoters and gene bodies of these genes are indeed more closed. In contrast, the genes capable of distinguishing the largest numbers of GACPs are neuronal effector genes like receptors, ion channels and cell adhesion molecules. Interestingly, genes defined by the robustness and patterns of their expression differences also differ in their transcript length. Genes with robust, low noise expression tend to be shorter, while genes with the greatest capacity to distinguish populations tend to be longer.

Here we provide important new resources for mapping brain cell types including a large set of low-noise profiles from genetically identified neurons, anatomical maps of their distributions, and a method to compare and contextualize single cell RNA-seq datasets. We implement a novel strategy to mine information from large surveys of cell types, and demonstrate the utility of this strategy in generating specific biological insights into the genes contributing to neuronal diversity.

## Results

### A dataset of genetically-identified neuronal transcriptomes

To identify genes contributing most to mammalian neuronal diversity, we collected transcriptomes from 179 genetically and anatomically identified populations of neurons and 15 populations of nonneuronal cells in mice (Table 1; Figure 1; Figure 1 Supplement 1; Supplementary File 1,2). The great majority (186/194) were identified both genetically and anatomically, with the remaining identified only anatomically, by their location and projection patterns. Each collected population represents a group of fluorescently labeled cells dissociated and sorted from a specific micro-dissected region of the mouse brain or other tissue. The pipeline for collecting GACP transcriptomes is depicted in Figure 1A (see Methods for additional details). Mouse lines were first characterized by generating a high-resolution atlas of reporter expression (Figure 1B) then, regions containing labeled cells with uniform morphology were chosen for sorting and RNA-seq. In total, we sequenced 2.3 trillion bp in 565 libraries. This effort (NeuroSeq) constitutes the largest and most diverse single collection of genetically identified cell populations profiled by RNA-seq. The raw data is deposited to NCBI GEO (GSE79238). The processed data, including anatomical atlases, RNA-seq coverage, and TPM are available at http://neuroseq.janelia.org (Figure 1C).

To determine the sensitivity of our transcriptional profiling, we used ERCC spike-ins. Amplified RNA libraries had an average sensitivity (50% detection) of 23 copy*kbp of ERCC spike-ins across all libraries (Figure 1D). Since manually sorted samples had 132±16 cells (mean± SEM), this indicates our pipeline had the sensitivity to detect a single copy of a transcript per cell 80% of the time. This high sensitivity allowed for deep transcriptional profiling in our diverse set of cell populations.

To assess the extent of contamination in the dataset, we checked expression levels of marker genes for several nonneuronal cell populations (Figure 1 Supplement 2B). As previously shown (*Okaty et al., 2011*), manual sorting produced, in general, extremely clean data.

To assess the homogeneity of the sorted, pooled samples, we compared our datasets to publicly available single cell (SC) datasets. To compare across different datasets, we used a method based on linear decomposition by non-negative least squares (NNLS) (See Figure 2 and Figure 2 Supplements 1-6). This method tests the degree to which individual profiles can be decomposed into linear mixtures of profiles from another dataset. Such mixtures or impurities can arise in at least two ways (Figure 2A): by pooling similar cell types prior to sequencing in the case of sorted datasets, or by pooling similar profiles after sequencing, at the clustering stage, in the case of SC datasets. Although NNLS is a widely used decomposition procedure, it has not previously been applied to expression profiles. Therefore, we performed a number of control experiments to validate its use. First, we cross-validated the decompositions by dividing each dataset in half and testing the ability to decompose one half by the other (Figure 2 supplement 1). This revealed that some NeuroSeq samples had overlapping coefficients and so could not well distinguished. For example, pairs of populations identified in layer 2/3 of two different regions in the same strain (AI.L23_glu_P157 / ORBm.L23_glu_P157) or by retrogradely labeled cells in the same layer and region from two different targets (SSp.L23_glu_M1.inj / SSp.L23_glu_S2.inj and SSp.L5_glu_BPn.inj / SSp.L5_glu_IRT.in) were hard to distinguish. On the other hand, overlapping coefficients were also present for some pairs of cell populations in the SC datasets (such as Oligo Serpinb1a / Oligo Synpr in the Tasic dataset and MGL1 / MGL2 / MGL3 in the Zeisel dataset). On average the purity, defined as how well a single sample can be decomposed into the most closely corresponding sample, was similar across the three datasets (Figure 2 Supplement 1D). As a second control, we demonstrated that NNLS decomposition could be used to recover the numbers of cell types isolated from distinct strains in a SC dataset, after mixing these profiles together, despite the fact that this information was not included in the fitting procedure (Figure 2 Supplement 2). Finally, NNLS (Figure 2B,C) produced comparable or cleaner decompositions than a competing Random Forest algorithm (Figure 2 Supplement 6). These results indicate that NNLS can be used to reliably decompose mixtures of cellular profiles. Similar average coefficients (i.e. similar purity) were obtained for decompositions of the NeuroSeq data by SC datasets and by decomposing these datasets by each other (Figure 2, Figure 2 Supplements 3-6). Hence our decomposition results indicate that although heterogeneity may exist in some of our sorted samples, it is comparable to the inaccuracies introduced by clustering SC profiles.

Since merging or splitting of closely related clusters either prior to sequencing or during the clustering process can lead to poor discrimination between samples, we also measured the separability of cell population profiles obtained in each study (Figure 2 Supplement 7). As expected, the clusters of sorted population samples, which are averages across one hundred cells or more, were much more cleanly separable than SC clusters. Taken together, NNLS decomposition and separability provide a quantitative framework for assessing the trade-offs between homogeneity and reproducibility when measuring population transcriptomes from GACPs and SCs.

To demonstrate the utility of the dataset, made possible by its broad sampling of neuronal populations, we extracted pan-neuronal genes (genes expressed commonly in all neuronal populations but expressed at lower levels or not at all in nonneuronal cell populations; Figure 1 Supplement 3). Here, broad sampling of cell populations is essential to avoid false positives (*Zhang et al., 2014b*; *Mo et al., 2015*; *Stefanakis et al., 2015*). Because of the high sensitivity and low noise, we were able to be conservative and exclude genes expressed in most but not all neuron types. Extracted pan-neuronal genes contain well known genes such as *Eno2* (Enolase2), which is the neuronal form of Enolase required for the Krebs cycle, *Slc2a3* (chloride transporter) required for inhibitory transmission, and *Atp1a3* (ATPase Na+/K+ transporting subunit alpha 3) which belongs to the complex responsible for maintaining electrochemical gradients across the membrane, as well as genes not previously known to be pan-neuronal, such as *2900011O08Rik* (now called Migration Inhibitory Protein; *Zhang et al. 2014a*). Synaptic genes are often differentially expressed among neurons, but interestingly, some were included in this pan-neuronal list such as *Syn1, Stx1b, Stxbp1, Sv2a*, and *Vamp2*. These appear to be common synaptic components, and highlight essential parts of

140 these complexes. Thus, the dataset should be useful for many other applications, especially those
141 requiring comparisons across a wide variety of neuronal cell types.

142 **Metrics to quantify diversity**

143 Analysis of expression differences between individual groups is the basis of most profiling efforts.
144 Variance-based metrics, such as Analysis of Variance (ANOVA) F-Value, or coefficient of variation (CV)
145 are commonly used for this purpose. However, these metrics are jointly affected by the pattern of
146 differential expression and the robustness of the differences, and so cannot readily separate these
147 two features (Figure 3,4; Figure 3 Supplement 1). Since these features may differ in their biological
148 significance, we searched for the simplest way to quantitatively separate them. This led us to adopt
149 two easily calculated variants of widely used metrics for differential expression and fold-change.

150 To quantify the contribution of each gene to cell type diversity, we measured the fraction of cell
151 population pairs in which the gene is differentially expressed. (For differential analysis, the limma-
152 voom framework was used, see Methods). This differentially expressed fraction (DEF) is closely
153 related to the Gini-Simpson diversity index (*Simpson, 1949*) widely used in ecology to measure
154 species diversity in a community (see Appendix 1). DEF ranges from 0 to 1. The maximum observed
155 value of 0.65 indicates that the gene distinguishes 65% of the pairs, while a value of 0 indicates that
156 the gene distinguishes none (i.e., it is expressed at similar levels in all cells). DEF is easy to calculate
157 and approximates the mutual information (MI) between expression levels and cell populations
158 (Appendix 1).

159 The robustness of an expression difference depends on its magnitude relative to the underlying
160 noise. Robustness is often quantified as a Signal-to-Noise-Ratio (SNR). Since the signals we are
161 interested in are the gene expression differences distinguishing cell types, we computed the ratio
162 of the mean fold-change expression differences between distinguished pairs to the mean fold-
163 change between undistinguished pairs. This fold-change ratio (FCR) indicates the robustness of
164 pair distinctions, but is independent of the number of pairs distinguished. High FCR genes robustly
165 distinguish cell populations and are therefore suitable as "marker genes".

166 Unlike DEF and FCR, variance-based methods like ANOVA F-values and CV are either affected
167 by both MI and SNR (ANOVA; Figure 4A-C and Figure 3 Supplement 1) or by neither (CV; Figure 3
168 Supplement 1). The fact that ANOVA does not distinguish between information content and SNR
169 can be appreciated from the fact that high-ANOVA genes (Figure 4A-C) include both high DEF and
170 high FCR genes. Therefore, DEF and FCR are useful because they provide independent measures of
171 the robustness and magnitude of differential expression between cell populations.

172 To determine the types of genes most differentially expressed (highest DEF) and most robustly
173 different (highest FCR) between cell populations, we performed over-representation analysis using
174 the HUGO Gene Groups (*Braschi et al. 2018*, Figure 4D,E). The most robust expression differences
175 (highest FCR) were those of homeobox transcription factors (TFs) and G-protein coupled receptors
176 (GPCRs; Figure 4D). High DEF genes are enriched for neuronal effector genes including receptors, ion
177 channels and cell adhesion molecules (Figure 4E). High FCR and High DEF enrichments were based
178 on the HUGO gene groups, but similar results were obtained using the PANTHER gene families
179 (*Mi et al., 2016*) and Gene Ontology annotations (*Ashburner et al. 2000*, Figure 4 Supplement 1).
180 In the case of the high FCR genes, the Gene Ontology categories differed, since this ontology
181 lacks a separate category for homeobox transcription factors. Instead multiple parent categories
182 (e.g. sequence-specific DNA binding, RNA polymerase II regulatory region DNA binding etc.) were
183 overrepresented.

184 Thus, using these two simple metrics we identify synaptic and signaling genes as the most
185 differentially expressed, and homeobox TFs and GPCRs as the most robustly distinguishing families
186 of genes. These two categories of genes drive neuronal diversity by endowing neuronal cell types
187 with specialized signaling and connectivity phenotypes, and by orchestrating cell type-specific
188 patterns of transcription. In addition, their distinct contributions to distinguishing neuronal types
189 suggests possible differences in the regulation of these two categories of genes.

**Table 1.** Summary of Profiled Samples.

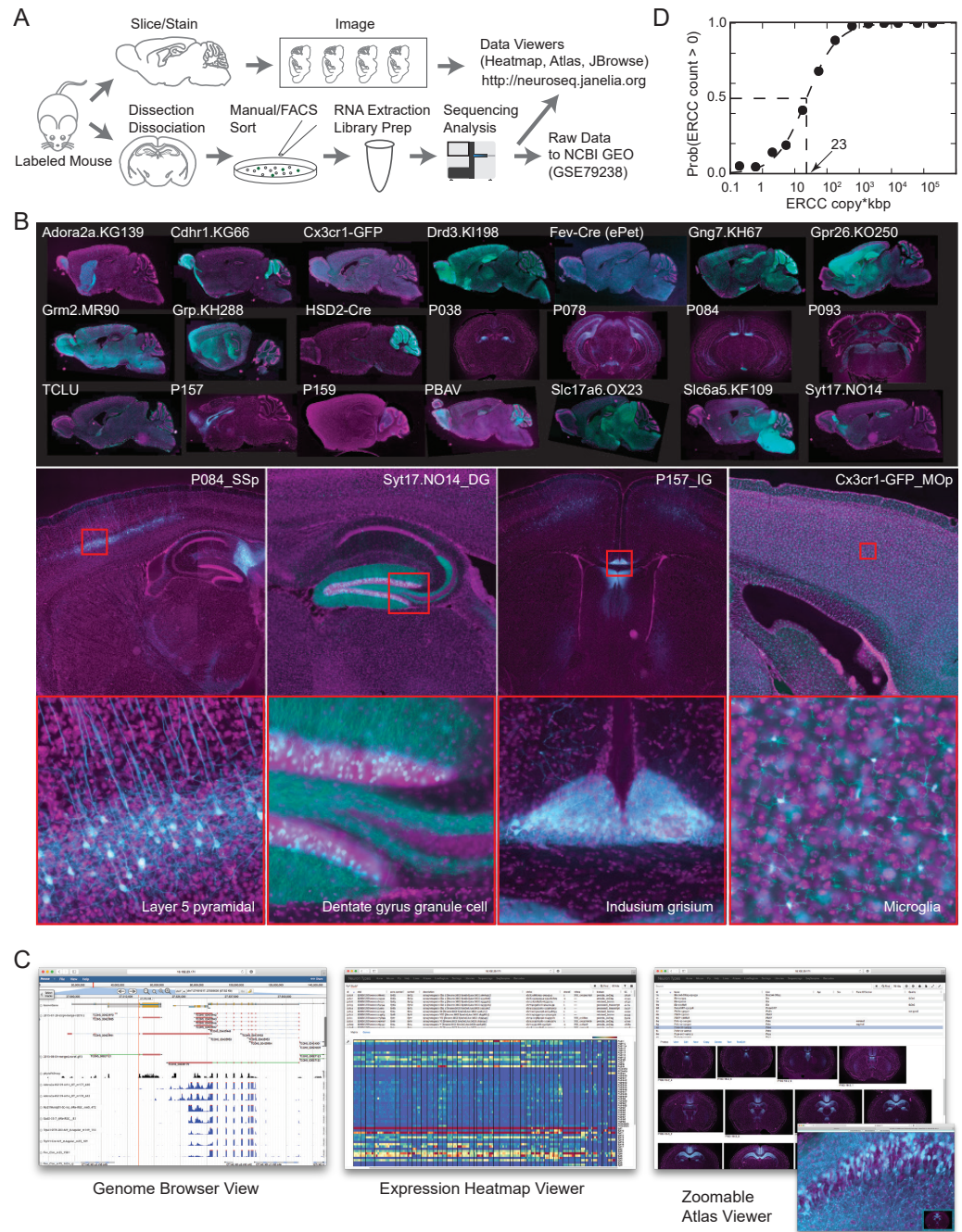| | region/type | transmitter | #groups | subregions | #samples |
|---|---|---|---|---|---|
| CNS neurons | Olfactory (OLF) | glu | 10 | AOBmi,MOBgl,PIR,AOB,COAp | 30 |
| | | GABA | 4 | AOBgr,MOBgr,MOBmi | 11 |
| | Isocortex | glu | 22 | VISp,AI,MOp5,MO,VISp6a,SSp,SSs,ECT,ORBm,RSPv | 68 |
| | | GABA | 3 | Isocortex,SSp (Sst+, Pvalb+) | 7 |
| | | glu,GABA | 1 | RSPv | 3 |
| | Subplate (CTXsp) | glu | 1 | CLA | 4 |
| | Hippocampus (HPF) | glu | 24 | CA1,CA1sp,CA2,CA3,CA3sp,DG,DG-sg,SUBd-sp,IG | 65 |
| | | GABA | 4 | CA3,CA,CA1 (Sst+, Pvalb+) | 12 |
| | Striatum (STR) | GABA | 12 | ACB,OT,CEAm,CEAl,islm,isl,CP | 33 |
| | Pallidum (PAL) | GABA | 1 | BST | 4 |
| | Thalamus (TH) | glu | 11 | PVT,CL,AMd,LGd,PCN,AV,VPM,AD | 29 |
| | Hypothalamus (HY) | glu | 11 | LHA,MM,PVHd,SO,DMHp,PVH,PVHp | 36 |
| | | GABA | 4 | ARH,MPN,SCH | 15 |
| | | glu,GABA | 2 | SFO | 3 |
| | Midbrain (MB) | DA | 2 | SNc,VTA | 5 |
| | | glu | 2 | SCm,IC | 6 |
| | | 5HT | 2 | DR | 10 |
| | | GABA | 1 | PAG | 4 |
| | | glu,DA | 1 | VTA | 3 |
| | Pons (P) | glu | 7 | PBl,PG | 22 |
| | | NE | 1 | LC | 2 |
| | | 5HT | 2 | CSm | 7 |
| | Medulla (MY) | GABA | 7 | AP,NTS,MV,NTSge,DCO | 18 |
| | | glu | 6 | NTSm,IO,ECU,LRNm | 20 |
| | | ACh | 2 | DMX,VII | 6 |
| | | 5HT | 1 | RPA | 3 |
| | | GABA,5HT | 1 | RPA | 4 |
| | | glu,GABA | 1 | PRP | 3 |
| | Cerebellum (CB) | GABA | 10 | CUL4, 5mo,CUL4, 5pu,CUL4, 5gr,PYRpu | 25 |
| | | glu | 4 | CUL4, 5gr,NODgr | 10 |
| | Retina | glu | 5 | ganglion cells (MTN,LGN,SC projecting) | 14 |
| | Spinal Cord | glu | 1 | Lumbar (L1-L5) dorsal part | 3 |
| | | GABA | 4 | Lumbar (L1-L5) dorsal part, central part | 12 |
| PNS | Jugular | glu | 2 | (TrpV1+) | 7 |
| | Dorsal root ganglion (DRG) | glu | 2 | (TrpV1+, Pvalb+) | 5 |
| | Olfactory sensory neurons (OE) | glu | 4 | MOE,VNO | 9 |
| nonneuron | Microglia | | 2 | MOp5(Isocortex),UVU(CB) (Cx3cr1+) | 6 |
| | Astrocytes | | 1 | Isocortex (GFAP+) | 4 |
| | Ependyma | | 1 | Choroid Plexus | 2 |
| | Ependyma | | 2 | Lateral ventricle (Rarres2+) | 6 |
| | Epithelial | | 1 | Blood vessel (Isocortex) (Apod+,Bgn+) | 3 |
| | Epithelial | | 1 | olfactory epithelium | 2 |
| | Progenitor | | 1 | DG (POMC+) | 3 |
| | Pituitary | | 1 | (POMC+) | 3 |
| non brain | Pancreas | | 2 | Acinar cell, beta cell | 7 |
| | Myofiber | | 2 | Extensor digitorum longus muscle | 7 |
| | Brown adipose cell | | 1 | Brown adipose cell from neck. | 4 |
| | | total | 194 | | 565 |

