

1 Cross-laboratory analysis of brain cell type transcriptomes with applications to interpretation of bulk
2 tissue data

3 B. Ogan Mancarci^{1,2,3}, Lilah Toker^{2,3}, Shreejoy J Tripathy^{2,3}, Brenna Li^{2,3}, Brad Rocco^{4,5}, Etienne
4 Sibille^{4,5}, Paul Pavlidis^{2,3*}

5 ¹Graduate Program in Bioinformatics, University of British Columbia, Vancouver, Canada

6 ²Department of Psychiatry, University of British Columbia, Vancouver, Canada

7 ³Michael Smith Laboratories, University of British Columbia, Vancouver, Canada

8 ⁴Campbell Family Mental Health Research Institute of CAMH

9 ⁵Department of Psychiatry and the Department of Pharmacology and Toxicology, University of Toronto, Toronto, Canada.

10

11 Address correspondence to;

12 Paul Pavlidis, PhD

13 177 Michael Smith Laboratories 2185 East Mall

14 University of British Columbia Vancouver BC V6T1Z4

15 604 827 4157 paul@msl.ubc.ca

16

17 Ogan Mancarci: ogan.mancarci@msl.ubc.ca

18 Lilah Toker: ltoker@msl.ubc.ca

19 Shreejoy Tripathy: stripathy@msl.ubc.ca

20 Brenna Li: brenna.li@msl.ubc.ca

21 Brad Rocco: Brad.Rocco@camh.ca

22 Etienne Sibille: Etienne.Sibille@camh.ca

23 **Abstract**

24 Establishing the molecular diversity of cell types is crucial for the study of the nervous system. We
 25 compiled a cross-laboratory database of mouse brain cell type-specific transcriptomes from 36
 26 major cell types from across the mammalian brain using rigorously curated published data from
 27 pooled cell type microarray and single cell RNA-sequencing studies. We used these data to
 28 identify cell type-specific marker genes, discovering a substantial number of novel markers, many
 29 of which we validated using computational and experimental approaches. We further demonstrate
 30 that summarized expression of marker gene sets in bulk tissue data can be used to estimate the
 31 relative cell type abundance across samples. Using this approach, we show that majority of genes
 32 previously reported as differentially expressed in Parkinson's disease can be attributed to the
 33 reduction in dopaminergic cell number rather than regulatory events. To facilitate use of this
 34 expanding resource, we provide a user-friendly web-interface at Neuroexpresso.org.

35 **Introduction**

36 Brain cells can be classified based on features such as their primary type (e.g. neurons vs. glia),
 37 location (e.g. cortex, hippocampus, cerebellum), electrophysiological properties (e.g. fast spiking
 38 vs. regular spiking), morphology (e.g. pyramidal cells, granule cells) or the
 39 neurotransmitter/neuromodulator they release (e.g. dopaminergic cells, serotonergic cells,
 40 GABAergic cells). Marker genes, genes that are expressed in a specific subset of cells, are often
 41 used in combination with other cellular features to define different types of cells (Hu et al., 2014;
 42 Margolis et al., 2006) and facilitate their characterization by tagging the cells of interest for further
 43 studies (Handley et al., 2015; Lobo et al., 2006; Tomomura et al., 2001). Marker genes have also
 44 found use in the analysis of whole tissue "bulk" gene expression profiling data, which can be
 45 challenging to interpret due to the difficulty to determine the source of the observed expressional
 46 change. For example, a decrease in a transcript level can indicate a regulatory event affecting the
 47 expression level of the gene, a decrease in the number of cells expressing the gene, or both. To

address this issue, computational methods have been proposed to estimate cell type specific proportion changes based on expression patterns of known marker genes (Chikina et al., 2015; Newman et al., 2015; Westra et al., 2015; Xu et al., 2013). Finally, marker genes are obvious candidates for having cell type specific functional roles.

An ideal cell type marker has a strongly enriched expression in a single cell type in the brain. However, this criterion can rarely be met, and for many purposes, cell type markers can be defined within the context of a certain brain region; namely, a useful marker may be specific for the cell type in one region but not necessarily in another region or brain-wide. For example, the calcium binding protein parvalbumin is a useful marker of both fast spiking interneurons in the cortex and Purkinje cells in the cerebellum (Celio and Heizmann, 1981; Kawaguchi et al., 1987). Whether the markers are defined brain-wide or in a region-specific context, the confidence in their specificity is established by testing their expression in as many different cell types as possible. This is important because a marker identified by comparing just two cell types might turn out to be expressed in a third, untested cell type, reducing its utility.

During the last decade, targeted purification of cell types of interest followed by gene expression profiling has been applied to many cell types in the brain. Such studies, targeted towards well-characterized cell types, have greatly promoted our understanding of the functional and molecular diversity of these cells (Cahoy et al., 2008; Chung et al., 2005; Doyle et al., 2008). However, individual studies of this kind are limited in their ability to discover specific markers as they often analyse only a small subset of cell types (Shrestha et al., 2015; Okaty et al., 2009; Sugino et al., 2006) or have limited resolution as they group subtypes of cells together (Cahoy et al., 2008). Recently, advances in technology have enabled the use of single cell transcriptomics as a powerful tool to dissect neuronal diversity and derive novel molecular classifications of cells (Poulin et al., 2016). However, with single cell analysis the classification of cells to different types is generally done post-hoc, based on the clustering similarity in their gene expression patterns. These molecularly defined cell types are often uncharacterized otherwise (e.g. electrophysiologically, morphologically), challenging their identification outside of the original study and understanding

their role in normal and pathological brain function. A notable exception is the single cell RNA-seq study of Tasic et al. (2016) analysing single labelled cells from transgenic mouse lines to facilitate matching of the molecularly defined cell types they discover to previously identified cell types. We hypothesized that aggregating cell type specific studies that analyse expression profiles of cell types previously defined in literature, a more comprehensive data set more suitable for marker genes could be derived.