**Figure 1. The NeuroSeq dataset. (A)** Schema of pipeline for anatomical and genomic data collection. **(B)** Example sections from atlases at low (top), medium (middle) and high (bottom) magnifications. **(C)** Web tools available at http://neuroseq.janelia.org **(D)** Sensitivity of library preparation measured from ERCC detection across all libraries. The 50% detection sensitivity of the assay itself was 23 copy*kbp.
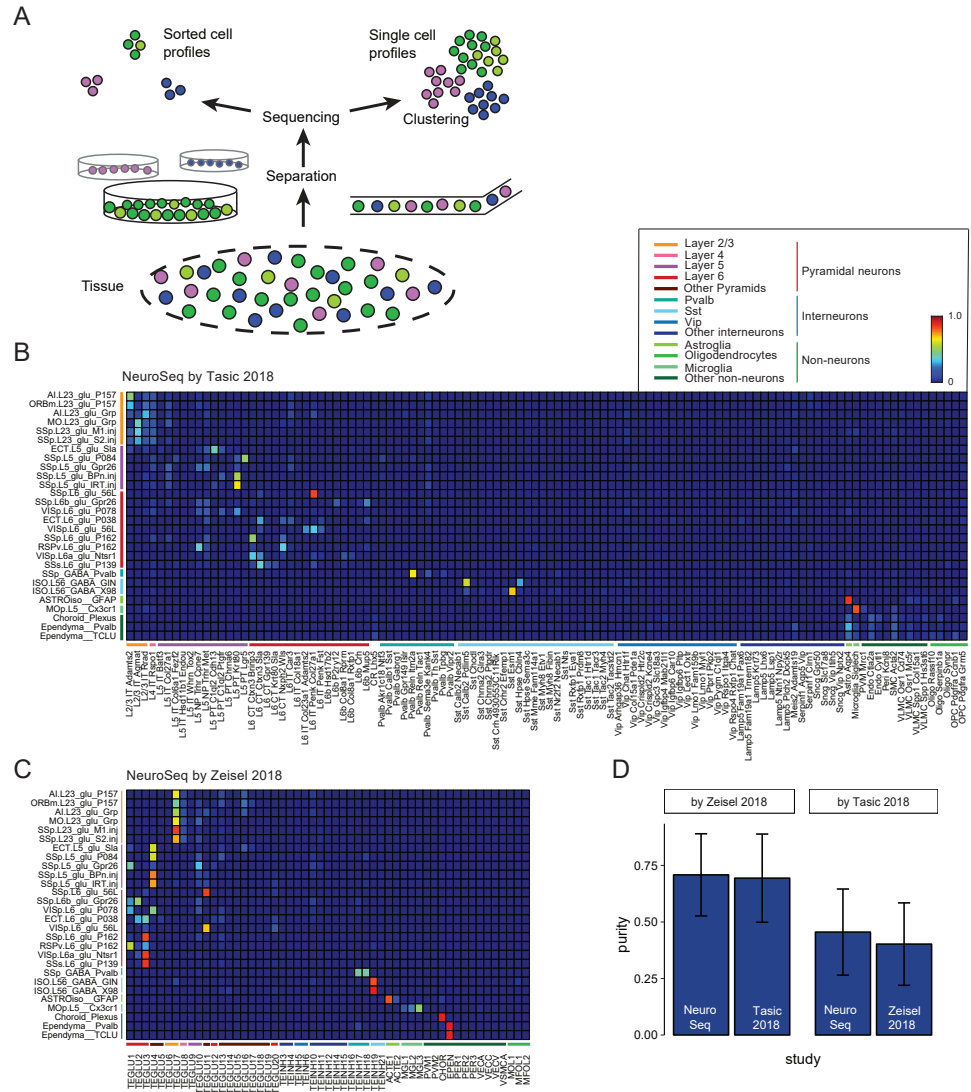
**Figure 2. Decomposition by non-negative least squares (NNLS) fitting. (A)** Diagram illustrating potential sources of heterogeneity at the separation phase in profiles from sorted cells (left) or at the clustering phase in profiles from single cells (right). **(B,C)** NNLS coefficients of NeuroSeq cell populations decomposed by two scRNA-seq datasets: (*Tasic et al., 2018*; *Zeisel et al., 2018*). **(D)** Mean purity scores for NeuroSeq and SC datasets. The purity score for a sample is defined as the ratio of the highest coefficient to the sum of all coefficients. Error bars are Std. Dev.
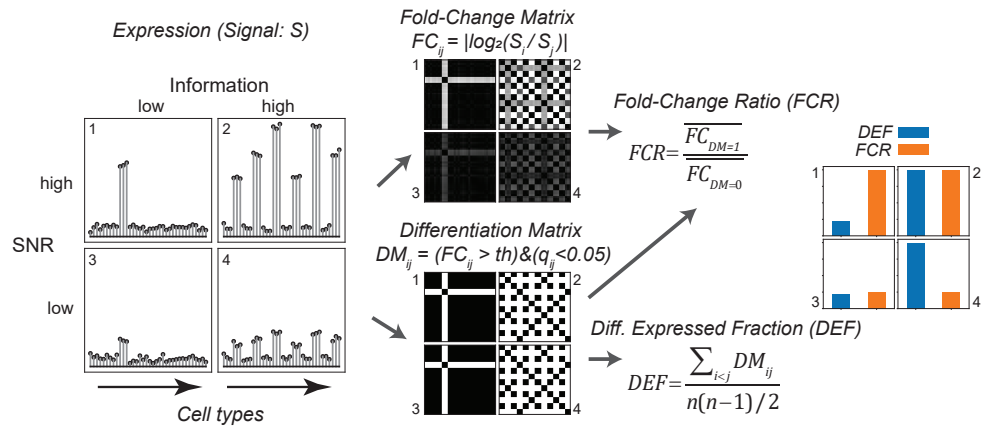
*Expression (Signal: S)*

Information
low   high

SNR

high

low

Cell types

*Fold-Change Matrix*
$FC_{ij} = |log_2(S_i/S_j)|$

*Differentiation Matrix*
$DM_{ij} = (FC_{ij} > th) \& (q_{ij} < 0.05)$

*Fold-Change Ratio (FCR)*

$$FCR = \frac{\overline{FC_{DM=1}}}{\overline{FC_{DM=0}}}$$

*Diff. Expressed Fraction (DEF)*

$$DEF = \frac{\sum_{i<j} DM_{ij}}{n(n-1)/2}$$

DEF
FCR

**Figure 3. Gene expression metrics related to information content and robustness (Left)** Cartoon illustrating the process of calculating fold-change ratio (FCR) and differentially expressed fraction (DEF) for four different hypothetical genes that differ in the information content (2&4 vs. 1&3) and signal-to-noise ratio (SNR; 1&2 vs. 3&4) of their expression patterns across cell populations. **(Middle)** Expression signals are used to construct matrices for each gene of the log fold-changes between populations (fold-change matrix) and the distinctions between populations based on those differences (Differentiation Matrix; DM; see Methods). **(Right)** The differentially expressed fraction (DEF) is the fraction of the total pairs of cell populations distinguished (i.e. of nonzero values in DM excluding diagonal). The fold-change ratio (FCR) is the average expression difference between distinguished pairs divided by the average expression difference between undistinguished pairs. Orange and blue bars show that the resulting DEF and FCR calculations capture the variations in information and SNR across the four genes.

## Homeobox TFs have the highest SNRs and can form a combinatorial code for cell populations

FCR, like SNR, is a ratio between signal and noise, and so can reflect high expression levels in most ON cell types (high signal), low expression levels in most OFF cell types (low noise), or both. Homeobox genes are not among the most abundantly expressed genes. Their average expression levels (~30 FPKM) are significantly lower than, for example, those of neuropeptides (~90 FPKM). This suggests that the high FCR of homeobox TFs depend more on low noise than high signal. In fact, many homeobox TFs have uniformly low expression in OFF cell types (Figure 5A top). We quantified this "OFF noise" for all genes and found that homeobox genes are enriched among genes that have both low OFF noise and at least moderate ON expression levels (red dashed region in Figure 5B; see also Figure 5 Supplements 1,2). Homeobox genes were not enriched in a group of high OFF noise genes (blue dashed region in Figure 5B; data not shown) that was matched for maximum expression level (Figure 5 Supplement 1C). The enrichment of homeoboxes was also observable in two of the single cell datasets encompassing multiple brain regions (Figure 5 Supplement 3).

Tight control of expression may reflect closed chromatin. To test this we measured chromatin accessibility using ATAC-seq (see Methods). As expected, compared to high-noise genes (Figure 5C bottom), genes with low OFF noise had fewer and smaller peaks within the vicinity of their transcription start site (TSS) and gene body (Figure 5C top, Figure 5D), consistent with the idea that chromatin accessibility contributes to their low OFF noise. Functionally, the tight control of homeobox TF expression levels may reflect their known importance as determinants of cell identity, and that establishing and maintaining robust differences between cell types may require tight ON/OFF regulation rather than graded regulation.

Homeobox containing TFs can be subdivided into subfamilies based on their structure. The different homeobox subfamilies differed in their OFF noise and hence in their FCR values. Some
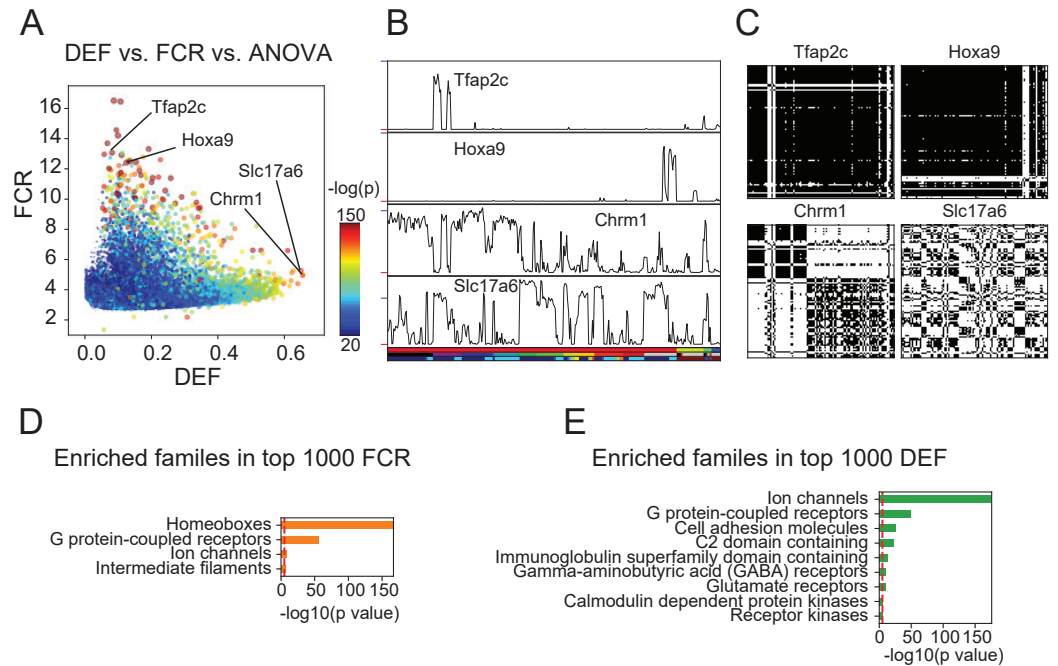
**Figure 4. DEF and FCR capture distinct aspects of expression diversity related to information content and robustness (A)** Highly variable genes (warm colored dots; color scale shows significance of ANOVA across cell populations) include both genes with high FCR and low DEF (like Tfap2c and Hoxa9) and genes with lower FCR and high DEF (like Chrm1 and Slc17a6). **(B)** Expression profiles of example genes labeled in A. Sample key in horizontal color bar as in Figure 1 Supplements 1-3. Red ticks at left indicates 0; Vertical scale is $log_2(FPKM + 1)$; blue ticks = 6) **(C)** DMs for example genes, calculated as shown in Figure 3. **(D),(E)** HUGO gene groups enriched in the top 1000 FCR and top 1000 DEF genes. Red lines indicate the p = $10^{-5}$ threshold used to judge significance.

214 families (e.g. HOXL, NKL, PRD) had very low OFF noise and high FCR, while others (e.g. CERS, PROS,
215 CUT) had higher OFF noise and lower FCR (Figure 5 Supplement 4).

216     The ability of gene families to provide information about cell identities reflects both how infor-
217 mative individual family members are, and the relationships between them. If the information
218 across family members is independent, the overall information is increased relative to the case in
219 which multiple members contain redundant information. This aspect of "family-wise" information
220 is not captured by "gene-wise" metrics like mean DEF, or by enrichment analysis (Figure 4D,E).
221 One means of capturing the additive, non-redundancy within a gene family is to measure the
222 orthogonality of expression patterns among the member genes. This analysis (Figure 5E) reveals
223 that homeobox TFs and GPCRs have the greatest orthogonality between cell types among HUGO
224 groups (as well as in PANTHER families, Figure 5 Supplement 1E). Related to this, we found that the
225 homeobox family can distinguish more than 99% of GACP pairs, suggesting these TFs comprise a
226 combinatorial code for the cell populations profiled. To illustrate this, we computed the minimum
227 set of homeobox TFs needed to distinguish the populations studied and found that a set of as few
228 as 8 could distinguish 99% of GACP pairs (Figure 5 Supplement 2B). Combinatorial codes could
229 also be produced from other highly orthogonal gene families, as illustrated for GPCRs Figure 5
230 Supplement 2C). As illustrated in these heat maps, expression differences for Homeobox TFs had
231 higher contrast, consistent with the fact that individually, homeobox TFs have the highest FCR
232 (Figure 4D) and lowest OFF noise (Figure 5B). In summary, we found that many homeobox genes are
233 expressed with a very high signal-to-noise ratio and are one of the groups of genes with the most
234 orthogonal expression patterns. This suggests that, similar to other tissues (*Kratsios et al., 2017*;
235 *Zheng et al., 2015*; *Dasen and Jessell, 2009*; *Philippidou and Dasen, 2013*), homeobox TFs play an
236 important role in specifying cell types in the brain.

### Diversity arising from alternative splicing

238 Alternative splicing is known to increase transcriptome diversity (*Andreadis et al., 1987*). To assess
239 the contribution of alternative splicing to diversifying transcriptomes across cell populations, we
240 quantified the branch probabilities at each alternative splice donor site within each gene (Figure 6A
241 top). The branch probabilities at each donor site are the relative frequencies with which particular
242 splice acceptors are chosen, and can be estimated from observed junction read counts. Branch
243 probabilities are highly bimodal (Figure 6A bottom), suggesting that most branch point choices are
244 made consistently, in an all-or-none fashion, for any given cell population.

245     To test the significance of differential splicing across cell populations, we utilized a statistical
246 test based on the Dirichlet-Multinomial distribution and the log-likelihood ratio test, developed in
247 LeafCutter (*Li et al., 2017*). We used pair-wise differential expression of each branch to calculate
248 a branch DEF, much as we previously calculated the differentially expressed fraction (DEF) from
249 expression values (Figure 3). Examples of branches with high DEFs are shown in Figure 6B. The list
250 includes known examples like the site of the flip and flop variants of the AMPA receptor subunit *Gria2*
251 (*Sommer et al., 1990*). Another previously known example is the splicing regulator muscleblind like
252 splicing factor 2 (*Mbnl2*), which is known to regulate splicing in the developing brain (*Charizanis
253 et al., 2012*) and is known to be spliced at multiple sites, including the one shown in Figure 6B
254 (*Pascual et al., 2006*).

255     In order to determine which families of genes are highly differentially spliced, we computed a
256 splice DEF per gene by combining the ability of a gene's alternatively spliced sites to distinguish a
257 pair of samples (i.e. a pair is distinguished by a gene if any alternatively spliced site in the gene can
258 distinguish the pair). Using this combined splice DEF, we found that RNA binding proteins, especially
259 splicing related factors (such as *Pcbp2* and *Mbnl2*) are highly alternatively spliced among neuronal
260 cell types (*Zheng and Black, 2013*), but over-represented categories also included other families
261 such as Glutamate receptors and G-protein modulators (Figure 6C).

262     To begin to assess the functional impact of alternative splicing, we determined which alternative
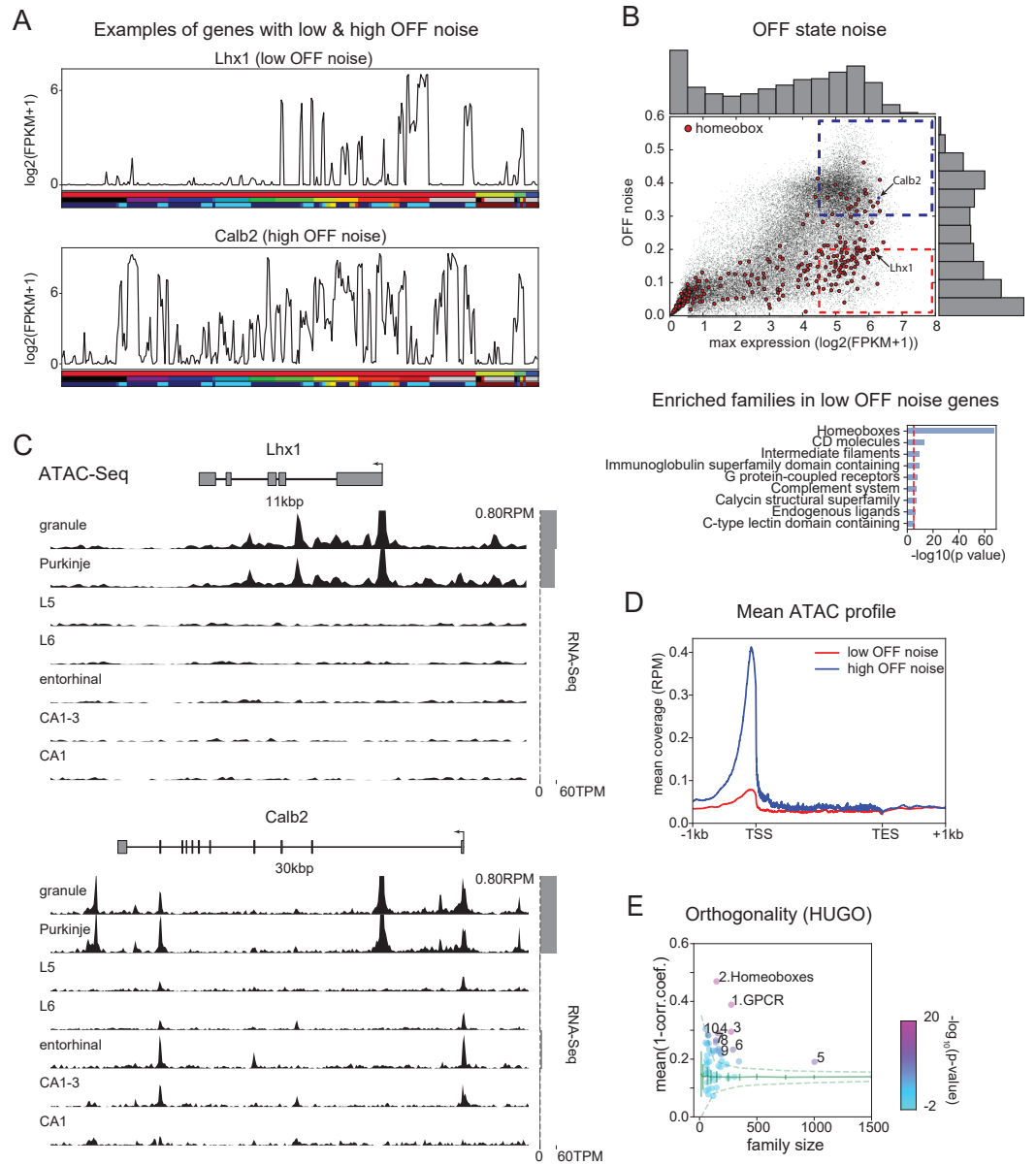263 sites lead to inclusion or exclusion of a known protein domain using the Pfam database (*Finn*