Here we report the analysis of an aggregated cross-laboratory dataset of cell type specific expression profiling experiments from mouse brain, composed both of pooled cell microarray data and single cell RNA-seq data. We used these data to identify sets of brain cell marker genes more comprehensive than any previously reported, and validated the markers genes in external mouse and human single cell datasets. We further show that the identified markers are applicable for the analysis of human brain and demonstrate the usage of marker genes in the analysis of bulk tissue data via the summarization of their expression into “marker gene profiles” (MGPs), which can be cautiously interpreted as correlates of cell type proportion. Finally, we made both the cell type expression profiles and marker sets available to the research community at neuroexpresso.org.

Results

Compilation of a brain cell type expression database

A key input to our search for marker genes is expression data from purified pooled brain cell types and single cells. Expanding on work from Okaty et al. (2011), we assembled and curated a database of cell type-specific expression profiles from published data (see Methods, Figure 1A). The database represents 36 major cell types from 12 brain regions (Figure 1B) from a total of 263 samples and 30 single cell clusters. Frontal cortex is represented by both microarray and RNA-seq data, with 5 of the 15 cortical cell types represented exclusively by RNA-seq data. We used rigorous quality control steps to identify contaminated samples and outliers (see Methods). In the microarray dataset, all cell types except for ependymal cells are represented by at least 3

replicates and in the entire database, 14/36 cell types are represented by multiple independent studies (Table 1). The database is in constant growth as more cell type data becomes available. To facilitate access to the data and allow basic analysis we provide a simple search and visualization interface on the web, www.neuroexpresso.org (Figure 1C). The app provides means of visualising gene expression in different brain regions based on the cell type, study or methodology, as well as differential expression analysis between groups of selected samples.

Cell Type	Sample count	Marker gene count	GEO accession and reference
Whole Brain			
Astrocyte	9 / 1*	94**	GSE9566 (Cahoy et al., 2008), GSE35338 (Zamanian et al., 2012), GSE71585 (Tasic et al., 2016)
Oligodendrocyte	25 / 1*	22**	GSE48369, (Bellesi et al., 2013), GSE9566 (Cahoy et al., 2008), GSE13379 (Doyle et al., 2008), GSE30016 (Fomchenko et al., 2011), GSE71585 (Tasic et al., 2016)
Microglia	3 / 1*	131**	GSE29949 (Anandasabapathy et al., 2011), GSE71585 (Tasic et al., 2016)
Cortex			
FS Basket (G42)	13 / 5*	18	GSE17806 (Okaty et al., 2009), GSE8720 (Sugino et al., 2014), GSE2882 (Sugino et al., 2006), GSE71585 (Tasic et al., 2016)
Martinotti (GIN)	3 / 1*	15	GSE2882 (Sugino et al., 2006), GSE71585 (Tasic et al., 2016)
VIPReIn (G30)	6 / 1*	33	GSE2882 (Sugino et al., 2006), GSE71585 (Tasic et al., 2016)
Pan-Pyramidal***	9 / 17 *	35	See below
Pyramidal cortico-thalamic	3 / 2*	2	GSE2882 (Schmidt et al., 2012), GSE71585 (Tasic et al., 2016)
Pyramidal Glt25d2	3 / 2*	3	GSE35758 (Schmidt et al., 2012), GSE71585 (Tasic et al., 2016)
Pyramidal S100a10	3 / 4*	2	GSE35751 (Schmidt et al., 2012), GSE71585 (Tasic et al., 2016)
Layer 2 3 Pyra	2*	3	GSE71585 (Tasic et al., 2016)
Layer 4 Pyra	3*	5	GSE71585 (Tasic et al., 2016)
Layer 6a Pyra	2*	6	GSE71585 (Tasic et al., 2016)
Layer 6b Pyra	2*	9	GSE71585 (Tasic et al., 2016)
Oligodendrocyte precursors	1*	184	GSE71585 (Tasic et al., 2016)
Endothelial	2*	178	GSE71585 (Tasic et al., 2016)
BasalForebrain			
Forebrain cholinergic	3	90	GSE13379 (Doyle et al., 2008)
Striatum			
Forebrain cholinergic	3	45	GSE13379 (Doyle et al., 2008)
Medium spiny neurons	39	74	GSE13379 (Doyle et al., 2008), GSE55096 (Heiman et al., 2014), GSE54656 (Maze et al., 2014), GSE48813 (C. L. Tan et al., 2013)
Amygdala			
Glutamatergic	3	10	GSE2882 (Sugino et al., 2006)
Pyramidal Thy1 Amyg	12	21	GSE2882 (Sugino et al., 2006)

106 Table 1: Cell types in the NeuroExpresso database

* The number of clusters from RNA-seq data.

** Marker genes for these cell types are identified in multiple regions displayed yet only the number of the genes that are found in the region specified on the table is shown for the sake of conservation of space. Astrocytes, microglia and oligodendrocyte markers are identified in the context of all other brain regions (except cerebellum for astrocytes) and dopaminergic markers are also identified for midbrain.

Hippocampus			
DentateGranule	3	17	GSE11147 (Perrone-Bizzozero et al., 2011)
GabaSSTReIn	3	54	GSE2882 (Sugino et al., 2006)
Pyramidal Thy1 Hipp	12	17	GSE2882 (Sugino et al., 2006)
Subependymal			
Ependymal	2	50	GSE18765 (Beckervordersandforth et al., 2010)
Thalamus			
GabaReIn	3	53	GSE2882 (Sugino et al., 2006)
Hypocretinergic	4	35	GSE38668 (Dalal et al., 2013)
Thalamus cholinergic	3	40	GSE43164 (Görllich et al., 2013)
Midbrain			
Midbrain cholinergic	3	34	GSE13379 (Doyle et al., 2008)
Serotonergic	3	18	GSE36068 (Dougherty et al., 2013)
Substantia nigra			
Dopaminergic	30	58**	No accession **** (Chung et al., 2005), GSE17542 (Phani et al., 2010)
LocusCoeruleus			
Noradrenergic	9	133	GSE8720 (Sugino et al., 2014), No accession**** (Sugino et al. Unpublished)
Cerebellum			
Basket	16	6	GSE13379 (Doyle et al., 2008), GSE37055 (Paul et al., 2012)
Bergmann	3	52	GSE13379 (Doyle et al., 2008)
Cerebral granule cells	3	11	GSE13379 (Doyle et al., 2008)
Golgi	3	26	GSE13379 (Doyle et al., 2008)
Purkinje	44	43	GSE13379 (Doyle et al., 2008), GSE57034 (Galloway et al., 2014), GSE37055 (Paul et al., 2012), No accession**** (Rossner et al., 2006), GSE8720 (Sugino et al., 2014), No accession**** Sugino et al. unpublished
SpinalCord			
Spinal cord cholinergic	3	124	GSE13379 (Doyle et al., 2008)

107 Table 1 Continued.

Sample count - number of samples that representing the cell type; Gene count - number of marker genes detected for cell type.

** Marker genes for these cell types are identified in multiple regions displayed yet only the number of the genes that are found in the region specified on the table is shown for the sake of conservation of space. Astrocytes, microglia and oligodendrocyte markers are identified in the context of all other brain regions (accept cerebellum for astrocytes) and dopaminergic markers are also identified for midbrain.

*** Pan-pyramidal is a merged cell type composed of all pyramidal samples.

**** Data obtained directly from authors.

Identification of cell type enriched marker gene sets

We used the NeuroExpresso data to identify marker gene sets (MGSs) for each of the 36 cell types. An individual MGS is composed of genes highly enriched in a cell type in the context of a brain region (Figure 2A). Marker genes were selected based on a) fold of change relative to other cell types in the brain region and b) a lack of overlap of expression levels in other cell types (see Methods for details). This approach captured previously known marker genes (e.g. Th for dopaminergic cells (Pickel et al., 1976), Tmem119 for microglia (Bennett et al., 2016) (corroborating previous reports Satoh et al. (2016), Erny et al. (2015), this gene was classified as downregulated in activated microglia in our analysis) along with numerous new candidate markers such as Cox6a2 for fast spiking parvalbumin (PV)⁺ interneurons. Some marker genes previously reported by individual studies whose data were included in our database, were not selected by our analysis. For example, Fam114a1 (9130005N14Rik), identified as a marker of fast spiking basket cells by Sugino et al. (2006), is highly expressed in oligodendrocytes and oligodendrocyte precursor cells (Figure 2B). These cell types were not considered in the Sugino et al. (2006) study, and thus the lack of specificity of Fam114a1 could not be observed by the authors. In total, we identified 2671 marker genes, with 3-186 markers per cell type (Table 1). The next sections focus on verification and validation of our proposed markers, using multiple methodologies.