**Figure 5. Mechanisms contributing to low noise and high information content of Homeobox TFs. (A)** Example expression patterns of a LIM class homeobox TF (Lhx1) and a calcium binding protein (Calb2) with similar overall expression levels. Sample key as in Figure 1 Supplements 1-3. **(B)** (upper) OFF state noise (defined as standard deviation (std) of samples with FPKM<1) plotted against maximum expression. (lower) HUGO gene groups enriched in the region indicated by red dashed lines in the upper panel (see Figure 5 Supplement 1 for PANTHER and Gene Ontology enrichments). **(C)** Average (replicate n=2) ATAC-seq profiles for the genes shown in A. Some peaks are truncated. Expression levels are plotted at right (grey bars). **(D)** Length-normalized ATAC profile for genes with high (> 0.3, blue dashed box in B, n=853) and low (< 0.2, red dashed box in B, n=1643) OFF state expression noise. **(E)** Each circle represents the orthogonality of expression patterns calculated using HUGO gene groups. Orthogonality is a measure of the degree of non-redundancy in a set of expression patterns. Since the dispersion of orthogonality depends on family size, results are compared to orthogonality calculated from randomly sampled groups of genes (green solid lines: mean and std. dev.; green dashed lines: 99% confidence interval). Families, Z-scores, family size: 1. GPCR: 17.1, n=277; 2. Homeoboxes: 16.6, n=148; 3. Ion channels: 10.7, n=275; 4. C2 domain containing: 7.8, n=159; 5. Zinc fingers: 6.9, n=1002; 6. Immunoglobulin superfamily domain containing: 6.7, n=292; 7. PDZ domain containing: 6.3, n=144; 8. Fibronectin type III domain containing: 5.9, n=143; 9. Endogenous ligands: 5.1, n=165; 10. Basic helix-loop-helix proteins: 4.9, n=77

*et al., 2015*). In addition to providing information relevant to the potential functions of many previously unknown isoforms, our analysis also provides a more comprehensive view of known splice events. Two examples are shown in Figure 6D. Alternative splicing of Amyloid precursor-like protein 2 (*Aplp2*) is known to regulate inclusion of a bovine pancreatic trypsin inhibitor (BPTI) Kunitz domain (*Sandbrink et al., 1997*) and this domain is known to regulate proteolysis of the related protein APP, the amyloid precursor protein implicated in Alzheimer's disease (*Beckmann et al., 2016*). Differential inclusion of this exon is known to occur between neurons and nonneurons. Intriguingly, we found that splicing at this site in hippocampal interneurons differs not only from that in forebrain excitatory neurons, but also from other forebrain inhibitory neurons in neocortex and striatum. Kalirin (*Kalrn*) is a RhoGEF kinase implicated in Huntington's disease, schizophrenia and synaptic plasticity (*Penzes and Jones, 2008*). Kalrn is known to be regulated via binding of adaptor proteins to its SH3 (SRC homology 3) domains (*Schiller et al., 2006*) which is regulated by alternative splicing of this domain. In addition to expanding the number of known variants (blue exons and junctions in Figure 6D) we reveal their detailed distribution across the profiled set of neural populations. In total, the data reveal a detailed quantitative view of hundreds of thousands of known and unknown cell type-specific splicing events, providing an unmatched resource for investigating their functional significance.

Not all splicing events alter the inclusion or exclusion of known protein domains. Many splicing events introduce frame shifts or new stop codons and hence are predicted to lead to nonsense-mediated decay (NMD). Coupling of regulated splicing to NMD is believed to be an important mechanism for regulating protein abundance (*Lewis et al., 2002*). Consistent with previous observations (*Yan et al., 2015*; *Mauger and Scheiffele, 2017*), we noticed that most alternative sites contain branches that can lead to NMD (Figure 6E). This suggests that alternative splicing may contribute not only to the diversity of isoforms present, but to diversity defined on the basis of transcript abundance.

The present results provide a comprehensive resource of known and novel splicing events across a large number of neuronal cell types. Altogether, nearly 70% of alternative sites lead to differential inclusion of a known Pfam domain or NMD (Figure 6E), and thus to functional or quantitative diversity across cell types.

## Long genes contribute disproportionately to neuronal diversity

We found that neuronal effector genes (ion channels, receptors and cell adhesion molecules, etc.) have the greatest ability to distinguish cell populations (Figure 4E). Previously, these categories of genes have been found to be selectively enriched in neurons and to share the physical characteristic of being long (*Sugino et al., 2014*; *Gabel et al., 2015*; *Zylka et al., 2015*). Consistent with this, DEF, which approximates the mutual information (MI) between expression levels and cell populations, is significantly correlated with length (Figure 7A; correlation coefficient=0.19; p=7.5e-189), reaching a maximum for the very longest genes. Long genes (≥100kb) have nearly twice the average ability to distinguish cell populations (DEF) as shorter genes (Figure 7A), and provide greater family-wise separation between cell types (Figure 7C). Analyzing publicly available single cell data confirms that this bias is broadly observable (Figure 7 Supplement 1). In contrast, FCR, which measures the signal-to-noise or robustness of expression differences, is higher for shorter genes, reaching a maximum for genes below 10 kbp in length (Figure 7B).

Recently, (*Raman et al. 2018*) have argued that many prior observations of long gene bias are not significant when controlling for baseline variability in length-dependent expression. In order to assess the applicability of this argument to the present observations, we compared the fold-changes across length between groups and within replicates of individual groups as in (*Raman et al., 2018*). An example of this test applied to two populations is shown in Figure 7-Supplement 2A,B. Even after applying corrections for multiple comparisons across all bins, the long gene bins (>100 kbp) are highly significant. Panels C,D of this figure illustrate the results of performing this comparison for all GACPs in our dataset. The median fraction of significant long gene bins (0.89) greatly exceeded
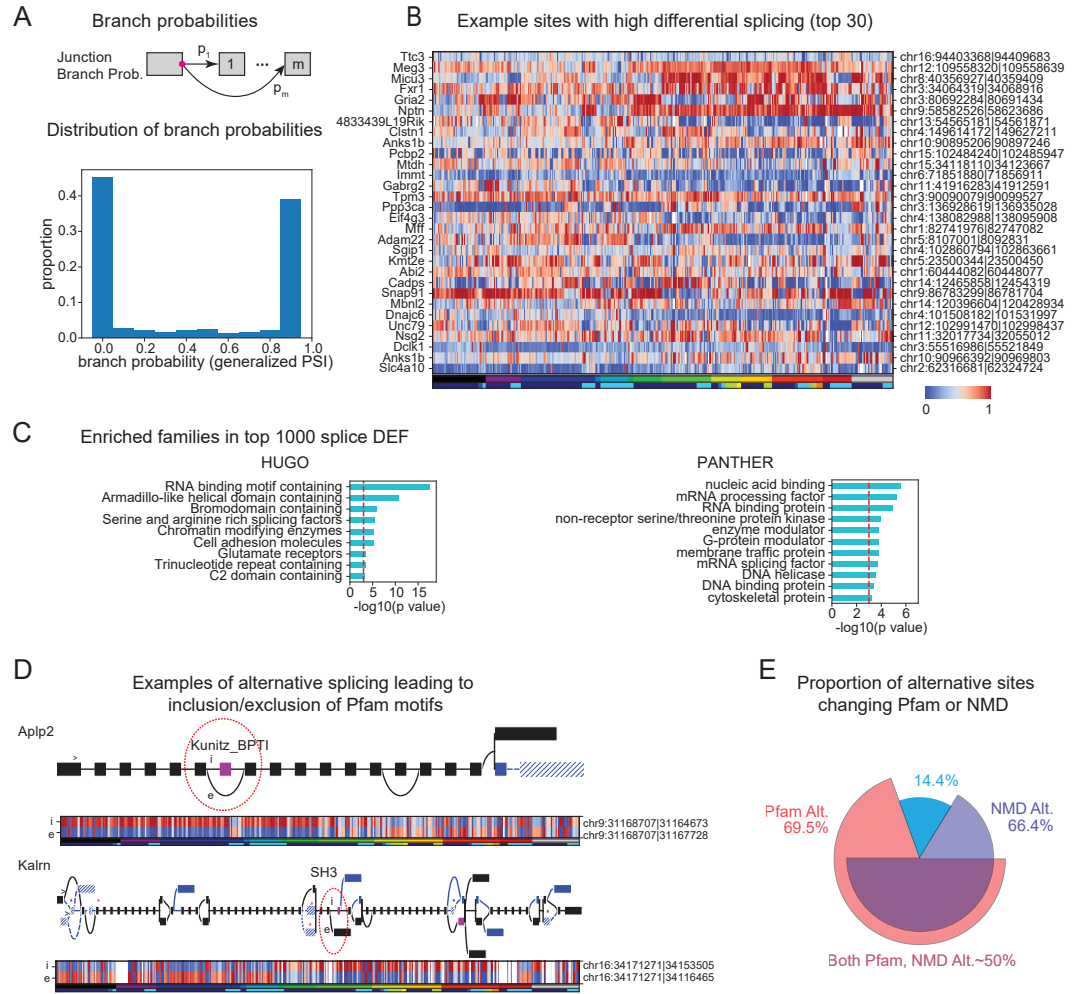
**Figure 6. Alternative splicing and neuronal diversity. (A)** (Top) Schematic representation of branch probabilities. Alternative donor sites (red dot) can be spliced to multiple acceptor sites $1, \dots, m$ with probabilities $p_1, \dots, p_m$. (Bottom) Distribution of branch probabilities across all samples and all alternative splice sites. **(B)** Heatmap showing branch probabilities across neuronal samples for branches with highest splice DEF. Each row corresponds to a branch within the indicated gene on the left and the location is indicated on the right. Samples without junctional reads at this branch are colored white. **(C)** Enriched HUGO gene groups and PANTHER protein classes for genes with top 1000 combined splice DEF. **(D)** Splice graphs illustrating examples of alternative splicing leading to inclusion or exclusion (marked "i","e") of Pfam domains (magenta exons) with branch probabilities shown in the heatmap below. Previously unannotated exons and junctions are blue; annotated are black. Dotted lines indicate branches predicted to lead to nonsense-mediated decay (NMD). A red star above an exon indicates existence of a premature termination codon (PTC) within the exon which satisfies the "50nt rule" for NMD (*Nagy and Maquat, 1998*) (i.e. more than 50bp upstream to the next junction), whereas a black star indicates existence of a PTC within 50bp of the next junction. Dashed lines and hatches indicate that there is no coding path through the element. (>) indicates an annotated translation start site. **(E)** Proportion of branch points predicted to lead to NMD (purple), altered Pfam inclusion (red), or both (overlapped region), at one or more of its branches.
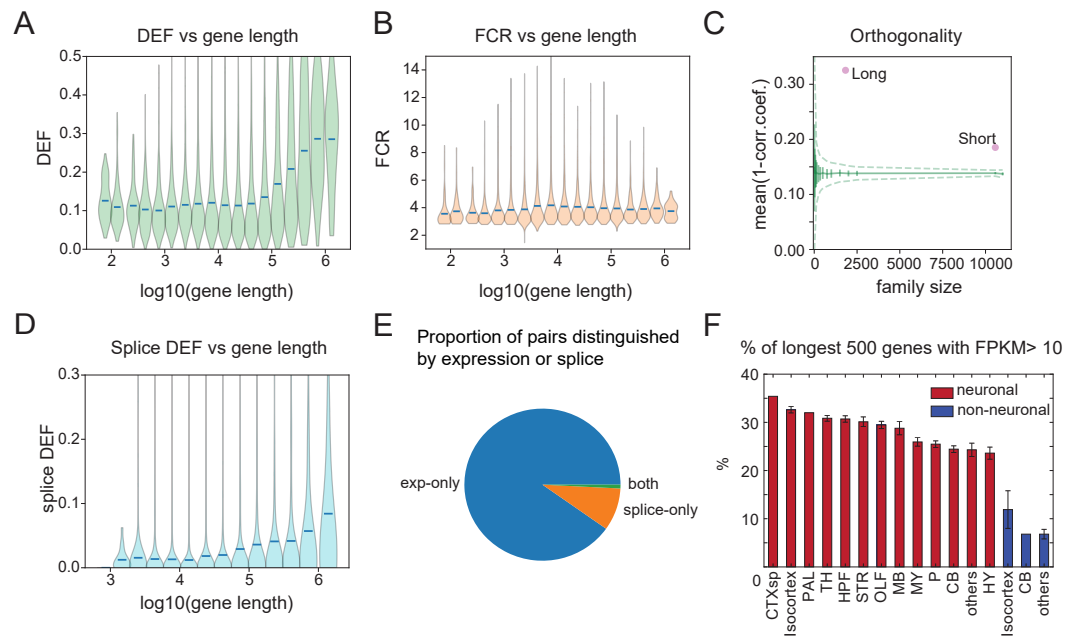
**Figure 7. Long genes have a greater capacity for distinguishing cell populations. (A)** DEF as a function of gene length. For violin plots in A, B, D, genes are sorted by length and binned (4 bins per log unit). **(B)** Robustness of expression difference (FCR) as a function of gene length. **(C)** Orthogonality of cell types calculated as in Figure 5E, but using long neuronal genes (n=1829, ≥100kb) and short neuronal genes (n=10572, <100kb) rather than functionally defined gene families. Z-score is 33.2 for long and 22.1 for short neuronal genes. Both are highly different from randomly sampled genes (green solid lines mean and std; dashed lines = 99% confidence interval), but long genes provide greater separation. **(D)** Splice DEF as a function of gene length. **(E)** Fraction of pairs distinguished by splicing (splice-only), transcript abundance (exp-only), or by both measures. **(F)** Variation in long gene expression in neuronal and nonneuronal populations across major brain regions studied. Error bars are SEM.

the fraction of short gene bins (0.1). A more detailed analysis of the test developed by Raman et al. and its application to other observations will be published elsewhere.

In addition to being differentially expressed, long genes are likely to have a larger number of exons and hence a greater potential for differential splicing. To evaluate the degree to which differential splicing of long genes contributes to distinguishing cell populations we plotted the splice DEF (Figure 6) as a function of gene length. As expected, DEF calculated from differential splicing also increased with gene length (Figure 7D) although the slope was more gradual and the maximum DEF value achieved was less than that for gene expression (Figure 7A). For each gene, we measured the fraction of cell populations pairs that could be distinguished on the basis of differential expression, differential splicing, or both. This revealed that for the current dataset, the average alternatively spliced gene distinguishes only 1.4 % of cell populations, but distinctions based on expression of these same genes were nearly ten times more common (13.9 %, Figure 7E).

Finally, to determine whether neuronal long gene expression contributes more to profiles in some anatomical regions than in others, we plotted the fraction of the longest genes expressed in neuronal and nonneuronal populations across each of the major brain regions studied. The results confirm strong differences between neurons and nonneurons and show the strongest long gene expression in forebrain regions, with weaker expression evident in hindbrain (Figure 7F). Analyses of single cell datasets revealed similar trends (Figure 7 Supplement 3).

## Discussion

### A resource of genetically identified neuronal transcriptomes

The dataset presented here is the largest collection of transcriptomes of anatomically and genetically specified neuronal cell types available in a mammalian species (Table 1). The approach employed in this study provides a complementary view of neuronal diversity to that afforded by SC sequencing. By sorting and pooling ~100 cells chosen based on genetic and anatomical similarity, we generated profiles with low noise and high depth, but, where tested, with a comparable degree of homogeneity, as that obtained in recent SC studies.

The fact that each transcriptome corresponds to a genetically (or retrogradely) labeled population will foster reproducible studies across investigators. The few profiles in our study that mapped to more than one SC profile (Figure 2), may represent cell types better distinguishable using SCs or improved genetic markers, or alternatively, may represent cell populations that are highly overlapping. The optimal granularity with which cell types may be distinguished remains an open question. Pooling cell profiles either prior to sequencing, as in this study, or after sequencing at the clustering phase, as in SC studies, risks compromising profile homogeneity. However, over-fragmenting clusters risks the opposite problem of reducing the reliability and reproducibility with which populations can be distinguished across studies. Given the complementary advantages of improved reproducibility and separability afforded by pooling profiles, and of reduced heterogeneity afforded by maximally separating profiles, further integration of these approaches with other modalities, such as FISH (*Moffitt et al., 2016*) are needed to accurately profile the full census of brain cell types. By linking these efforts to genetically identified neurons, the present dataset provides a useful resource for these efforts.

### A transcriptional code for neuronal diversity

We utilized easily calculated metrics that capture essential features of the robustness and information content of transcriptome diversity. These measures are simply versions of Fold-Change (FCR) and Differential Expression (DEF) adapted to the analysis of many separate populations simultaneously. Importantly, they capture independent components of the differences captured by variance-based metrics like ANOVA and CV (Figure 4A, Figure 3 Supplement 1). Metrics like ANOVA are influenced jointly by signal-to-noise and mutual information, while FCR and DEF better separate them (Figure 3 Supplement 1) and so these metrics may be more broadly useful when making genome-wide comparisons across many populations. In the present dataset, FCR and DEF identified two very different sets of genes contributing to neuronal diversity: high FCR, low-noise genes, exemplified by homeobox transcription factors, and high DEF, long neuronal effector genes like ion channels, receptors and cell adhesion molecules.

The homeobox family of TFs exhibited the most robust (high FCR) expression differences across cell types (Figure 4D). These ON/OFF differences were characterized by extremely low expression in the OFF state (Figure 5). Mechanistically, the low expression was associated with reduced genome accessibility measured by ATAC-seq (Figure 5C,D), presumably reflecting epigenetic regulation of the OFF state, known to occur for example at the clustered Hox genes via Polycomb group (PcG) proteins (*Montavon and Soshnikova, 2014*). Although this regulation has been studied most extensively at Hox genes, genome-wide ChIP studies reveal that PcG proteins are bound to over 100 homeobox TFs in ES cells (*Boyer et al., 2006*). Our results indicate that strong cell type-specific repression persists in the adult brain, presumably due to the continued functional importance of preventing even partial activation of inappropriate programs of neuronal identity.

Although individually, homeobox TFs contain less information about cell types than long neuronal effector genes, their patterns of expression are highly orthogonal and therefore their joint expression pattern is highly informative. As a group, homeobox TFs distinguished more than 99% of neuronal cell types profiled (Figure 5 Supplement 2). (Note this includes several Purkinje and Hippocampal pyramidal cell groups that may actually represent duplicate examples of the same

381 cell types). Historically, homeobox TFs are well known to combinatorially regulate neuronal identity
382 in *Drosophila* and *C. elegans* (***Kratsios et al., 2017***) and the vertebrate brainstem and spinal cord
383 (***Dasen and Jessell, 2009***; ***Philippidou and Dasen, 2013***). Our results suggest a broader importance
384 of homeobox TFs throughout the mammalian nervous system. Continued expression of these
385 factors in adult neurons suggests they likely also contribute to the maintenance of neuronal identity.

386 **Long genes and neuronal diversity**

387 Our study suggests that long neuronal effector genes contribute disproportionately to neuronal
388 transcriptional diversity (Figure 7). Previously, it was reported that differences in transcript length
389 can bias differential expression analysis of RNA-seq data (***Oshlack and Wakefield, 2009***). To ensure
390 that we avoided this bias, we used counts of reads only from within the 1 kbp-long 3′ ends of
391 the genes for calculating expression values. Recently, an alternative statistical analysis has been
392 used to argue that some of these length biases may be artefactual (***Raman et al., 2018***). Despite
393 concerns about the rigor of this analysis (manuscript in preparation), we found that the observed
394 length biases remain highly significant, even within this statistical framework (Figure 7 Supplement
395 2), suggesting that they are robust features of the transcriptional differences between neuronal
396 populations.

397 Long genes are expressed at higher levels in neurons than in nonneuronal cells in the nervous
398 system, a bias that was also present in SC datasets (Figure 7 Supplement 1,2) and that has been
399 reported previously (***Sugino et al., 2014***; ***Gabel et al., 2015***; ***Zylka et al., 2015***). These differences are
400 greatest in the forebrain (Figure 7F; Figure 7 Supplement 2), perhaps reflecting the large numbers
401 of distinct cell types in these regions and the enhanced ability of these genes to distinguish GACPs
402 based on their expression. However, we and others did not measure cell type-specific protein
403 expression, and so cannot be sure that the long gene bias extends to the level of neuronal proteins.

404 Long genes tend to have larger numbers of exons and therefore are likely to be expressed
405 in a larger number of distinct isoforms as a result of alternative splicing (alternative start sites
406 also contribute). We quantified differential splicing from analysis of junctional reads. Interestingly,
407 branch probabilities at most sites of alternative splicing were highly bimodal (Figure 6A), suggesting
408 that within each GACP, splicing is largely all or none, a finding previously reported in single immune
409 cells (***Shalek et al., 2013***) but not found in some single neuron studies (***Gokce et al., 2016***). This
410 led to patterns that often flipped between high and low probabilities for a given branch as one
411 traversed major brain region boundaries (Figure 6B). More than two thirds of these splicing events
412 lead to inclusion or exclusion of known protein domains (Figure 6E), but many of these, as well as
413 some of the remaining events that do not modify domain structure, also introduce a frame shift or
414 premature stop codon, and so are predicted to lead to nonsense mediated decay (NMD). We did not
415 directly test the contribution of NMD to transcript abundance, but our splicing results are consistent
416 with the idea that this may be an important mechanism for regulating transcript stability and hence
417 transcript abundance across different cell populations (***Yan et al., 2015***; ***Traunmuller et al., 2014***).
418 While differential splicing is able to distinguish fewer GACPs than transcript abundance (Figure 7E),
419 this may be an underestimate for two reasons. First, as just noted, splicing may influence transcript
420 abundance through NMD, and second, the sensitivity to detect splicing differences depends on an
421 adequate number of junctional reads. Deeper sequencing could increase the apparent contribution
422 of this component of neuronal diversity.

423 Long genes are enriched in the signaling molecules, receptors and ion channels responsible for
424 input/output transformations in neurons, and the cell adhesion molecules that specify neuronal
425 connectivity. The finding that these genes play an important role in diversifying cortical interneurons
426 (***Paul et al., 2017***), as well as distinguishing the larger set of populations studied here, is sensible
427 in light of the phenotypic diversity required for neuronal communication and connectivity. These
428 genes are long because of long introns that are rich in sequences derived from transposons and
429 other retroelements (***Grishkevich and Yanai, 2014***). Whether or how this increased length has
430 any functional significance for the regulation of these genes is unclear from our studies, but it is

intriguing that these long genes are disrupted in forms of autism spectrum disorder (*Zylka et al., 2015*; *Wei et al., 2016*) and in the related developmental disorder Rett Syndrome (*Sugino et al., 2014*; *Gabel et al., 2015*), where loss of the chromatin protein Mecp2 leads to selective upregulation of long neuronal genes in a highly cell type-specific fashion. These studies suggest the possibility that long neuronal genes are subject to distinct modes of regulation, with particular significance for neuronal diversity.

In contrast to long neuronal effector genes, which tend to be expressed later in development as neurons mature phenotypically (*Okaty et al., 2009*), low noise, high FCR genes are frequently critical for early development. These genes, such as many of the homeobox TFs, are often quite short and, at least in the case of the Hox genes, are known to be remarkably transposon impoverished (*Chinwalla et al., 2002*; *Simons, 2005*). This may reflect selection against transposon insertion, but may also reflect chromatin that is non-permissive for insertion in germ cells and the early embryo, where heritable transposition occurs. The high FCR/low OFF noise of many of these genes detected here may reflect a transcriptional signature of this class of genes. Consistent with this view, low OFF noise genes were nearly six times shorter than high OFF noise genes (Figure 5 Supplement 1D). Highly restrictive chromatin at these genes may be established early in development to protect them from disruptive transposition (*Montavon and Soshnikova, 2014*). If so, this tightly closed state is maintained in postmitotic neurons where it may also prevent transcriptional signals associated with inappropriate neural identities. This feature was not uniformly present across all subfamilies of homeobox transcription factors. Interestingly, however, the families with the highest FCR and lowest noise also had the shortest length, while those with higher noise expression (and lower FCR) were longer (Figure 5 Supplement 4).

The observation that long genes contribute disproportionately to neuronal transcriptional diversity is surprising both because of the increased metabolic cost of expressing them (*Castillo-Davis et al., 2002*), and since these genes are frequent sites of genome instability associated with genetic lesions leading to autism and other developmental disorders (*Wei et al., 2016*). These apparent disadvantages may be too weak to lead to selection against long gene expression in mammalian neurons. If this is not the case, however, it raises the question of why the mechanisms used to prevent elongation of shorter, low OFF noise genes were not also applied to neuronal effector genes. This could simply reflect developmental or later functional constraints that exclude the use of these epigenetic protection mechanisms. Alternatively, length itself may confer some advantages that outweigh other disadvantages. This could occur either through benefits provided by the diversification of alternative splicing, or through regulatory features contained within intronic sequences (*Zhao et al., 2018*).

## Acknowledgments

## Competing Interests

The authors declare no competing interests.

## Materials and Methods

### Cell types and mouse lines

We assume that cell types are organized hierarchically in a tree-like fashion proceeding from major branches (e.g. "cortical excitatory neuron") to more specialized subtypes, with the terminal "leaf-level" branches comprising "atomic" cell types. Profiled cell populations are defined operationally by the intersection of a transgenic mouse strain (or in some cases anatomical projection target) and a

brain region. Mouse lines profiled in this study are summarized in Supplementary File 1. Most were obtained from GENSAT (*Gong et al., 2007*) or from the Brandeis Enhancer Trap Collection (*Shima et al., 2016*). For Cre-driver lines, the Ai3, Ai9 or Ai14 reporter (*Madisen et al., 2009*) was crossed and offspring hemizygous for Cre and the reporter gene were used for profiling. Information on samples profiled is in Supplementary File 2. Populations profiled are designed to sample regions and cell types across the mouse brain within the limits of available resources. In addition several non-brain samples were profiled as out-groups. Replicate numbers (averaging 3 across all populations) are in Supplementary File 2. Replicates were obtained in single animals, except for a few cases in which pooling across animals was needed due to difficulty in sorting. Our study used a small number of replicates (n=2-4) per population to maximize the number of populations studied, while still allowing calculation of summary statistics. No explicit power analysis was performed. No attempt was made to remove outliers. Sequenced libraries were not used when total reads were low (<5M reads). Out of 179 neuronal GACPs, there are 165 groups which have more than one replicate. Of these, 14 were recent additions, and most analyses were performed with the remaining 151 groups. All experiments were conducted in accordance with the requirements of the Institutional Animal Care and Use Committees at Janelia Research Campus and Brandeis University.

**Tissue data**

In addition to cell type-specific data obtained in this study, we analyzed publicly available RNA-seq and DNase-seq data using tissue samples. Information on these samples are described in Supplementary File 3.

**Atlas**

Animals were anesthetized and perfused with 4% paraformaldehyde and brains were sectioned at $50\mu m$ thickness. Every fourth section was mounted on slides and imaged with a slide scanner equipped with a 20x objective lens (3DHISTECH; Budapest, Hungary). In house programs were used to adjust contrast and remove shading caused by uneven lighting. Images were converted to a zoomify-compatible format for web delivery and are available at http://neuroseq.janelia.org.

**Cell sorting**

Manual cell sorting was performed as described (*Hempel et al., 2007*; *Sugino et al., 2014*). Briefly, animals were sacrificed following isoflurane anesthesia, and $300\mu m$ slices were digested with pronase E (1mg/ml, P5147; Sigma-Aldrich) for 1 hour at room temperature, in artificial cerebrospinal fluid (ACSF) containing 6,7-dinitroquinoxaline-2,3-dione ($20\mu M$; Sigma-Aldrich), D-(–)-2-amino-5-phosphonovaleric acid ($50\mu M$; Sigma-Aldrich), and tetrodotoxin ($0.1\mu M$; Alomone Labs). Desired brain regions were micro-dissected and triturated with Pasteur pipettes of decreasing tip size. Dissociated cell suspensions were diluted 5-20 fold with filtered ACSF containing fetal bovine serum (1%; HyClone) and poured over Petri dishes coated with Sylgard (Dow Corning). For dim cells, Petri dishes with glass bottoms were used. Fluorescent cells were aspirated into a micropipette (tip diameter 30-50$\mu m$) under a fluorescent stereomicroscope (M165FC; Leica), and were washed 3 times by transferring to clean dishes. After the final wash, pure samples were aspirated in a small volume (1~3$\mu l$) and lysed in 47$\mu l$ XB lysis buffer (Picopure Kit, KIT0204; ThermoFisher) in a 200$\mu l$ PCR tube (Axygen), incubated for 30min at 40°C on a thermal cycler and then stored at -80°C. Detailed information on profiled samples are provided in Supplementary File 2.

**RNA-seq**

Total RNA was extracted using the Picopure kit (KIT0204; ThermoFisher). Either 1$\mu l$ total, or 1$\mu l$ per 50 sorted cells of $10^{-5}$ dilution of ERCC spike-in control (#4456740; Life Technologies) was added to the purified RNA and vacuum concentrated to 5$\mu l$ and immediately processed for reverse transcription using the NuGEN Ovation RNA-Seq System V2 (#7102; NuGEN) which yielded 4~8$\mu g$ of amplified DNA. Amplified DNA was fragmented (Covaris E220) to an average of ~200bp and ligated

525  to Illumina sequencing adaptors with the Encore Rapid Kit (0314; NuGEN). Libraries were quantified
526  with a KAPA Library Quant Kit (KAPA Biosystems) and sequenced on an Illumina HiSeq 2500 with 4
527  to 32-fold multiplexing (single end, usually 100bp read length, see Supplemental Table 2).

### RNA-seq analysis

529  Adaptor sequences (AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC for Illumina sequencing and
530  CTTTGTGTTTGA for NuGEN SPIA) were removed from de-multiplexed FASTQ data using cutadapt
531  v1.7.1 (http://dx.doi.org/10.14806/ej.17.1.200) with parameters "–overlap=7 –minimum-length=30".
532  Abundant sequences (ribosomal RNA, mitochondrial, Illumina phiX and low complexity sequences)
533  were detected using bowtie2 (*Langmead and Salzberg, 2012*) v2.1.0 with default parameters. The
534  remaining reads were mapped to the UCSC mm10 genome using STAR (*Dobin et al., 2012*) v2.4.0i
535  with parameters "–chimSegmentMin 15 –outFilterMismatchNmax 3". Mapped reads are quantified
536  with HTSeq (*Anders et al., 2015*) using Gencode.vM13 (*Harrow et al., 2012*).

### Annotations

538  For reference annotations we used Gencode.vM13 (*Harrow et al., 2012*) downloaded from
539  http://www.gencodegenes.org/, and NCBI RefSeq (*Pruitt et al., 2013*) downloaded from the UCSC
540  genome browser.

### Pan-neuronal genes

542  Pan-neuronal genes satisfied the following conditions: 1) mean neuronal expression level (NE)> 20
543  FPKM, 2) minimum NE > 5 FPKM, 3) mean NE > maximum nonneuronal expression level (NNE), 4)
544  minimum NE > mean NNE, 5) mean NE > 4x mean NNE, 6) mean NE > mean NNE + 2x standard
545  deviation of NNE, 7) mean NE − 2x standard deviation of NE > mean NNE.

### DEF/FCR/DM calculation

547  To calculate DEF, the following criteria were used to assign a "1" or "0" to each element in the
548  differentiation matrix (DM): log fold change > 2 and q-value <0.05. Q-values were calculated using
549  the limma package including the voom method (*Law et al., 2014*). To adjust the power to be similar
550  across cell types, two replicates (the most recent two) were used for all cell populations with more
551  than two replicates. We have tried the same calculations with 3 replicates (using a fewer number of
552  cell populations) and obtained similar results (data not shown). To avoid possible bias in variances
553  due to transcript length differences (*Oshlack and Wakefield, 2009*), we quantified counts using
554  reads from within the 3' 1 kbp of each gene. For genes with transcript lengths shorter than 1 kbp,
555  we used the whole gene length. We also calculated DEF and FCR across five SC datasets: For (*Zeisel
556  et al., 2015*), (*Tasic et al., 2016*) and (*Tasic et al., 2018*), we used log fold change > 1 and q-value
557  <0.05 calculated using the limma/voom method for differential gene expression. For (*Saunders
558  et al., 2018*) and (*Zeisel et al., 2018*), only cluster average expression was available, and log fold
559  change > 1 was defined as the criterion for differential expression.

### Overrepresentation, Orthogonality and Minimal gene sets

561  Overrepresentation analysis was performed using the top-level HUGO gene groups (Figures 4-6)
562  and was supplemented (Figure 6, Figure 4 Supplement 1, Figure 5 Supplements 1,3) using the
563  PANTHER Classification System and the Molecular Function component of the Gene Ontology
564  Annotation (GOM). Orthogonality quantifies the non-redundancy across expression patterns. We
565  calculated orthogonality (Figure 5E) as the mean pairwise decorrelation (1- Pearson's corr. coef.)
566  over a family of genes. Gene groups with less than 50 members were excluded, since variance of
567  this measure was much larger in small groups of randomly selected genes (dashed green lines in
568  Figure 5E). Minimal gene sets capable of serving as combinatorial codes across cell populations
569  (Figure 5 Supplement 2) were calculated by a greedy algorithm using the Differentiation Matrix (DM)
570  defined in Figure 3. Specifically, from a set of genes (such as homeobox TFs or other families), the

gene with the highest DEF was chosen as the first member of the set. Successive members were chosen, irrespective of their individual DEF, so as to maximize the combined DEF of the set. The combined DEF is the fraction of pairs distinguished by any gene in the group, and is calculated from the combined DM, which is the logical OR of the individual DMs for each gene in the group. This procedure continued until the combined DEF exceeded the desired threshold (0.99 in the case of Figure 5 Supplement 2). The homeoboxes set was constructed by merging the HUGO Homeoboxes gene group and the PANTHER homeobox protein TFs (PC00119) and had 156 genes. The GPCRs set is a merging of G protein-coupled receptors in HUGO and G-protein coupled receptors (PC00021) in PANTHER and has 347 genes.

## Calculation of differential splicing

To identify differential splicing, we utililzed a statistical test based on the Dirichlet-Multinomial distribution and the log-likelihood ratio test, developed in LeafCutter (*Li et al., 2017*). However, instead of using a group of connected introns as a unit for tests (as done in LeafCutter), we used a group of introns originating from an alternative donor site. Total junctional reads at an alternative donor > 10 was a prerequisite for testing. DM for alternative donors were then calculated as 1 for pairs of cell populations with $p < 0.05$ and maximum delta-PSI $> 0.1$, and 0 for others.

## NNLS/Random forest decomposition

The following single-cell datasets were downloaded and used for decomposition: (*Zeisel et al., 2015*) (NCBI GEO GSE60361), (*Tasic et al., 2016*) (NCBI GEO GSE71585), (*Tasic et al., 2018*) (http://celltypes.brain-map.org/rnaseq), (*Zeisel et al., 2018*) (http://mousebrain.org/), (*Saunders et al., 2018*) (dropviz.org). Deposited count data were converted to $log_2(CPM + 1)$ and used for comparison. The NeuroSeq dataset was quantified using RefSeq and featurecount (*Liao et al., 2013*) and converted into $log_2(CPM + 1)$. Subsets of genes common to NeuroSeq, Tasic 2018 and Zeisel 2018 datasets were used for decomposition. To account for differences in distributions of logCPM values between datasets, they were quantile-normalized to an average profile generated from the decomposed dataset. Since most genes in the single-cell profiles exhibited noisy expression patterns, using the entire gene set for decomposition was not feasible. Therefore, we selected genes deemed most informative for distinguishing cell classes based on the ANOVA F-statistic across cell classes (obtained using limma/voom in R). However, simply taking the top ANOVA genes led to highly biased gene selection since some cell types exhibited much larger transcriptional differences than others (e.g. many ANOVA selected genes were specific to microglia). We therefore selected genes to reduce the redundancy between distinguished cell populations. Beginning with the highest ANOVA gene (highest ANOVA F-value), genes were selected only if their DM (Differentiation Matrix defined in Figure 3) differed from those previously selected, enforced by requiring a Jaccard index threshold of 0.5, across all studies. We chose the top 500 genes meeting this criterion. Decompositions were performed on average profiles created by averaging NeuroSeq replicates or by averaging single-cell profiles using cluster assignments provided by the authors. NNLS was implemented using the R nnls library. For Random forest, the randomForest R package was used.