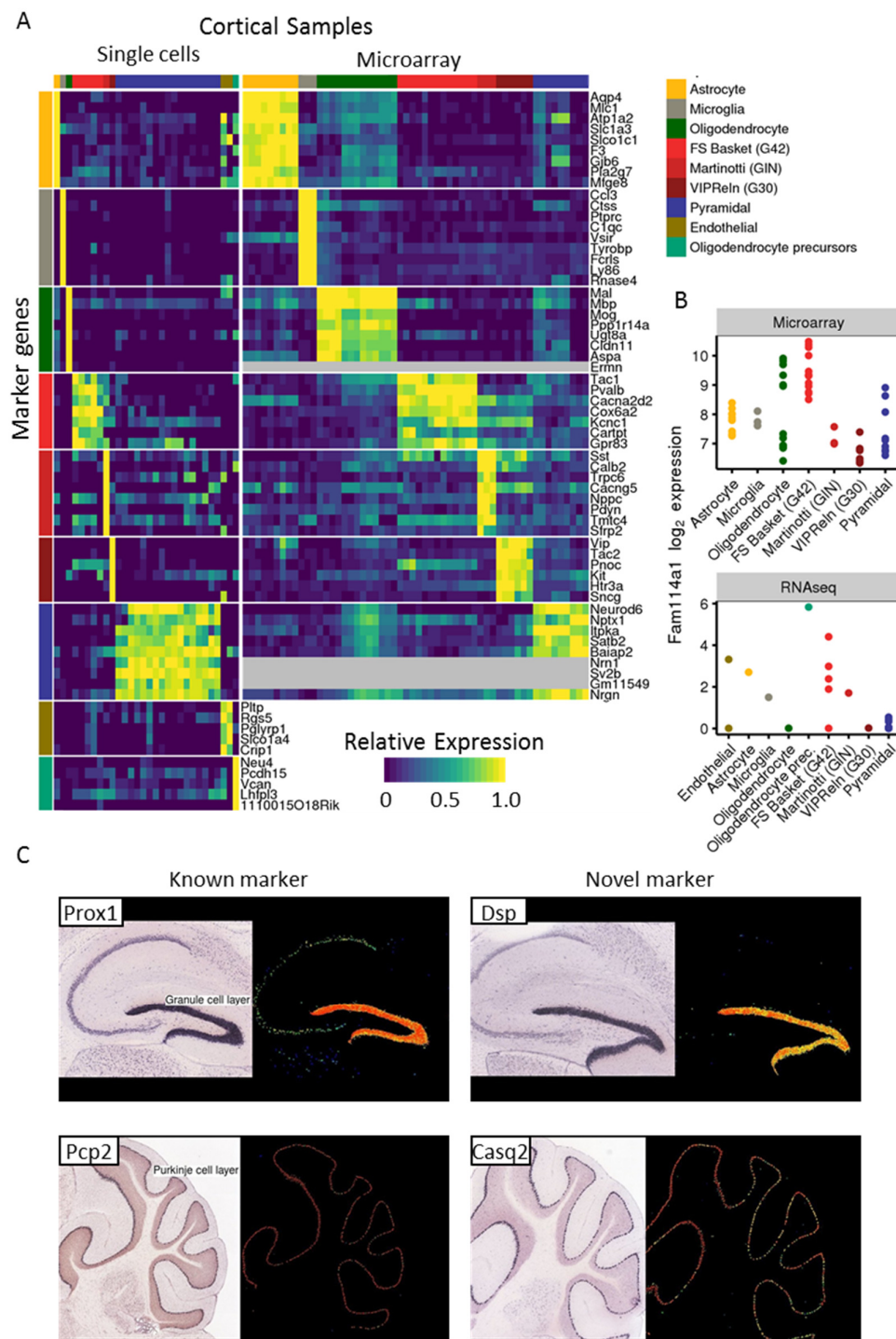


Figure 2: Marker genes are selected for mouse brain cell types and used to estimate cell type profiles. **(A)** Expression of top marker genes selected for cell cortical cell types in cell types represented by RNA-seq (left) and microarray (right) data in NeuroExpresso. Expression levels were normalized per gene to be between 0-1 for each dataset. **(B)** Expression of Fam114a1 in frontal cortex in microarray (left) and RNA-seq (right) datasets. Fam114a1 is a proposed fast spiking basket cell marker. It was not selected as a marker in

this study due to its high expression in oligodendrocytes and S100a10 expressing pyramidal cells that were both absent from the original study. **(C)** In situ hybridization images from the Allen Brain Atlas. Rightmost panels show the location of the image in the brain according to the Allen Brain mouse reference atlas. Panels on the left show the ISH image and normalized expression level of known and novel dentate granule (upper panels) and Purkinje cell (lower panels) markers.

Verification of markers by in situ hybridization

Two cell types in our database (Purkinje cells of the cerebellum and hippocampal dentate gyrus granule cells) are organized in well-defined anatomical structures that can be readily identified in tissue sections. We exploited this fact to use in situ hybridization (ISH) data from the Allen Brain Atlas (ABA) (<http://mouse.brain-map.org>) (Sun et al., 2013) to verify co-localization of known and novel markers for these two cell types. There was a high level of agreement (Figure S1-S2) that the markers were correctly localized to the corresponding brain structures, and by implication, cell types. For dentate granule (DG) cell markers, all 16 genes were represented in ABA. Of these, 14 specifically co-localized with known markers (that is, had the predicted expression pattern confirming our marker selection), one marker exhibited non-specific expression and one marker showed no signal. For Purkinje cell markers, 41/43 genes were represented in ABA. Of these, 37 specifically co-localized with known markers, one marker exhibited non-specific expression and three markers showed no signal in the relevant brain structure. Figure 2C shows representative examples for the two cell types (details of our ABA analysis, including images for all the genes examined and validation status of the genes, are provided in the supplement (Figure S1-S2, Table S1-S2). The four markers for which no signal was detected (one marker of dentate gyrus granule cells and three markers of Purkinje cells) underwent additional scrutiny. For one of the markers of Purkinje cells (Eps8l2), the staining of cerebellar sections was inconsistent, with some sections showing no staining, some sections showing nonspecific staining and several sections showing the predicted localization. The three remaining genes had no signal in ABA ISH data brain-wide. We considered such absence or inconsistency of ISH signal inconclusive. Further analysis of these cases (one DG marker, three Purkinje) suggests that the ABA data is the outlier. As part of our marker selection procedure, Pter, the DG cell marker in question, was found to have high

expression in granule cells both within NeuroExpresso and Hipposeq – a data set that is not used for primary selection of markers(see methods). In addition, Hipposeq indicates specificity to DG cells relative to the other neuron types in Hipposeq. For the Purkinje markers, specific expression for one gene (Sycp1) was supported by the work of Rong et al. (2004), who used degeneration of Purkinje cells to identify potential markers of these cells (20/43 Purkinje markers identified in our study were also among the list of potential markers reported by Rong et al. (Table S3)). We could not find data to further establish expression for the two remaining markers of Purkinje cells (Eps8l2 and Smpx). However, we stress that the transcriptomic data for Purkinje cells in our database are from five independent studies using different methodologies for cell purification, all of which support the specific expression of Eps8l2 and Smpx in Purkinje cells. Overall, through a combination of examination of ABA and other data sources, we were able to find confirmatory evidence of cell-type-specificity for 53/57 genes, with two false positives, and inconclusive findings for two genes.

We independently verified Cox6a2 as a marker of cortical fast spiking PV+ interneurons using triple label in situ hybridization of mouse cortical sections for Cox6a2, Pvalb and Slc32a1 (a pan-GABAergic neuronal marker) transcripts. As expected, we found that approximately 25% of all identified neurons were GABAergic (that is, Slc32a1 positive), while 46% of all GABAergic neurons were also Pvalb positive. 80% of all Cox6a2+ neurons were Pvalb and Slc32a1 positive whereas Cox6a2 expression outside GABAergic cells was very low (1.65% of Cox6a2 positive cells), suggesting high specificity of Cox6a2 to PV+ GABAergic cells (Figure 3).

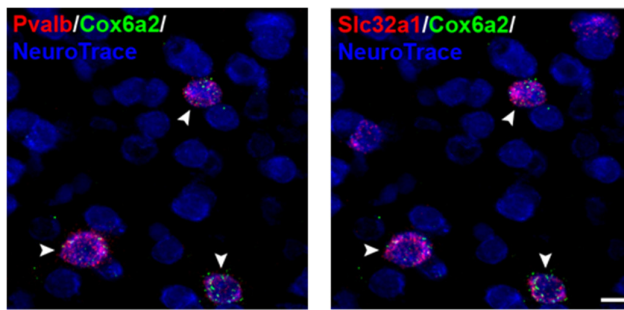


Figure 3: Single-plane image of mouse sensorimotor cortex labeled for Pvalb, Slc32a1, and Cox6a2 mRNAs and counterstained with NeuroTrace. Arrows indicate Cox6a2+ neurons. Bar = 10 μ m.

Verification of marker gene sets in single-cell RNA-seq data

As a further validation of our marker gene signatures, we analysed their properties in recently published single cell RNA-seq datasets derived from mouse cortex (Zeisel et al., 2015) and human cortex. We could not directly compare our MGSs to markers of cell type clusters identified in the studies producing these datasets since their correspondence to the cell types in NeuroExpresso was not clear. However, since both datasets represent a large number of individual cells, they are likely to include cells that correspond to the cortical cell types in our database. Thus if our MGSs are cell type specific, and the corresponding cells are present in the single cell datasets, MGS should have a higher than random rate of being co-detected in the same cells, relative to non-marker genes. A weakness of this approach is that a failure to observe a correlation might be due to absence of the cell type in the data set rather than a true shortcoming of the markers. Overall, all MGSs for all cell types with the exception of oligodendrocyte precursor cells were successfully validated ($p < 0.001$, Wilcoxon rank sum test) in both single cell datasets (Table 2).

	Zeisel et al. (mouse)		Darmanis et al. (human)	
Cell Types	p-value	Gene Count	p-value	Gene Count
Endothelial	p<0.001	180	p<0.001	157
Astrocyte	p<0.001	282	p<0.001	239
Microglia	p<0.001	248	p<0.001	201
Oligodendrocyte	p<0.001	156	p<0.001	201
Oligodendrocyte precursors	0.831	193	0.999	203
FS Basket (G42)	p<0.001	26	p<0.001	26
Martinotti (GIN)	p<0.001	21	p<0.001	20
VIPReIn (G30)	p<0.001	43	p<0.001	36
Pyramidal	p<0.001	34	p<0.001	27

Table 2: Coexpression of cortical MGSs in single cell RNA-seq data.

NeuroExpresso as a tool for understanding the biological diversity and similarity of brain cells

One of the applications of NeuroExpresso is as an exploratory tool for exposing functional and biological properties of cell types. In this section, we highlight three examples we encountered: We observed high expression of genes involved in GABA synthesis and release (Gad1, Gad2 and Slc32a1) in forebrain cholinergic neurons, suggesting the capability of these cells to release GABA in addition to their cognate neurotransmitter acetylcholine (Figure 4A). Indeed, co-release of GABA and acetylcholine from forebrain cholinergic cells was recently demonstrated by Saunders et al. (2015). Similarly, the expression of the glutamate transporter Slc17a6, observed in thalamic (habenular) cholinergic cells suggests co-release of glutamate and acetylcholine from these cells, recently supported experimentally (Ren et al., 2011) (Figure 4A). We observed consistently high expression of Ddc (Dopa Decarboxylase), responsible for the second step in the monoamine synthesis pathway. This surprising result is suggestive of a previously unknown ability of oligodendrocytes to produce monoamine neurotransmitters upon exposure to appropriate precursor, as previously reported for several populations of cells in the brain (Ren et al., 2016; Ugrumov, 2013). Alternatively, this finding might indicate a previously unknown function of Ddc. Lastly, we found overlap between the markers of spinal cord and brainstem cholinergic cells, and of midbrain noradrenergic cells, suggesting previously unknown functional similarity between cholinergic and noradrenergic cell types. The common markers included Chodl, Calca, Cda and

Hspb8, which were recently confirmed to be expressed in brainstem cholinergic cells (Enjin et al., 2010), and Phox2b, a known marker of noradrenergic cells (Pattyn et al., 1997).