## ATAC-seq

7 cell types, Purkinje and granule cells from cerebellum, excitatory layer 5, 6 and entorhinal pyramidal cells from cortex, excitatory CA1, or CA1-3 pyramidal cells from hippocampus, labeled in mouse lines P036, P033, P078, 56L, P038, P064, and P036 respectively (all from *Shima et al., 2016*) were profiled with ATAC-seq. They were isolated by FACS to obtain ~40,000 labeled neurons. ATAC libraries for Illumina next-generation sequencing were prepared in accordance with a published protocol (*Buenrostro et al., 2013*). Briefly, collected cells were lysed in buffer containing 0.1% IGEPAL CA-630 (I8896, Sigma-Aldrich) and nuclei pelleted for resuspension in tagmentation DNA buffer with Tn5 (FC-121-1030, Illumina). Nuclei were incubated for 20-30 min at 37°C. Library amplification was monitored by real-time PCR and stopped prior to saturation (typically 8-10 cycles). Library

quality was assessed prior to sequencing using BioAnalyzer estimates of fragment size distributions looking for a ladder pattern indicative of fragmentation at nucleosome intervals as well as qPCR to determine relative enrichment at two housekeeping genes compared to background (specifically the TSS of *Gapdh* and *Actb* were assessed relative to the average of three intergenic regions). For sequencing, Illumina HiSeq 2500 with 2 to 4-fold multiplexing and paired end 100bp read length was used. In addition to ATAC-seq, RNA-seq was performed on replicate samples of ~2,000 cells collected in a similar way, and library prepared using the same method described above.

**ATAC-seq analysis**

Nextera adaptors (CTGTCTCTTATACACATCT) were trimmed from both ends from de-multiplexed FASTQ files using cutadapt with parameters "-n 3 -q 30,30 -m 36". Reads were then mapped to UCSC mm10 genome using bowtie2 (*Langmead and Salzberg, 2012*) with parameters "-X2000 –no-mixed –no-discordant". PCR duplicates were removed using Picard tools (http://broadinstitute.github.io/picard, v2.8.1) and reads mapping to mitochondrial DNA, scaffolds, and alternate loci were discarded. BigWig genomic coverage files were generated using bedtools (*Quinlan and Hall, 2010*) and scaled by the total number of reads per million.

**Anatomical region abbreviations**

Region abbreviations:

ACB: Nucleus accumbens

AD: Anterodorsal nucleus

AI: Agranular insular area

AMd: Anteromedial nucleus, dorsal part

AOBgr: Accessory olfactory bulb, granular layer

AOBmi: Accessory olfactory bulb, mitral layer

AP: Area postrema

ARH: Arcuate hypothalamic nucleus

AV: Anteroventral nucleus of thalamus

CA: Hippocampus Ammon's horn

CA1: Hippocampus field CA1

CA1sp: Hippocampus field CA1, pyramidal layer

CA3: Hippocampus field CA3

CEAm: Central amygdalar nucleus, medial part

CEAl: Central amygdalar nucleus, lateral part

CL: Central lateral nucleus of the thalamus

COAp: Cortical amygdalar area, posterior part

CP: Caudoputamen

CSm: Superior central nucleus raphe, medial part

CUL4,5gr: Cerebellum lobules IV-V, granular layer

CUL4,5mo: Cerebellum lobules IV-V, molecular layer

CUL4,5pu: Cerebellum lobules IV-V, Purkinje layer

DCO: Dorsal cochlear nucleus

DG: Hippocampus dentate gyrus

DMHp: Dorsomedial nucleus of the hypothalamus, posterior part

DMX: Dorsal motor nucleus of the vagus nerve

DR: Dorsal nucleus raphe

ECT: Ectorhinal area

IC: Inferior colliculus

IG: Induseum griseum

IO: Inferior olivary complex

isl: Islands of Calleja

668 islm: Major island of Calleja

669 LC: Locus ceruleus

670 LGd: Dorsal part of the lateral geniculate complex

671 LHA: Lateral hypothalamic area

672 MM, Medial mammillary nucleus

673 MO: Somatomotor area

674 MOBgl: Main olfactory bulb, glomerular layer

675 MOBgr: Main olfactory bulb, granular layer

676 MOBmi: Main olfactory bulb, mitral layer

677 MOE: main olfactory epithelium

678 MOp5: Primary motor area, layer 5

679 MV: Medial vestibular nucleus

680 NTS: Nucleus of the solitary tract

681 NTSge: Nucleus of the solitary tract, gelatinous part

682 NTSm: Nucleus of the solitary tract, medial part

683 ORBm: Orbital area, medial part

684 OT: Olfactory tubercle

685 PAG: Periaqueductal gray

686 PBl: Parabrachial nucleus, lateral division

687 PCN: Paracentral nucleus

688 PG: Pontine gray

689 PIR: Piriform area

690 PRP: Nucleus prepositus

691 PVH, Paraventricular hypothalamic nucleus

692 PVHd: Paraventricular hypothalamic nucleus, descending division

693 PVHp, Paraventricular hypothalamic nucleus, parvicellular division

694 PVT: Paraventricular nucleus of the thalamus

695 PYRpu: Cerebellum Pyramus (VIII), Purkinje layer

696 RPA: Nucleus raphe pallidus

697 RSPv: Retrosplenial area, ventral part

698 RT, Reticular nucleus of the thalamus

699 SCH: Suprachiasmatic nucleus

700 SCm: Superior colliculus, motor related

701 SFO: Subfornical organ

702 SNc: Substantia nigra, compact part

703 SO: Supraoptic nucleus

704 SSp: Primary somatosensory area

705 SSs: Supplemental somatosensory area

706 SUBd-sp: Subiculum, dorsal part, pyramidal layer

707 VII: Facial motor nucleus

708 VISp: Primary visual area

709 VISp6a: Primary visual area, layer 6a

710 VNO: vemoronasal organ

711 VPM: Ventral posteromedial nucleus of the thalamus

712 VTA: Ventral tegmental area

713

714 ## Appendix 1

715 **Relationship between DEF and Gini-Simpson index or MI** Here we explore in more detail the

716 relationship between DEF (differentially expressed fraction of populations) and Gini-Simpson index

717 (GSI) or MI (mutual information). DEF of a gene is equivalent to the Gini-Simpson index calculated

using distinguishable levels of expression of the gene and it is also closely related to mutual information between (discretized) expression levels and cell population labels.

Assume there are $N_e$ distinguishable expression levels of a gene and there are $n_i$ cell population groups in level $i$. Then, the Gini-Simpson index (GSI) is:

$$GSI = 1 - \sum_{i=1}^{N_e} p_i^2 \tag{1}$$

$$= 1 - \frac{\sum_{i=1}^{N_e} n_i(n_i - 1)}{N(N - 1)} \tag{2}$$

Where $p_i$ is the probability of randomly selected element being in expression level $i$ and $N = \sum_{i=1}^{N_e} n_i$ is the total number of groups. The second equation holds since $p_i^2 = n_i(n_i - 1)/N(N - 1)$ for sampling without replacement.

Since $n_i(n_i - 1)/N(N - 1) = (n_i(n_i - 1)/2)/(N(N - 1)/2)$, this term is the fraction of pairs in level $i$. So the sum of these are the total fraction of indistinguishable pairs and one minus this sum equals the fraction of distinguishable pairs, which is DEF. Thus, DEF is equivalent to the Gini-Simpson index calculated using distinguishable levels of expression.

To calculate mutual information between expression levels and cell populations, we discretize expression levels into $N_e$ levels. Let $N_s$ be the number of samples. Let $n_{ij}$ be counts in the contingency table where $i = 1, ..., N_e$ and $j = 1, ..., N_s$. Then the joint probability distribution and the marginal probability distribution can be written as:

$$p(i, j) = \frac{n_{ij}}{N_s} \tag{3}$$

$$p(i) = \frac{\sum_j n_{ij}}{N_s} = \frac{n_i}{N_s} \tag{4}$$

$$p(j) = \frac{\sum_i n_{ij}}{N_s} = \frac{n_j}{N_s} \tag{5}$$

Where $n_i = \sum_j n_{ij}$ and $n_j = \sum_i n_{ij}$. $n_i$ is the number of samples in level $i$ and $n_j$ is the number of replicates in cell type $j$. The mutual information between expression level (E) and samples (S) is:

$$I(E; S) = \sum_{i,j} p(i, j) \log \frac{p(i, j)}{p(i)p(j)} \tag{6}$$

$$= \sum_{i,j} p(i, j) \log \frac{p(i, j)}{p(j)} - \sum_{i,j} p(i, j) \log p(i) \tag{7}$$

$$= \sum_{i,j} p(j)p(i|j) \log p(i|j) - \sum_{i,j} p(i, j) \log p(i) \tag{8}$$

$$= \sum_j p(j) \sum_i p(i|j) \log p(i|j) - \sum_i \log p(i) \sum_j p(i, j) \tag{9}$$

$$= -\sum_j p(j)H(E|S = j) - \sum_i p(i) \log p(i) \tag{10}$$

$$= -H(E|S) + H(E) \tag{11}$$

$H(E|S = j)$ is the entropy of expression levels in cell population j, which represents the expression noise in cell population j, and $H(E|S)$ is the average of these across all cell populations. When there are no replicates, $H(E|S)$ is zero. When there are replicates, $H(E|S = j)$ represents how noisy the expression is. This may depend on expression level, and $H(E|S)$, the average of $H(E|S = j)$ may depend on expression prevalence (i.e., how widely the gene is expressed), but in any case, the first term $-H(E|S)$ represents reduction of the mutual information by noise.

The second term $H(E)$ is the entropy of the marginal distribution $p(i)$ and represents the main information content about cell groups encoded in expression levels. This can be rewritten using counts in the contingency table as:

$$H(E) = -\sum_i p(i) \log p(i) \tag{12}$$

$$= -\sum_i \frac{n_i}{N_s} \log \frac{n_i}{N_s} \tag{13}$$

$$= -\sum_i \frac{n_i}{N_s} \log n_i + \sum_i \frac{n_i}{N_s} \log N_s \tag{14}$$

$$= -\frac{1}{N_s} \sum_i n_i \log n_i + \log N_s \tag{15}$$

Thus, it is maximized when all $n_i$'s are 0 or 1, which corresponds to the case in which one expression level corresponds to one cell population, making all cell populations distinguishable by the expression levels. This is true when the number of discretization levels exceeds the number of samples. When the number of discretization levels ($N_e$) is less than the number of samples ($N_s$), $H(E)$ takes the maximum value of $\log N_e$ when all the samples are distributed equally across each bin.

To explore the relationship between $H(E)$ and DEF, the $\log n_i$ in the first term is replaced (approximated) by $(n_i - 1)$ (first two terms in the Taylor expansion of $\log n_i$ around $n_i = 1$.):

$$H(E) \sim -\frac{1}{N_s} \sum_i n_i(n_i - 1) + \log N_s \tag{16}$$

$$= -\frac{2}{N_s} \sum_i n_i(n_i - 1)/2 + \log N_s \tag{17}$$

$$= \frac{2}{N_s} \left\{ N_s(N_s - 1)/2 - \sum_i n_i(n_i - 1)/2 \right\} - (N_s - 1) + \log N_s \tag{18}$$

$$= (N_s - 1)DEF - (N_s - 1) + \log N_s \tag{19}$$

Since $n_i$ is the number of samples in one expression level, $n_i(n_i - 1)/2$ is the number of indistinguishable pairs in that expression level when there are no replicates. The term within the curly bracket is then the number of distinguishable pairs, leading to eq.(19).

More formally, since both $h(p) = \sum n_i \log n_i$ and $d(p) = \sum n_i(n_i - 1) = \sum n_i^2 - N_s$ are Schur-convex functions[1] on partitions of $N_s$, $p = (n_1, n_2, ..., n_k)$, when partition $p_1$ majorizes $p_2$ then, $h(p_1) \geq h(p_2)$ and $d(p_1) \geq d(p_2)$. When the partition length is 2, that is, when expression levels are discretized into only 2 levels, corresponding to ON and OFF, then, all of the partitions can be ordered with respect to majorization, therefore, $h(p)$ and $d(p)$ are order-preserved transformations of each other (Figure 3 Supplement 1C left). When the partition length is greater than 2, this relationship is not satisfied. However, they are still highly correlated to each other (Figure 3 Supplement 1C right).

When DEF is calculated from global discretization (as in the above case), the maximum number of pairs distinguishable occurs when all samples are equally distributed across bins and the number of distinguishable pairs is $\left( \frac{N_s}{N_e} \right)^2 N_e(N_e - 1)/2$. Therefore,

$$max(DEF) = \left( \frac{N_s}{N_e} \right)^2 \frac{N_e(N_e - 1)/2}{N_s(N_s - 1)/2} \tag{20}$$

$$= \left( 1 - \frac{1}{N_e} \right) / \left( 1 - \frac{1}{N_s} \right) \tag{21}$$

$$\sim 1 - \frac{1}{N_e} \quad (when \quad N_s \gg 1) \tag{22}$$

---

[1]A Schur-convex function is a function $f : \mathbb{R}^k \to \mathbb{R}$ which satisfies $f(x) \geq f(y)$ for all $x, y$ where $x$ majorizes $y$. For $x = (x_1, x_2, ..., x_k) \in \mathbb{R}^k where(x_1 \geq x_2 \geq ... \geq x_k)$ and $y = (y_1, y_2, ..., y_k) \in \mathbb{R}^k where(y_1 \geq y_2 \geq ... \geq y_k)$. $x$ majorizes $y$ when $\sum_{i=1}^k x_i = \sum_{i=1}^k y_i and \sum_{i=1}^j x_i \geq \sum_{i=1}^j y_i for all j = 1, ..., k$. When $x$ majorizes $y$, it follows $x_i \geq y_i$ for all $i$, so it is easy to see $h(x) \geq h(y)$ and $d(x) \geq d(y)$.

As stated above, this is also when the entropy $H(E)$ takes the maximum value of $\log_2 N_e$ in the unit of bits. (Figure 3 Supplement 1C)

## References

**Anders S**, Pyl PT, Huber W. HTSeq–a Python framework to work with high-throughput sequencing data. Bioinformatics. 2015; 31(2):166–9. https://www.ncbi.nlm.nih.gov/pubmed/25260700, doi: 10.1093/bioinformatics/btu638.

**Andreadis A**, Gallego ME, Nadal-Ginard B. Generation of Protein Isoform Diversity by Alternative Splicing: Mechanistic and Biological Implications. Annual Review of Cell Biology. 1987 nov; 3(1):207–242. http://dx.doi.org/10.1146/annurev.cb.03.110187.001231, doi: 10.1146/annurev.cb.03.110187.001231.

**Arendt D**, Musser JM, Baker CVH, Bergman A, Cepko C, Erwin DH, Pavlicev M, Schlosser G, Widder S, Laubichler MD, Wagner GP. The origin and evolution of cell types. Nature Reviews Genetics. 2016 nov; 17(12):744–757. https://doi.org/10.1038%2Fnrg.2016.127, doi: 10.1038/nrg.2016.127.

**Ashburner M**, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM, Sherlock G. Gene Ontology: tool for the unification of biology. Nature Genetics. 2000 may; 25(1):25–29. https://doi.org/10.1038%2F75556, doi: 10.1038/75556.

**Beckmann AM**, Glebov K, Walter J, Merkel O, Mangold M, Schmidt F, Becker-Pauly C, Gütschow M, Stirnberg M. The intact Kunitz domain protects the amyloid precursor protein from being processed by matriptase-2. Biological Chemistry. 2016 jan; 397(8). http://dx.doi.org/10.1515/hsz-2015-0263, doi: 10.1515/hsz-2015-0263.

**Boyer LA**, Plath K, Zeitlinger J, Brambrink T, Medeiros LA, Lee TI, Levine SS, Wernig M, Tajonar A, Ray MK, Bell GW, Otte AP, Vidal M, Gifford DK, Young RA, Jaenisch R. Polycomb complexes repress developmental regulators in murine embryonic stem cells. Nature. 2006 apr; 441(7091):349–353. https://doi.org/10.1038%2Fnature04733, doi: 10.1038/nature04733.

**Braschi B**, Denny P, Gray K, Jones T, Seal R, Tweedie S, Yates B, Bruford E. Genenames.org: the HGNC and VGNC resources in 2019. Nucleic Acids Research. 2018 oct; https://doi.org/10.1093%2Fnar%2Fgky930, doi: 10.1093/nar/gky930.

**Buenrostro JD**, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin DNA-binding proteins and nucleosome position. Nature Methods. 2013 oct; 10(12):1213–1218. https://doi.org/10.1038%2Fnmeth.2688, doi: 10.1038/nmeth.2688.

**Ramon y Cajal S**. La fine structure des centres nerveux. The croonian lecture. Proc R Soc Lond B Biol Sci. 1894; 55:443–468.

**Castillo-Davis CI**, Mekhedov SL, Hartl DL, Koonin EV, Kondrashov FA. Selection for short introns in highly expressed genes. Nature Genetics. 2002 jul; 31(4):415–418. https://doi.org/10.1038%2Fng940, doi: 10.1038/ng940.

**Charizanis K**, Lee KY, Batra R, Goodwin M, Zhang C, Yuan Y, Shiue L, Cline M, Scotti MM, Xia G, Kumar A, Ashizawa T, Clark HB, Kimura T, Takahashi MP, Fujimura H, Jinnai K, Yoshikawa H, Gomes-Pereira M, Gourdon G, et al. Muscleblind-like 2-Mediated Alternative Splicing in the Developing Brain and Dysregulation in Myotonic Dystrophy. Neuron. 2012 aug; 75(3):437–450. https://doi.org/10.1016%2Fj.neuron.2012.05.029, doi: 10.1016/j.neuron.2012.05.029.

**Chinwalla AT**, Cook LL, Delehaunty KD, Fewell GA, Fulton LA, Fulton RS, Graves TA, Hillier LW, Mardis ER, McPherson JD, Miner TL, Nash WE, Nelson JO, Nhan MN, Pepin KH, Pohl CS, Ponce TC, Schultz B, Thompson J, Trevaskis E, et al. Initial sequencing and comparative analysis of the mouse genome. Nature. 2002 dec; 420(6915):520–562. https://doi.org/10.1038%2Fnature01262, doi: 10.1038/nature01262.

**Dasen JS**, Jessell TM. Chapter Six Hox Networks and the Origins of Motor Neuron Diversity. In: *Current Topics in Developmental Biology* Elsevier; 2009.p. 169–200. https://doi.org/10.1016%2Fs0070-2153%2809%2988006-x, doi: 10.1016/s0070-2153(09)88006-x.

**Dobin A**, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. STAR: ultrafast universal RNA-seq aligner. Bioinformatics. 2012 oct; 29(1):15–21. https://doi.org/10.1093%2Fbioinformatics%2Fbts635, doi: 10.1093/bioinformatics/bts635.

**Doyle JP**, Dougherty JD, Heiman M, Schmidt EF, Stevens TR, Ma G, Bupp S, Shrestha P, Shah RD, Doughty ML, Gong S, Greengard P, Heintz N. Application of a Translational Profiling Approach for the Comparative Analysis of CNS Cell Types. Cell. 2008 nov; 135(4):749–762. https://doi.org/10.1016%2Fj.cell.2008.10.029, doi: 10.1016/j.cell.2008.10.029.

**Finn RD**, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC, Punta M, Qureshi M, Sangrador-Vegas A, Salazar GA, Tate J, Bateman A. The Pfam protein families database: towards a more sustainable future. Nucleic Acids Res. 2015 dec; 44(D1):D279–D285. http://dx.doi.org/10.1093/nar/gkv1344, doi: 10.1093/nar/gkv1344.