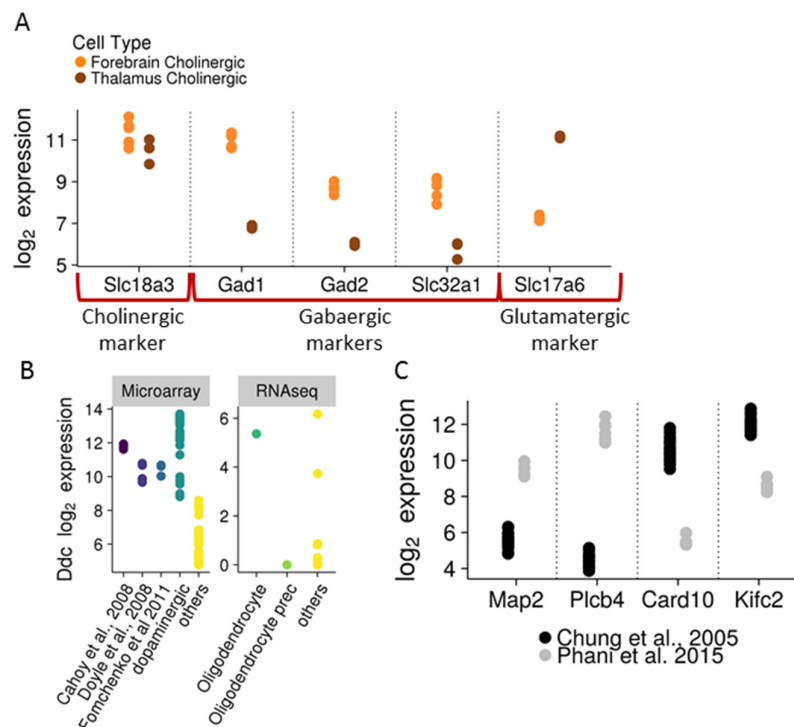


Figure 4: NeuroExpresso reveals gene expression patterns. (A) Expression of cholinergic, GABAergic and glutamatergic markers in cholinergic cells from forebrain and thalamus. Forebrain cholinergic neurons express GABAergic markers while thalamus (hubenular) cholinergic neurons express glutamatergic markers. **(B)** (Left) Expression of *Ddc* in oligodendrocyte samples from Cahoy et al., Doyle et al. and Fomchenko et al. datasets and in comparison to dopaminergic cells and other (non-oligodendrocyte) cell types from the frontal cortex in the microarray dataset. In all three datasets expression of *Ddc* in oligodendrocytes is comparable to expression in dopaminergic cells and is higher than in any of the other cortical cells. Oligodendrocyte samples show higher than background levels of expression across datasets. (Right) *Ddc* expression in oligodendrocytes, oligodendrocyte precursors, and other cell types from Tasic et al. single cell dataset. **(C)** Bimodal gene expression in two dopaminergic cell isolates by different labs. Genes shown are labeled as marker genes in the context of midbrain if the two cell isolates are labeled as different cell types.

Marker Gene Profiles can be used to infer changes in cellular proportions in the brain

Marker genes are by definition cell type specific, and thus changes in their expression observed in bulk tissue data can represent either changes in the number of cells or cell type specific

transcriptional changes (or a combination). Marker genes of four major classes of brain cell types (namely neurons, astrocytes, oligodendrocytes and microglia) were previously used to gain cell type specific information from brain bulk tissue data (Bowling et al., 2016; Hagenauer et al., 2016; Kuhn et al., 2011; Sibille et al., 2008; Skene and Grant, 2016; P. P. C. Tan et al., 2013), and generally interpreted as reflecting changes in proportions. To apply this approach using our markers, following the practice of others, we summarize the expression profiles of marker genes as the first principal component of their expression (see Methods) (Chikina et al., 2015; Westra et al., 2015; Xu et al., 2013). We refer to these summaries as Marker Gene Profiles (MGPs).

In order to validate the use of MGPs as surrogates for relative cell type proportions, we used bulk tissue expression data from conditions with known changes in cellular proportions. Firstly, we calculated MGPs for human white matter and frontal cortex using data collected by (Trabzuni et al., 2013). Comparing the MGPs in white vs. grey matter, we observed the expected increase in oligodendrocyte MGP, as well as increase in oligodendrocyte progenitor cell, endothelial cell, astrocyte and microglia MGPs, corroborating previously reported higher number of these cell types in white vs. grey matter (Gudi et al., 2009; Ogura et al., 1994; Williams et al., 2013). We also observed decrease in MGPs of all neurons, corroborating the low neuronal cell body density in white vs. grey matter (Figure 5A).

A more specific form of validation was obtained from a pair of studies done on the same cohort of subjects, with one study providing expression profiles (study 2 from SMRI microarray database, see Methods) and another providing stereological counts of oligodendrocytes (Uranova et al., 2004), for similar brain regions. We calculated oligodendrocyte MGPs based on the expression data and compared the results to experimental cell counts from Uranova et al. (2004). The MGPs were consistent with the reduction of oligodendrocytes observed by Uranova et al. in schizophrenia, bipolar disorder and depression patients. (Figure 5B; direct comparison between MGP and experimental cell count at a subject level was not possible, as Uranova et al. did not provide subject identifiers corresponding to each of the cell count values).

To further assess and demonstrate the ability of MGPs to correctly represent cell type specific changes in neurological conditions, we calculated dopaminergic profiles of substantia nigra samples in three expression data sets of Parkinson's disease (PD) patients and controls from Moran et al. (2006) (GSE8397), Lesnick et al. (2007) (GSE7621) and Zhang et al. (2005) (GSE20295). We tested whether the well-known loss of dopaminergic cells in PD could be detected using our MGP approach. MGP analysis correctly identified reduction in dopaminergic cells in substantia nigra of Parkinson's disease patients (Figure 5C). Further application of the MGPs allowed us to re-test a hypothesis raised by earlier work on these data. Moran et al. (2006) identified 22 genes consistently differentially expressed in PD patients in their study and in the study of Zhang et al. (2005), suggesting these genes can be considered as a "PD expression signature". We hypothesized that differential expression of some of these genes might be better explained by a decrease in dopaminergic cells rather than intracellular regulatory changes. We performed a principal component analysis of the expression profiles of the 22 genes, in the Moran, Lesnick and Zhang datasets. We then calculated the correlation between the first principal component (PC1) of the 22 genes and dopaminergic MGPs, separately for control and PD subjects. Our results show that in all datasets PC1 of the proposed PD signature genes is highly correlated with dopaminergic MGPs in both control (Moran: $\rho = 0.88$; Lesnick: $\rho = 0.95$; Zhang $\rho = 0.85$) and PD (Moran: $\rho = 0.52$; Lesnic: $\rho = 0.42$; Zhang $\rho = 0.430$) subjects (Figure 5D). Examination of individual expression patterns of the proposed signature genes revealed that in both datasets a majority of these genes is positively correlated to the dopaminergic MGPs, regardless of the disease state in both datasets (Figure 5E). These results suggest that majority of the genes identified by Moran et al. are more simply interpreted as reflecting changes in dopaminergic cell number rather than disease-induced transcriptional changes. This further emphasizes the need to account for cellular changes in gene expression analyses.

Discussion

Cell type specific expression database as a resource for neuroscience

We present NeuroExpresso, a rigorously curated database of brain cell type specific gene expression data (www.neuroexpresso.org), and demonstrate its utility in identifying cell-type markers and in the interpretation of bulk tissue expression profiles. To our knowledge, NeuroExpresso is the most comprehensive database of expression data for identified brain cell types. The database will be expanded as more data become available.

NeuroExpresso allows simultaneous examination of gene expression associated with numerous cell types across different brain regions. This approach promotes discovery of cellular properties that might have otherwise been unnoticed or overlooked when using gene-by-gene approaches or pathway enrichment analysis. For example, a simple examination of expression of genes involved in biosynthesis and secretion of GABA and glutamate, suggested the co-release of these neurotransmitters from forebrain and habenular cholinergic cells, respectively.

Studies that aim to identify novel properties of cell types can benefit from our database as an inexpensive and convenient way to seek novel patterns of gene expression. For instance, our database shows significant bimodality of gene expression in dopaminergic cell types from the midbrain (Figure 4C). The observed bimodality might indicate heterogeneity in the dopaminergic cell population, which could prove a fruitful avenue for future investigation. Another interesting finding from NeuroExpresso is the previously unknown overlap of several markers of motor cholinergic and noradrenergic cells. While the overlapping markers were previously shown to be expressed in spinal cholinergic cells, to our knowledge their expression in noradrenergic (as well as brain stem cholinergic) cells was previously unknown.

NeuroExpresso can be also used to facilitate interpretation of genomics and transcriptomics studies. Recently (Pantazatos et al., 2016) used an early release of the databases to interpret expression patterns in the cortex of suicide victims, suggesting involvement of microglia. Moreover,

this database has further applications beyond the use of marker genes, such as interpreting other features of brain cell diversity, like in electrophysiological profiles (Tripathy et al., submitted). Importantly, NeuroExpresso is a cross-laboratory database. A consistent result observed across several studies raises the certainty that it represents a true biological finding rather than merely an artefact or contamination with other cell types. This is specifically important for unexpected findings such as the expression of Ddc in oligodendrocytes (Figure 4B).