**Gabel HW**, Kinde B, Stroud H, Gilbert CS, Harmin DA, Kastan NR, Hemberg M, Ebert DH, Greenberg ME. Disruption of DNA-methylation-dependent long gene repression in Rett syndrome. Nature. 2015 mar; 522(7554):89–93. https://doi.org/10.1038%2Fnature14319, doi: 10.1038/nature14319.

**Gokce O**, Stanley GM, Treutlein B, Neff NF, Camp JG, Malenka RC, Rothwell PE, Fuccillo MV, Südhof TC, Quake SR. Cellular Taxonomy of the Mouse Striatum as Revealed by Single-Cell RNA-Seq. Cell Reports. 2016 jul; 16(4):1126–1137. https://doi.org/10.1016%2Fj.celrep.2016.06.059, doi: 10.1016/j.celrep.2016.06.059.

**Gong S**, Doughty M, Harbaugh CR, Cummins A, Hatten ME, Heintz N, Gerfen CR. Targeting Cre Recombinase to Specific Neuron Populations with Bacterial Artificial Chromosome Constructs. Journal of Neuroscience. 2007 sep; 27(37):9817–9823. https://doi.org/10.1523%2Fjneurosci.2707-07.2007, doi: 10.1523/jneurosci.2707-07.2007.

**Gong S**, Zheng C, Doughty ML, Losos K, Didkovsky N, Schambra UB, Nowak NJ, Joyner A, Leblanc G, Hatten ME, Heintz N. A gene expression atlas of the central nervous system based on bacterial artificial chromosomes. Nature. 2003 oct; 425(6961):917–925. https://doi.org/10.1038%2Fnature02033, doi: 10.1038/nature02033.

**Grishkevich V**, Yanai I. Gene length and expression level shape genomic novelties. Genome Research. 2014 jul; 24(9):1497–1503. https://doi.org/10.1101%2Fgr.169722.113, doi: 10.1101/gr.169722.113.

**Harrow J**, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D, Zadissa A, et al SS. GENCODE: The reference human genome annotation for The ENCODE Project. Genome Research. 2012 sep; 22(9):1760–1774. http://dx.doi.org/10.1101/gr.135350.111, doi: 10.1101/gr.135350.111.

**Hempel CM**, Sugino K, Nelson SB. A manual method for the purification of fluorescently labeled neurons from the mammalian brain. Nat Protoc. 2007 nov; 2(11):2924–2929. http://dx.doi.org/10.1038/nprot.2007.416, doi: 10.1038/nprot.2007.416.

**Kratsios P**, Kerk SY, Catela C, Liang J, Vidal B, Bayer EA, Feng W, Cruz EDDL, Croci L, Consalez GG, Mizumoto K, Hobert O. An intersectional gene regulatory strategy defines subclass diversity of C. elegans motor neurons. eLife. 2017 jul; 6. https://doi.org/10.7554%2Felife.25751, doi: 10.7554/elife.25751.

**Langmead B**, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nature Methods. 2012 mar; 9(4):357–359. http://dx.doi.org/10.1038/nmeth.1923, doi: 10.1038/nmeth.1923.

**Law CW**, Chen Y, Shi W, Smyth GK. voom: precision weights unlock linear model analysis tools for RNA-seq read counts. Genome Biology. 2014; 15(2):R29. https://doi.org/10.1186%2Fgb-2014-15-2-r29, doi: 10.1186/gb-2014-15-2-r29.

**Lein ES**, Hawrylycz MJ, Ao N, Ayres M, Bensinger A, Bernard A, Boe AF, Boguski MS, Brockway KS, Byrnes EJ, Chen L, Chen L, Chen TM, Chin MC, Chong J, Crook BE, Czaplinska A, Dang CN, Datta S, Dee NR, et al. Genome-wide atlas of gene expression in the adult mouse brain. Nature. 2006 dec; 445(7124):168–176. https://doi.org/10.1038%2Fnature05453, doi: 10.1038/nature05453.

**Lewis BP**, Green RE, Brenner SE. Evidence for the widespread coupling of alternative splicing and nonsense-mediated mRNA decay in humans. Proceedings of the National Academy of Sciences. 2002 dec; 100(1):189–192. https://doi.org/10.1073%2Fpnas.0136770100, doi: 10.1073/pnas.0136770100.

**Li YI**, Knowles DA, Humphrey J, Barbeira AN, Dickinson SP, Im HK, Pritchard JK. Annotation-free quantification of RNA splicing using LeafCutter. Nature Genetics. 2017 dec; http://dx.doi.org/10.1038/s41588-017-0004-9, doi: 10.1038/s41588-017-0004-9.

**Liao Y**, Smyth GK, Shi W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. Bioinformatics. 2013 nov; 30(7):923–930. https://doi.org/10.1093%2Fbioinformatics%2Fbtt656, doi: 10.1093/bioinformatics/btt656.

**Madisen L**, Zwingman TA, Sunkin SM, Oh SW, Zariwala HA, Gu H, Ng LL, Palmiter RD, Hawrylycz MJ, Jones AR, Lein ES, Zeng H. A robust and high-throughput Cre reporting and characterization system for the whole mouse brain. Nature Neuroscience. 2009 dec; 13(1):133–140. https://doi.org/10.1038%2Fnn.2467, doi: 10.1038/nn.2467.

**Mauger O**, Scheiffele P. Beyond proteome diversity: alternative splicing as a regulator of neuronal transcript dynamics. Current Opinion in Neurobiology. 2017 aug; 45:162–168. https://doi.org/10.1016%2Fj.conb.2017.05.012, doi: 10.1016/j.conb.2017.05.012.

**Mi H**, Huang X, Muruganujan A, Tang H, Mills C, Kang D, Thomas PD. PANTHER version 11: expanded annotation data from Gene Ontology and Reactome pathways and data analysis tool enhancements. Nucleic Acids Research. 2016 nov; 45(D1):D183–D189. https://doi.org/10.1093%2Fnar%2Fgkw1138, doi: 10.1093/nar/gkw1138.

**Mo A**, Mukamel EA, Davis FP, Luo C, Henry GL, Picard S, Urich MA, Nery JR, Sejnowski TJ, Lister R, Eddy SR, Ecker JR, Nathans J. Epigenomic Signatures of Neuronal Diversity in the Mammalian Brain. Neuron. 2015 jun; 86(6):1369–1384. https://doi.org/10.1016%2Fj.neuron.2015.05.018, doi: 10.1016/j.neuron.2015.05.018.

**Moffitt JR**, Hao J, Bambah-Mukku D, Lu T, Dulac C, Zhuang X. High-performance multiplexed fluorescence in situ hybridization in culture and tissue with matrix imprinting and clearing. Proceedings of the National Academy of Sciences. 2016 nov; 113(50):14456–14461. https://doi.org/10.1073%2Fpnas.1617699113, doi: 10.1073/pnas.1617699113.

**Montavon T**, Soshnikova N. Hox gene regulation and timing in embryogenesis. Seminars in Cell & Developmental Biology. 2014 oct; 34:76–84. https://doi.org/10.1016%2Fj.semcdb.2014.06.005, doi: 10.1016/j.semcdb.2014.06.005.

**Muotri AR**, Gage FH. Generation of neuronal variability and complexity. Nature. 2006 jun; 441(7097):1087–1093. https://doi.org/10.1038%2Fnature04959, doi: 10.1038/nature04959.

**Nagy E**, Maquat LE. A rule for termination-codon position within intron-containing genes: when nonsense affects RNA abundance. Trends in Biochemical Sciences. 1998 jun; 23(6):198–199. http://dx.doi.org/10.1016/s0968-0004(98)01208-0, doi: 10.1016/s0968-0004(98)01208-0.

**Okaty BW**, Miller MN, Sugino K, Hempel CM, Nelson SB. Transcriptional and Electrophysiological Maturation of Neocortical Fast-Spiking GABAergic Interneurons. Journal of Neuroscience. 2009 may; 29(21):7040–7052. https://doi.org/10.1523%2Fjneurosci.0105-09.2009, doi: 10.1523/jneurosci.0105-09.2009.

**Okaty BW**, Sugino K, Nelson SB. A Quantitative Comparison of Cell-Type-Specific Microarray Gene Expression Profiling Methods in the Mouse Brain. PLoS ONE. 2011 jan; 6(1):e16493. https://doi.org/10.1371%2Fjournal.pone.0016493, doi: 10.1371/journal.pone.0016493.

**Oshlack A**, Wakefield MJ. Transcript length bias in RNA-seq data confounds systems biology. Biology Direct. 2009; 4(1):14. https://doi.org/10.1186%2F1745-6150-4-14, doi: 10.1186/1745-6150-4-14.

**Pascual M**, Vicente M, Monferrer L, Artero R. The Muscleblind family of proteins: an emerging class of regulators of developmentally programmed alternative splicing. Differentiation. 2006 mar; 74(2-3):65–80. https://doi.org/10.1111%2Fj.1432-0436.2006.00060.x, doi: 10.1111/j.1432-0436.2006.00060.x.

**Paul A**, Crow M, Raudales R, He M, Gillis J, Huang ZJ. Transcriptional Architecture of Synaptic Communication Delineates GABAergic Neuron Identity. Cell. 2017 oct; 171(3):522–539.e20. https://doi.org/10.1016%2Fj.cell.2017.08.032, doi: 10.1016/j.cell.2017.08.032.

**Penzes P**, Jones KA. Dendritic spine dynamics – a key role for kalirin-7. Trends in Neurosciences. 2008 aug; 31(8):419–427. http://dx.doi.org/10.1016/j.tins.2008.06.001, doi: 10.1016/j.tins.2008.06.001.

**Philippidou P**, Dasen JS. Hox Genes: Choreographers in Neural Development Architects of Circuit Organization. Neuron. 2013 oct; 80(1):12–34. https://doi.org/10.1016%2Fj.neuron.2013.09.020, doi: 10.1016/j.neuron.2013.09.020.

**Pruitt KD**, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, Farrell CM, Hart J, Landrum MJ, McGarvey KM, Murphy MR, O'Leary NA, Pujar S, Rajput B, Rangwala SH, Riddick LD, Shkeda A, Sun H, Tamez P, Tully RE, et al. RefSeq: an update on mammalian reference sequences. Nucleic Acids Research. 2013 nov; 42(D1):D756–D763. http://dx.doi.org/10.1093/nar/gkt1114, doi: 10.1093/nar/gkt1114.

**Quinlan AR**, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics. 2010 jan; 26(6):841–842. http://dx.doi.org/10.1093/bioinformatics/btq033, doi: 10.1093/bioinformatics/btq033.

905 **Raman AT**, Pohodich AE, Wan YW, Yalamanchili HK, Lowry WE, Zoghbi HY, Liu Z. Apparent bias toward long gene
906 misregulation in MeCP2 syndromes disappears after controlling for baseline variations. Nature Communica-
907 tions. 2018; 9(1):3225. https://doi.org/10.1038%2Fs41467-018-05627-1, doi: 10.1038/s41467-018-05627-1.

908 **Sandbrink R**, Mönning U, Masters CL, Beyreuther K. Expression of the APP Gene Family in Brain Cells Brain
909 Development and Aging. Gerontology. 1997; 43(1-2):119–131. https://doi.org/10.1159%2F000213840, doi:
910 10.1159/000213840.

911 **Saunders A**, Macosko EZ, Wysoker A, Goldman M, Krienen FM, de Rivera H, Bien E, Baum M, Bortolin L, Wang S,
912 Goeva A, Nemesh J, Kamitaki N, Brumbaugh S, Kulp D, McCarroll SA. Molecular Diversity and Specializations
913 among the Cells of the Adult Mouse Brain. Cell. 2018; 174(4):1015–1030 e16. https://www.ncbi.nlm.nih.gov/
914 pubmed/30096299, doi: 10.1016/j.cell.2018.07.028.

915 **Schiller MR**, Chakrabarti K, King GF, Schiller NI, Eipper BA, Maciejewski MW. Regulation of RhoGEF Activity
916 by Intramolecular and Intermolecular SH3 Domain Interactions. Journal of Biological Chemistry. 2006 apr;
917 281(27):18774–18786. http://dx.doi.org/10.1074/jbc.m512482200, doi: 10.1074/jbc.m512482200.

918 **Shalek AK**, Satija R, Adiconis X, Gertner RS, Gaublomme JT, Raychowdhury R, Schwartz S, Yosef N, Malboeuf C, Lu
919 D, Trombetta JJ, Gennert D, Gnirke A, Goren A, Hacohen N, Levin JZ, Park H, Regev A. Single-cell transcriptomics
920 reveals bimodality in expression and splicing in immune cells. Nature. 2013 may; 498(7453):236–240.
921 https://doi.org/10.1038%2Fnature12172, doi: 10.1038/nature12172.

922 **Shima Y**, Sugino K, Hempel CM, Shima M, Taneja P, Bullis JB, Mehta S, Lois C, Nelson SB. A Mammalian
923 enhancer trap resource for discovering and manipulating neuronal cell types. eLife. 2016 mar; 5. https:
924 //doi.org/10.7554%2Felife.13503, doi: 10.7554/elife.13503.

925 **Simons C**. Transposon-free regions in mammalian genomes. Genome Research. 2005 dec; 16(2):164–172.
926 https://doi.org/10.1101%2Fgr.4624306, doi: 10.1101/gr.4624306.

927 **Simpson EH**. Measurement of Diversity. Nature. 1949 apr; 163(4148):688–688. https://doi.org/10.1038%
928 2F163688a0, doi: 10.1038/163688a0.

929 **Sommer B**, Keinanen K, Verdoorn T, Wisden W, Burnashev N, Herb A, Kohler M, Takagi T, Sakmann B, Seeburg P.
930 Flip and flop: a cell-specific functional switch in glutamate-operated channels of the CNS. Science. 1990 sep;
931 249(4976):1580–1585. https://doi.org/10.1126%2Fscience.1699275, doi: 10.1126/science.1699275.

932 **Stefanakis N**, Carrera I, Hobert O. Regulatory Logic of Pan-Neuronal Gene Expression in C. el-
933 egans. Neuron. 2015 aug; 87(4):733–750. https://doi.org/10.1016%2Fj.neuron.2015.07.031, doi:
934 10.1016/j.neuron.2015.07.031.

935 **Sugino K**, Hempel CM, Okaty BW, Arnson HA, Kato S, Dani VS, Nelson SB. Cell-Type-Specific Repression by Methyl-
936 CpG-Binding Protein 2 Is Biased toward Long Genes. Journal of Neuroscience. 2014 sep; 34(38):12877–12883.
937 https://doi.org/10.1523%2Fjneurosci.2674-14.2014, doi: 10.1523/jneurosci.2674-14.2014.

938 **Sugino K**, Hempel CM, Miller MN, Hattox AM, Shapiro P, Wu C, Huang ZJ, Nelson SB. Molecular taxonomy
939 of major neuronal classes in the adult mouse forebrain. Nature Neuroscience. 2005 dec; 9(1):99–107.
940 http://dx.doi.org/10.1038/nn1618, doi: 10.1038/nn1618.

941 **Taniguchi H**, He M, Wu P, Kim S, Paik R, Sugino K, Kvitsani D, Fu Y, Lu J, Lin Y, Miyoshi G, Shima Y, Fishell
942 G, Nelson SB, Huang ZJ. A Resource of Cre Driver Lines for Genetic Targeting of GABAergic Neurons in
943 Cerebral Cortex. Neuron. 2011 sep; 71(6):995–1013. https://doi.org/10.1016%2Fj.neuron.2011.07.026, doi:
944 10.1016/j.neuron.2011.07.026.

945 **Tasic B**, Menon V, Nguyen TN, Kim TK, Jarsky T, Yao Z, Levi B, Gray LT, Sorensen SA, Dolbeare T, Bertagnolli
946 D, Goldy J, Shapovalova N, Parry S, Lee C, Smith K, Bernard A, Madisen L, Sunkin SM, Hawrylycz M, et al.
947 Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. Nature Neuroscience. 2016 jan;
948 19(2):335–346. https://doi.org/10.1038%2Fnn.4216, doi: 10.1038/nn.4216.

949 **Tasic B**, Yao Z, Graybuck LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J, Garren E, Economo MN, Viswanathan
950 S, Penn O, Bakken T, Menon V, Miller J, Fong O, Hirokawa KE, Lathia K, Rimorin C, Tieu M, Larsen R, et al.
951 Shared and distinct transcriptomic cell types across neocortical areas. Nature. 2018 oct; 563(7729):72–78.
952 https://doi.org/10.1038%2Fs41586-018-0654-5, doi: 10.1038/s41586-018-0654-5.

953 **Traunmuller L**, Bornmann C, Scheiffele P. Alternative Splicing Coupled Nonsense-Mediated Decay Generates
954 Neuronal Cell Type-Specific Expression of SLM Proteins. Journal of Neuroscience. 2014 dec; 34(50):16755–
955 16761. https://doi.org/10.1523%2Fjneurosci.3395-14.2014, doi: 10.1523/jneurosci.3395-14.2014.

**Wei PC**, Chang AN, Kao J, Du Z, Meyers RM, Alt FW, Schwer B. Long Neural Genes Harbor Recurrent DNA Break Clusters in Neural Stem/Progenitor Cells. Cell. 2016 feb; 164(4):644–655. https://doi.org/10.1016%2Fj.cell.2015.12.039, doi: 10.1016/j.cell.2015.12.039.

**Yan Q**, Weyn-Vanhentenryck SM, Wu J, Sloan SA, Zhang Y, Chen K, Wu JQ, Barres BA, Zhang C. Systematic discovery of regulated and conserved alternative exons in the mammalian brain reveals NMD modulating chromatin regulators. Proceedings of the National Academy of Sciences. 2015 mar; 112(11):3445–3450. http://dx.doi.org/10.1073/pnas.1502849112, doi: 10.1073/pnas.1502849112.

**Zeisel A**, Hochgerner H, Lonnerberg P, Johnsson A, Memic F, van der Zwan J, Haring M, Braun E, Borm LE, La Manno G, Codeluppi S, Furlan A, Lee K, Skene N, Harris KD, Hjerling-Leffler J, Arenas E, Ernfors P, Marklund U, Linnarsson S. Molecular Architecture of the Mouse Nervous System. Cell. 2018; 174(4):999–1014 e22. https://www.ncbi.nlm.nih.gov/pubmed/30096314, doi: 10.1016/j.cell.2018.06.021.

**Zeisel A**, Munoz-Manchado AB, Codeluppi S, Lonnerberg P, Manno GL, Jureus A, Marques S, Munguba H, He L, Betsholtz C, Rolny C, Castelo-Branco G, Hjerling-Leffler J, Linnarsson S. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. Science. 2015 feb; 347(6226):1138–1142. https://doi.org/10.1126%2Fscience.aaa1934, doi: 10.1126/science.aaa1934.

**Zhang S**, Kanemitsu Y, Fujitani M, Yamashita T. The newly identified migration inhibitory protein regulates the radial migration in the developing neocortex. Scientific Reports. 2014 aug; 4(1). https://doi.org/10.1038%2Fsrep05984, doi: 10.1038/srep05984.

**Zhang Y**, Chen K, Sloan SA, Bennett ML, Scholze AR, O'Keeffe S, Phatnani HP, Guarnieri P, Caneda C, Ruderisch N, Deng S, Liddelow SA, Zhang C, Daneman R, Maniatis T, Barres BA, Wu JQ. An RNA-Sequencing Transcriptome and Splicing Database of Glia Neurons, and Vascular Cells of the Cerebral Cortex. Journal of Neuroscience. 2014 sep; 34(36):11929–11947. https://doi.org/10.1523%2Fjneurosci.1860-14.2014, doi: 10.1523/jneurosci.1860-14.2014.

**Zhao YT**, Kwon DY, Johnson BS, Fasolino M, Lamonica JM, Kim YJ, Zhao BS, He C, Vahedi G, Kim TH, Zhou Z. Long genes linked to autism spectrum disorders harbor broad enhancer-like chromatin domains. Genome Research. 2018; 28:933–942. https://doi.org/10.1101%2Fgr.233775.117, doi: 10.1101/gr.233775.117.

**Zheng C**, Diaz-Cuadros M, Chalfie M. Hox Genes Promote Neuronal Subtype Diversification through Posterior Induction in Caenorhabditis elegans. Neuron. 2015 nov; 88(3):514–527. https://doi.org/10.1016%2Fj.neuron.2015.09.049, doi: 10.1016/j.neuron.2015.09.049.

**Zheng S**, Black DL. Alternative pre-mRNA splicing in neurons: growing up and extending its reach. Trends in Genetics. 2013 aug; 29(8):442–448. https://doi.org/10.1016%2Fj.tig.2013.04.003, doi: 10.1016/j.tig.2013.04.003.

**Zylka MJ**, Simon JM, Philpot BD. Gene Length Matters in Neurons. Neuron. 2015 apr; 86(2):353–355. https://doi.org/10.1016%2Fj.neuron.2015.03.059, doi: 10.1016/j.neuron.2015.03.059.

## Supplementary Materials

**Supplementary Files**

Supplementary File 1

Table listing information for mouse lines. Information (columns) includes regions profiled, source of the mouse line, repository ID and URL, whether atlas is available via the Janelia viewer, URL for other atlases, and relevant references.

Supplementary File 2

Table for sample information. Included fields are,

1. sample_id: Sample ID;
2. sample_name: Sample Name;
3. group: Sample Group ID;
4. group_label: Label for Group;
5. sample_label: Label for Sample;
6. seqlane: Sequencing Lane ID;
7. mouseline: Mouse Line ID;

8. sample_code: Type of sample, cs.n: cell-type-specific neuronal sample; cs.o: cell-type-specific nonneuronal sample; ti.b: tissue sample from brain; ti.o: sample from non-brain tissue; cs.p: cell-type-specific progenitor sample;

9. region: Anatomical Region (large structure);

10. transmitter: Transmitter;

11. allenregion: Region using Allen Reference Atlas notation;

12. num_cells: Number of cells used in the sample;

13. age_(day): Postnatal age (in days) of the mouse;

14. sex: Sex of the mouse;

15. weight_(g): Weight (g) of the mouse;

16. ercc(10ᵉ5 dilution ul): Amount of added ERCC in ul. ($10^{-5}$ diluted);

17. ercc_mix: Which ERCC mix is used;

18. adaptor: Which Illumina (Solexa) sequencing adaptor is used;

19. total_reads: Total number of sequencing reads;

20. total_wo_ERCC: Total number of sequencing reads without reads mapping to ERCC;

21. read_length: Sequencing read length;

22. ercc%: Percentage of ERCC reads;

23. ribosomal_etc%: Percentage of reads mapping to ribosomal or other abundant sequences (phiX, polyC, polyA);

24. unmapped_reads%: Percentage of reads not mapped to mm10 genome;

25. unique_reads%: Percentage of reads uniquely mapped;

26. nonunique_reads%: Percentage of non-uniquely mapped reads;

27. short_insert%: Percentage of short (<30bp) reads;

28. mapped_reads: Number of mapped reads;

29. comments: Comments;


## Supplementary File 3

Table listing public tissue samples used in analyses.

**Figure 1–Supplement 1.**
**GACP samples.** Sample groups color coded by type (left color bar), region (middle color bar) and transmitter phenotype (right color bar). Transmitter phenotype was determined from transmitter synthesis and storage enzyme expression. Abbreviations: OLF: olfactory regions; CTXsp.CLA: Claustrum; HPF: hippocampal formation; STR: Striatum and related ventral forebrain structures; PAL: pallidum; TH: thalamus; HY: hypothalamus; MB: midbrain; MY: medulla; P: pons; CB: cerebellum; RE: retina; OE: olfactory epithelium; SP: spinal cord; @: other non-brain regions. For additional abbreviations see Methods.

**Figure 1–Supplement 2.**

**Quality control measures. (A)** (Top) Total reads for each of the libraries. Samples are color coded by type, region and transmitter, as shown in Figure 1 Supplement 1. (Bottom) Categories of reads in each library: unmapped: reads that did not map to the mm10 genome including chimeric and back-spliced reads; short: reads less than 30bp in length after removing adaptor sequences; non-unique: reads mapping to multiple locations; abundant: reads containing ribosomal RNA, polyA, polyC and phiX sequences, and unique: uniquely mapped reads. For further analyses, abundant, short and unmapped reads were not used. **(B)** Contaminating transcripts from nonneuronal cell populations. Samples with significant expression of these transcripts (at right) include tissue samples and nonneuronal samples. Each row is normalized by the maximum value.

**Figure 1–Supplement 3.**
**Pan-neuronal genes.** Genes expressed in all neuronal GACPs, but not (or at much lower levels) in nonneurons within the dataset. Heat-map shows log expression levels and the color at the right side indicates fold-change of the expression level between neurons and nonneurons. Criteria for extracting these genes are listed in the Methods.

A

Tasic 2018



B

Zeisel 2018



C

NeuroSeq



1035

D

Mean Purity Scores



**Figure 2–Supplement 1.**
**Self decompositions by NNLS**. Each dataset is randomly divided into two groups and one is used to decompose the other. Coefficients matrix with perfect decomposition would be diagonal. Non-diagonal elements indicate limitation of the decomposition method due to having a subset of cell groups too similar to each other. **(A-C)** Heatmaps illustrate NNLS coefficients for subsets of samples in each dataset. Column order is same as row order. **(A)** 25 neocortical samples from *Tasic et al.* (*2018*) **(B)** 25 neocortical samples from *Zeisel et al.* (*2018*) **(C)** 28 neocortical samples from present study. **(D)** Mean purity scores (as defined in Figure 2) for cross-validation (calculated over all neocortical samples) were comparable in each dataset. Error bars are Std. Dev.

**Figure 2–Supplement 2.**

**A validation of NNLS decomposition. (Left)** Single cell profiles from *Tasic et al.* (*2016*) were merged according to which of the 17 transgenic strains and sub-dissected layers they originated from (row labels). Merged profiles were then decomposed using individual cell type cluster profiles defined in *Tasic et al.* (*2016*) (column labels). **(Right)** The reported proportion of single cell profiles according to the author's classification. The close similarity between left and right matrices indicates an accurate NNLS decomposition of the merged clusters. Note that information about which and how many individual cell types were sorted from each line and set of layers was not explicitly provided to the decomposition algorithm, but were accurately deduced from the merged expression profiles.

**Figure 2–Supplement 3.**
**NNLS decomposition of SC datasets: Tasic by Zeisel**. The same neocortical samples from (*Tasic et al., 2018*; *Zeisel et al., 2018*) used in Figure 2 to decompose NeuroSeq neocortical samples were used to decompose each other. See Figure 2 for further details of cell identity. Order of samples listed is as in Figure 2. Presumably because Tasic et al. samples are more finely sub-clustered, individual Zeisel et al. samples (horizontal) frequently map to multiple Tasic samples (vertical).

**Figure 2–Supplement 4.**
**NNLS decomposition of SC datasets: Zeisel by Tasic**. The same neocortical samples from (*Tasic et al., 2018*; *Zeisel et al., 2018*) used in Figure 2 to decompose NeuroSeq neocortical samples were used to decompose each other, but in the reverse order from the preceding supplementary figure. See Figure 2 for further details of cell identity. Order of samples listed is as in Figure 2.

**Figure 2–Supplement 5.**
**NNLS decomposition of interneuron datasets**. Data from (*Paul et al., 2017*), a third recent single cell study focusing on neocortical interneurons, was used to decompose the cortical interneuron samples from **(A)** (*Tasic et al., 2018*), **(B)** (*Zeisel et al., 2018*), and **(C)** NeuroSeq. In addition, this data set was decomposed using the interneuron samples from the two other single cell data sets **(D,E)**.

**Figure 2–Supplement 6.**
**Random forest decomposition.** A random forest classifier (500 decision trees) was trained from single cell profiles (column labels) and then used to decompose NeuroSeq cell populations (row labels). Coefficients are the ratio of the votes from the 500 trees (coefficient ranges from 0 to 1 and 1 indicates all trees vote for a single class). The pattern of coefficients is similar to that obtained by NNLS (Figure 2) suggesting the decomposition is relatively robust and does not reflect a peculiarity of the NNLS algorithm.

**A**

$$separability = \frac{d_{12}}{d_1 + d_2}$$

$$d_i = mean + 3\sigma$$

1041

**B** Separability of clusters

Tasic 2018     Zeisel 2018     NeuroSeq

**Figure 2–Supplement 7.**

**Separability of cell population clusters**. **(A)** Definition of separability. Cartoon represents two different single cell clusters as distributions of points. The separability is the ratio of the distance between the centroids to the sum of the "diameter" of each cluster. The diameter of a cluster is calculated as the mean distance to the centroid of the cluster + 3 times the standard deviation of the distances of each point in the cluster. With this definition, two clusters are "touching" when separability =1, overlapping when <1, and separate when >1. The multi-dimensional distance is computed as 1- Pearson's corr.coef. Note that averaging is expected to improve separability by roughly the square root of the number of cells averaged, hence most of the improved separability in the NeuroSeq data likely reflects averaging. **(B)** Separabilities between cell population clusters for three datasets shown with two different dynamic ranges (color scale; 0-1 for upper row and 0-10 for lower row). The order of cell population clusters are the same as in Figure 2.

1042

**Figure 3–Supplement 1.**
**Simulated data reveal features of expression metrics**. **(A)** (Upper) An example of simulated binary and graded expression patterns with added noise. X-axis indicates cell populations. (Lower) Various average metrics calculated from the simulated expression patterns (100 individual simulations; error bars are standard deviations). Values are normalized within each metric across binary expression group or graded expression group. **(B)** Summary of each metric's correlation with Mutual Information (MI) and SNR: check mark–correlated, X–uncorrelated, triangle–partially correlated. **(C)** DEF and MI are highly correlated. The relationship between DEF, calculated without considering replicates, and MI with expression levels discretized into 2 levels (left) and 5 levels (right). Although increasing the number of discrete expression levels decreases the degree of correlation, they remain closely related.

**A**  Top 1000 FCR PANTHER enrichment

homeobox transcription factor
transcription factor
DNA binding protein
G-protein coupled receptor
receptor

0  20  40  60
-log10(p value)

**C**  Top1000 FCR GOM enrichment

olfactory receptor activity
transmembrane receptor activity
receptor activity
molecular transducer activity
transmembrane signaling receptor activity
signaling receptor activity
signal transducer activity
seq-specific DNA binding
RNA pol II regulatory region seq-specific DNA binding
RNA pol II regulatory region DNA binding
RNA pol II transcription factor activity, sequence-specific DNA binding
transcription regulatory region seq-specific DNA binding
transcription factor activity, seq-specific DNA binding
nucleic acid binding transcription factor activity
seq-specific double-stranded DNA binding
transcription regulatory region DNA binding
regulatory region DNA binding
regulatory region nucleic acid binding
double-stranded DNA binding
transcriptional activator activity, RNA pol II transcription regulatory region seq-specific binding
transcription factor activity, RNA pol II core promoter proximal region seq-specific binding
transcriptional activator activity, RNA pol II core promoter proximal region seq-specific binding
DNA binding
core promoter proximal region seq-specific DNA binding
G-protein coupled receptor activity
core promoter proximal region DNA binding
RNA pol II core promoter proximal region seq-specific DNA binding

0  25  50  75
-log10(p value)

**B**  Top 1000 DEF PANTHER enrichment

receptor
ion channel
ligand-gated ion channel
GABA receptor
G-protein coupled receptor
transporter
signaling molecule
voltage-gated ion channel
membrane-bound signaling molecule
cell adhesion molecule
voltage-gated potassium channel

0  10  20
-log10(p value)

1043

**D**  Top1000 DEF GOM enrichment

gated channel activity
ion channel activity
substrate-specific channel activity
channel activity
passive transmembrane transporter activity
cation channel activity
transmembrane receptor activity
voltage-gated ion channel activity
voltage-gated channel activity
metal ion transmembrane transporter activity
signaling receptor activity
receptor activity
molecular transducer activity
voltage-gated cation channel activity
inorganic cation transmembrane transporter activity
transmembrane signaling receptor activity
neurotransmitter receptor activity
signal transducer activity
ion transmembrane transporter activity
ligand-gated ion channel activity
ligand-gated channel activity
cation transmembrane transporter activity
substrate-specific transmembrane transporter activity
potassium channel activity
potassium ion transmembrane transporter activity
extracellular ligand-gated ion channel activity
transmembrane transporter activity
substrate-specific transporter activity
calcium ion binding
voltage-gated potassium channel activity
G-protein coupled receptor activity
monovalent inorganic cation transmembrane transporter activity
transporter activity
calcium ion transmembrane transporter activity
ligand-gated cation channel activity
calmodulin binding
calcium channel activity
channel regulator activity
divalent inorganic cation transmembrane transporter activity
sodium ion transmembrane transporter activity
molecular function regulator
receptor binding
peptide receptor activity
protein dimerization activity
sulfur compound binding
protein heterodimerization activity

0  10  20  30
-log10(p value)

**Figure 4–Supplement 1.**

**PANTHER and GO enrichment analysis for high FCR and high DEF genes. (A),(B)** Enrichment using PANTHER gene families. **(C),(D)** Enrichment using Gene Ontology Molecular Function (GOM) categories. Note that GOM does not contain a separate category for homeobox transcription factor, but that these are contained within the parent category: "sequence-specific DNA binding." Red lines indicate the $p = 10^{-5}$ threshold used to judge significance.
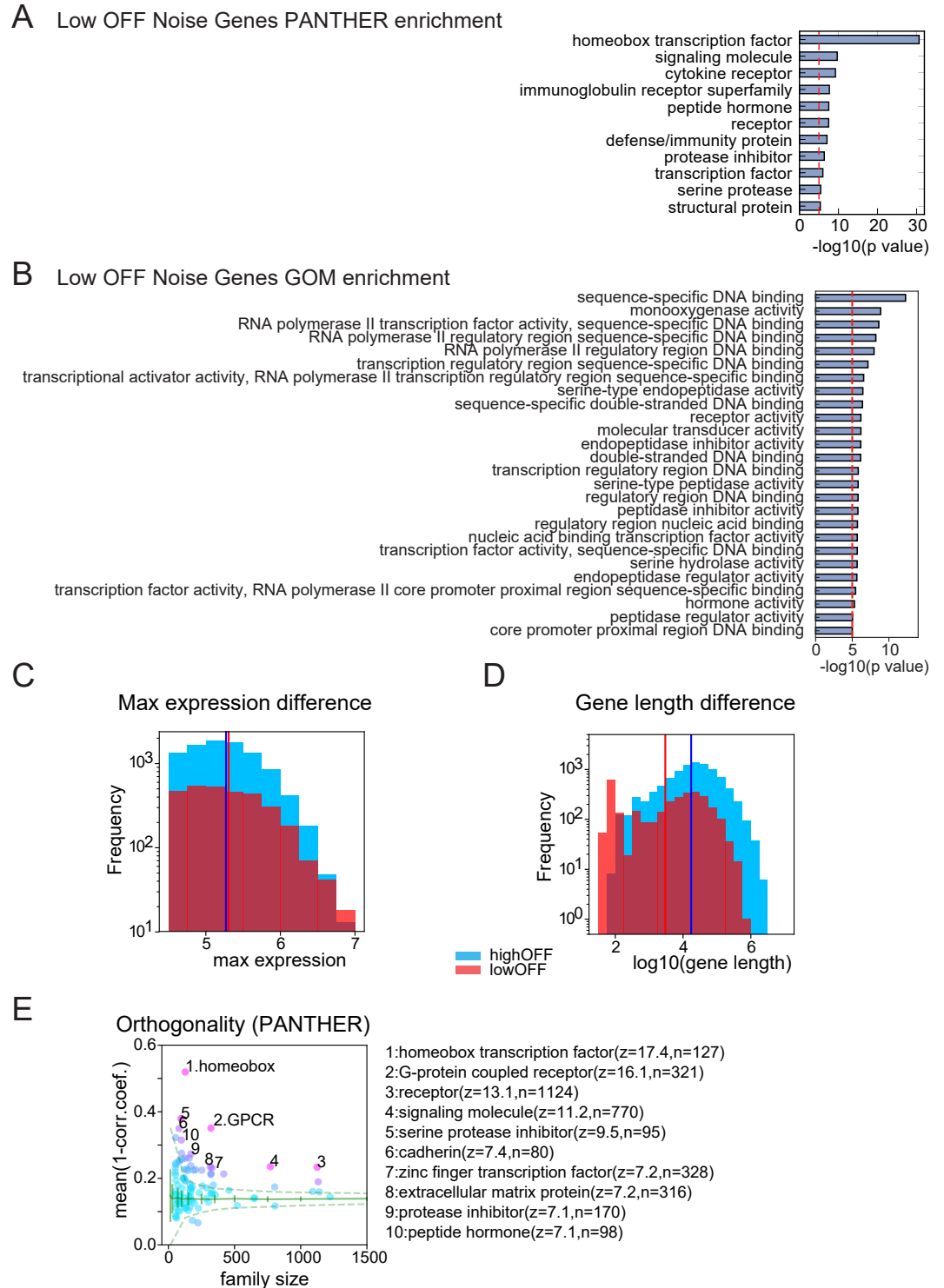
**A** Low OFF Noise Genes PANTHER enrichment

**B** Low OFF Noise Genes GOM enrichment

**C** Max expression difference

**D** Gene length difference

highOFF
lowOFF

**E** Orthogonality (PANTHER)

1:homeobox transcription factor(z=17.4,n=127)
2:G-protein coupled receptor(z=16.1,n=321)
3:receptor(z=13.1,n=1124)
4:signaling molecule(z=11.2,n=770)
5:serine protease inhibitor(z=9.5,n=95)
6:cadherin(z=7.4,n=80)
7:zinc finger transcription factor(z=7.2,n=328)
8:extracellular matrix protein(z=7.2,n=316)
9:protease inhibitor(z=7.1,n=170)
10:peptide hormone(z=7.1,n=98)

1044

**Figure 5–Supplement 1.**

**Properties of Low OFF noise genes.** PANTHER **(A)** and Gene Ontology (GOM: Gene Ontology Molecular functions category) **(B)** enrichments for low OFF noise genes defined by red dashed region in Figure 5B. **(C)** Histogram of max expression for Low OFF noise genes and high OFF noise genes (genes in blue dashed region in Figure 5B). Low OFF genes have slightly higher max expression values than high OFF genes, p=0.002, Students' t-test. Red and blue vertical lines indicate mean values (5.31 and 5.27 respectively). **(D)** Histogram of gene length for Low/high OFF genes. Low OFF noise genes are significantly shorter than high OFF noise genes, p=0 (below machine precision), Student's t-test. Red/blue vertical lines indicate mean values (3.47 and 4.24 respectively). **(E)** Orthogonality, calculated as in Figure 5E, but using the PANTHER gene families.
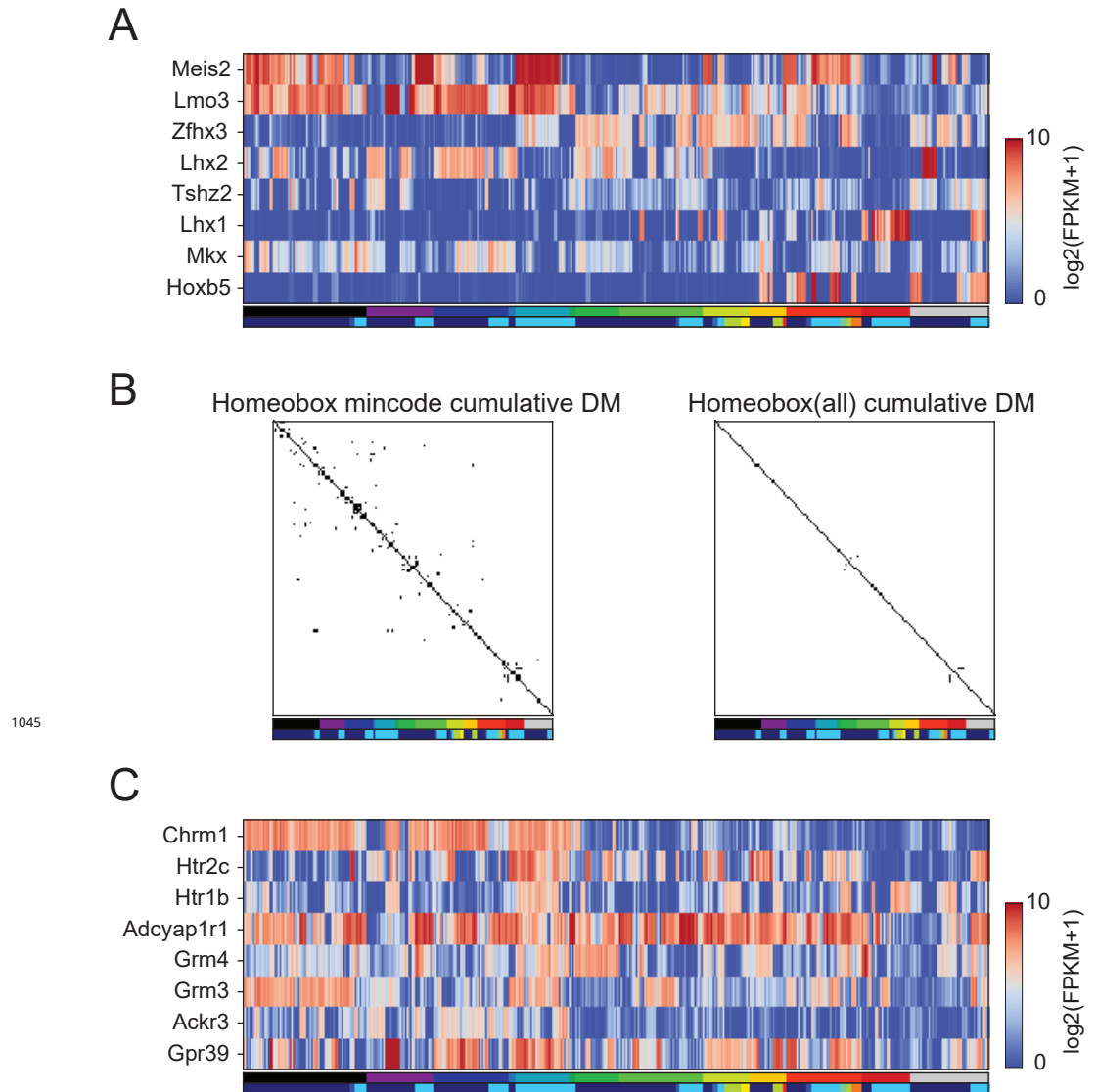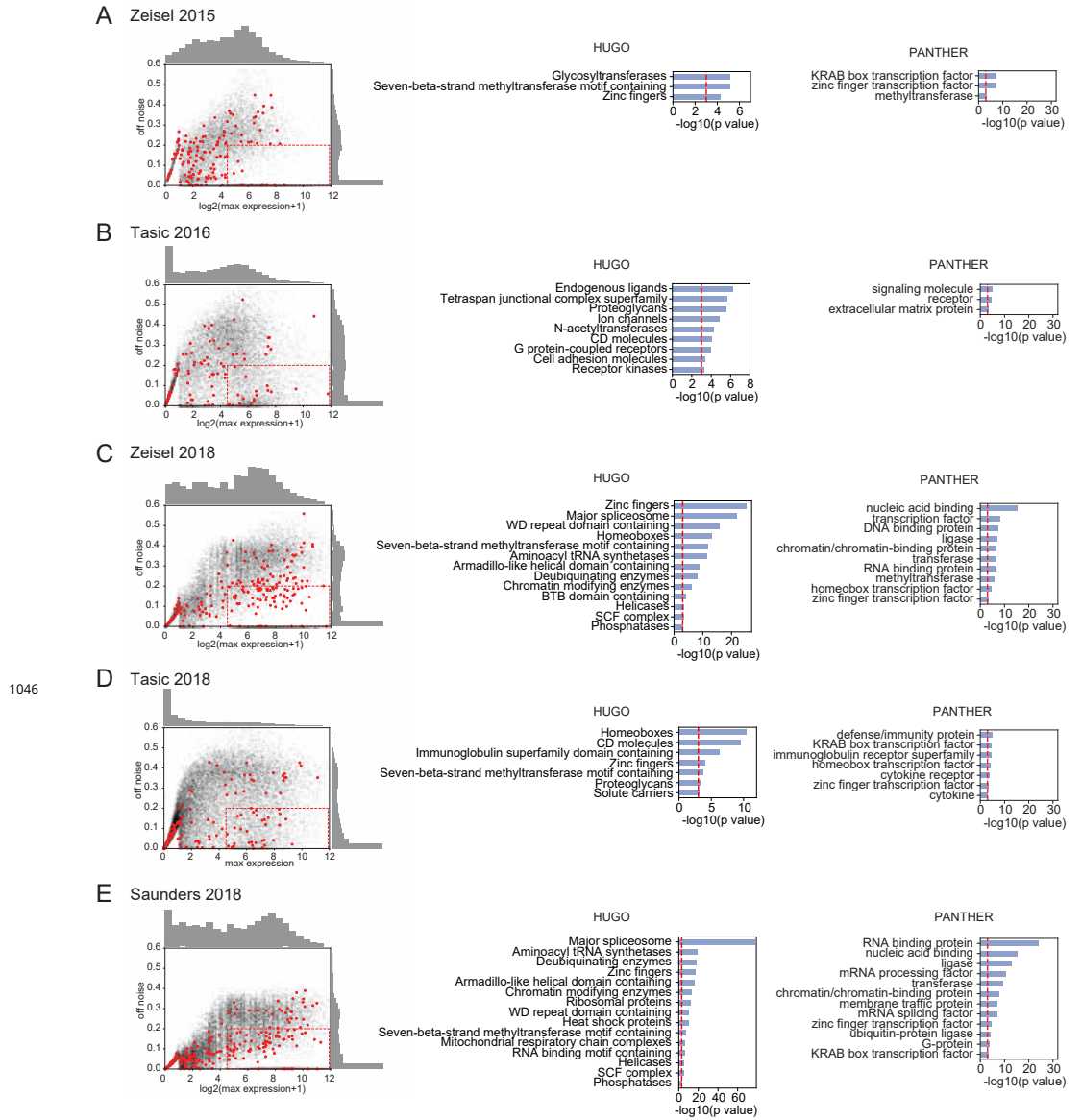
1045

**Figure 5–Supplement 2.**
**Homeobox TFs form a combinatorial code**. **(A)** Heatmap showing expression patterns of 8 homeobox TFs that distinguished 99% of pairs. A minimal gene set algorithm (see Materials and Methods) was used to select these TFs. Each GACP expressed an average of 4.1±1.3 (Std. Dev.) of these TFs. **(B)** Combined DM (differentiation matrix, see Figures 3 and 4C) constructed by allowing GACP pairs to be distinguished on the basis of expression of any of 8 homeobox TFs in the minimal set (left) or by any homeobox TFs (right). White indicates distinguishable pairs and black indicates indistinguishable pairs. **(C)** Heatmap showing expression patterns of minimal gene sets for GPCR capable of distinguishing 99% of pairs.

1046

**Figure 5–Supplement 3.**

**OFF noise in single cell datasets**. **(Left column)** OFF noise calculated as in Figure 5B from the standard deviation of cluster averages, plotted against the maximum expression. Red dots are homeobox transcription factors, black dots are all other genes. **(Middle, Right columns)**. HUGO gene groups and PANTHER protein families over-represented in the dashed red boxes in the OFF noise plots. Datasets are from (*Zeisel et al., 2015*; *Tasic et al., 2016*; *Zeisel et al., 2018*; *Saunders et al., 2018*; *Tasic et al., 2018*).

# Homeobox subfamilies differ in FCR/OFF noise



**Figure 5–Supplement 4.**
**OFF noise and gene length in Homeobox subfamilies**. **(Left)** Scatter plot of mean gene length and mean FCR for homeobox subfamilies. Subfamilies are as defined according to HUGO gene groups. **(Right)** Scatter plot of mean gene length and mean OFF noise for homeobox subfamilies.

**Figure 7–Supplement 1.**
**DEF length bias in SC datasets**. DEF is plotted against $\log_{10}$(gene length) for five SC RNA-seq datasets from (*Zeisel et al., 2015*; *Tasic et al., 2016*; *Zeisel et al., 2018*; *Saunders et al., 2018*; *Tasic et al., 2018*). Red bars represent average DEF for genes binned by gene length (4 bins per log unit), sorted by length.
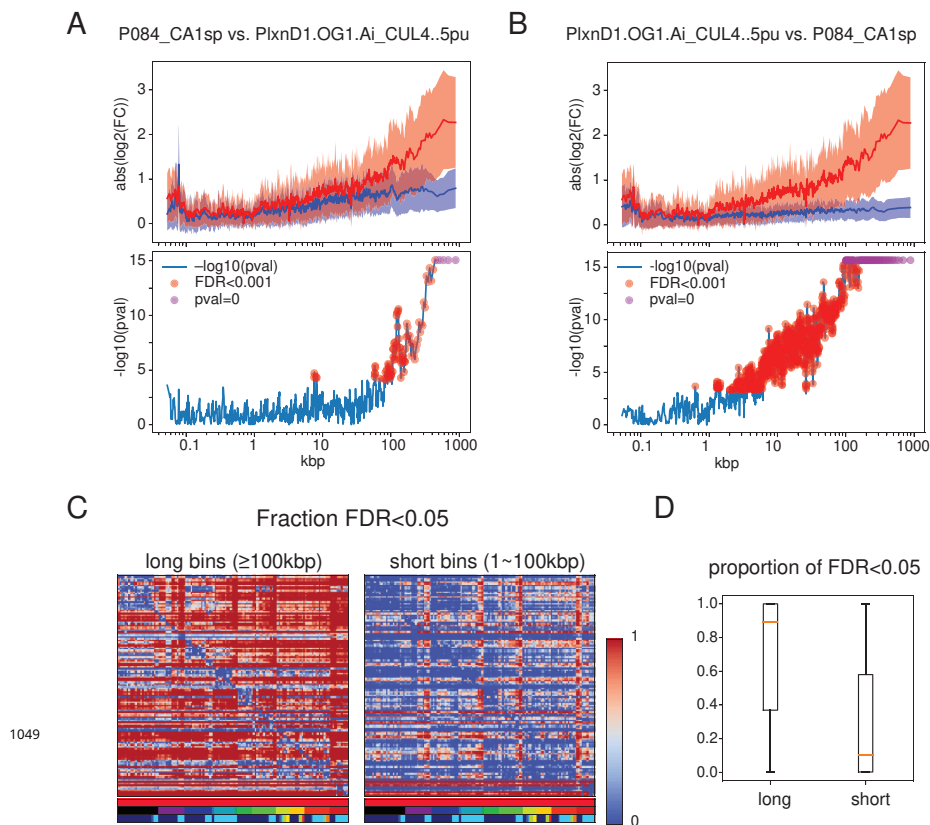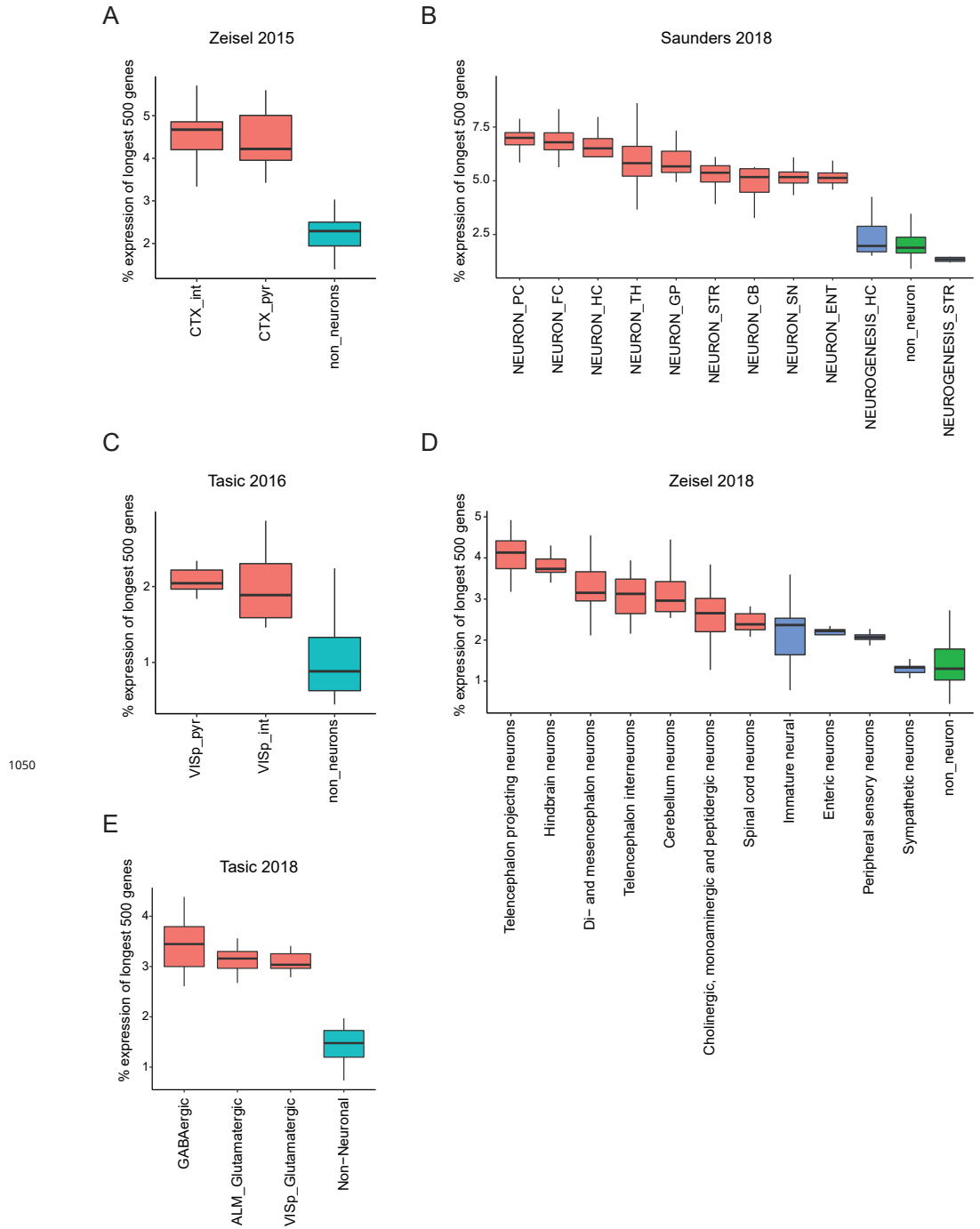
**Figure 7–Supplement 2.**

**Significant length differences using the test proposed by Raman et al.** Raman et al. (2018) propose evaluating length dependent differences by comparing expression ratios between groups to those within a single group. **(A,B)** (Top panels) Mean and standard deviation of $abs(log_2((mean(Grp2)+1)/(mean(Grp1)+1))$. Blue: same for $abs(log_2((Grp1_1+1)/(Grp1_2+1))$. For **A**, Grp1 is P084_CAsp and Grp2 is PlxnD1.OG1.Ai_CUL4..5pu. For **B**, groups are reversed. Note that the results are not symmetric because the proposed test makes use of baseline variance in only one of the two groups. (Bottom panels) Negative $log_{10}(pvalue)$ for each bin. P-values are calculated by Student's t-test (two-sided, unequal variance). Red dots indicate bins with FDR<0.001. FDR (multiple tests correction) is calculated using all bins (n=1245). Some bins have p-values below the machine precision (double float; ∼1e-308) indicated as pval=0 (magenta dots). **(C)** Matrices of the fraction of significant long (Left) and short (right) bins calculated using the Raman et al. test. Horizontal color legends below each matrix label populations as in Figures 1 Supplement 1, and Figures 4, 5, 6: top row:sample type (red indicates all are sorted neurons), second row: brain region, third row: transmitter. Vertical color bar indicates fraction of gene bins that are significant. The matrices are asymmetric because test significance can vary depending on which population is used to calculate baseline FC. **(D)** Boxplots showing median (orange bar), and first to third interquartile ranges (boxes) for the same data shown in matrix form above.

**Figure 7–Supplement 3.**
**Regional bias of long gene expression in SC datasets**. Percentage of expression of the longest 500 genes in four single-cell datasets. Boxes show median and quartiles. Whiskers extend to 1.5 x inter-quartile range. CNS neurons are shown in red. Immature or PNS neurons are shown in blue. Nonneurons shown in green. Abbreviations: CTX: cortex, pyr: pyramidal, int: interneuron, PC: posterior cortex, FC: frontal cortex, HC: hippocampus, TH: thalamus, GP: globus pallidus externus & nucleus basalis, STR: striatum, CB: cerebellum, SN: substantia nigra and ventral tegmental area, Ent: Enteropeduncular nucleus and subthalamic nucleus, VisP: primary visual area, ALM: anterior lateral motor cortex.