Validation of cell type markers

To assess the quality of the marker genes, a subset of our cell type markers were validated by in situ hybridization (Cox6a2 as a marker of fast spiking basket cells, and multiple Purkinje and DG cell markers). Further validation was performed with computational methods in independent single cell datasets from mouse and human. This analysis validated all cortical gene sets except Oligodendrocyte precursors (OP). In their paper, Zeisel et al. (2015) stated that none of the oligodendrocyte sub-clusters they identified were associated with oligodendrocyte precursor cells, which likely explains why we were not able to validate the OP MGP in their dataset. The Darmanis dataset however, is reported to include oligodendrocyte precursors (18/466 cells) (Darmanis et al., 2015), but again our OPC MGP did not show good validation. In this case the reason for negative results could be changes in the expression of the mouse marker gene orthologs in human, possibly reflecting functional differences between the human and mouse cell types (Shay et al., 2013; Zhang et al., 2016). Further work will be needed to identify a robust human OPC signature. However, since most MGSs did validate between mouse and human data, it suggests that most marker genes preserve their specificity despite cross-species gene expression differences.

Improving interpretation of bulk tissue expression profiles

Marker genes can assist with the interpretation of bulk tissue data in the form of marker gene profiles (MGPs). A parsimonious interpretation of a change in an MGP is a change in the relative abundance of the corresponding cell type. Similar summarizations of cell type specific is often used

to analyse gene expression (Chikina et al., 2015; Newman et al., 2015; Westra et al., 2015; Xu et al., 2013) and methylation data (Jones et al., 2017; Shannon et al., 2017). It remains possible that changes could have other explanations, but because our approach focuses on the overall trend of MGS expression levels, based on multiple markers, it should be relatively insensitive to within-cell-type expression changes for a subset of genes. Still, we prefer to refer the term “MGP expression” rather than “cell type proportions”, to emphasize the indirect nature of the approach.

Our results show that MGPs based on NeuroExpresso marker gene sets (MGSs) can reliably recapitulate relative changes in cell type abundance across different conditions. Direct validation of cell count estimation based on MGSs in human brain was not feasible due to the unavailability of cell counts coupled with expression data. Instead, we compared oligodendrocyte MGPs based on a gene expression dataset available through the SMRI database to experimental cell counts taken from a separate study (Uranova et al., 2004) of the same cohort of subjects and were able to recapitulate the reported reduction of oligodendrocyte proportions in patients with schizophrenia, bipolar disorder and depression. Based on analysis of dopaminergic MGPs we were able to capture the well-known reduction in dopaminergic cell types in PD patients. Moreover, we show that multiple genes previously reported as differentially expressed in PD are highly correlated with dopaminergic MGPs in both control and PD subjects. This high correlation suggests that the observed differential expression of these genes might merely reflect changes in dopaminergic cell populations rather than PD related regulatory changes

Limitations and caveats

While we took great care in the assembly of NeuroExpresso, there remain a number of limitations and room for improvement. First, the NeuroExpresso database was assembled from multiple datasets, based on different mouse strains and cell type extraction methodologies, which may lead to undesirable heterogeneity. We attempted to reduce inter-study variability by combined pre-processing of the raw data and normalization. However, due to insufficient overlap between cell types represented by different studies, many of the potential confounding factors such as age, sex

and methodology could not be explicitly corrected for. Thus, it is likely that some of the expression values in NeuroExpresso may be affected by confounding factors. While our confidence in the data is increased when expression signals are robust across multiple studies, many of the cell types in NeuroExpresso are represented by a single study. Hence, we advise that small differences in expression between cell types as well as previously unknown expression patterns based on a single data source should be treated with caution. In our analyses, we address these issues by enforcing a stringent set of criteria for the marker selection process, reducing the impact of outlier samples and ignoring small changes in gene expression.

An additional limitation of our study is that the representation for many of the brain cell types is still lacking in the NeuroExpresso database. Therefore, despite our considerable efforts to ensure cell type-specificity of the marker genes, we cannot rule out the possibility that some of them are also expressed in one or more of the non-represented cell types. This problem is partially alleviated in cortex due to the inclusion of single cell data. As more such datasets are available, it will be easier to create a more comprehensive database.

In summary, we believe that NeuroExpresso is a valuable resource for neuroscientists. We identified numerous novel markers for 36 major cell types and used them to estimate cell type profiles in bulk tissue data, demonstrating high correlation between our estimates and experiment-based cell counts. This approach can be used to reveal cell type specific changes in whole tissue samples and to re-evaluate previous analyses on brain whole tissues that might be biased by cell type-specific changes. Information about cell type-specific changes is likely to be very valuable since conditions like neuron death, inflammation, and astrogliosis are common hallmarks of in neurological diseases.

Materials and methods

Figure 1A depicts the workflow and the major steps of this study. All the analyses were performed in R version 3.3.2; the R code and data files can be accessed through neuroexpresso.org or

directly from <https://github.com/oganm/neuroexpressoAnalysis>.

Pooled cell type specific microarray data sets

We began with a collection of seven studies of isolated cell types from the brain, compiled by Okaty et al. (2011). We expanded this by querying PubMed (<http://www.ncbi.nlm.nih.gov/pubmed>) and Gene Expression Omnibus (GEO) (<http://www.ncbi.nlm.nih.gov/geo/>) (Barrett et al., 2013; Edgar et al., 2002) for cell type-specific expression datasets from the mouse brain that used Mouse Expression 430A Array (GPL339) or Mouse Genome 430 2.0 Array (GPL1261) platforms. These platforms were our focus as together, they are the most popular platforms for analysis of mouse samples and are relatively comprehensive in gene coverage, and using a reduced range of platforms reduced technical issues in combining studies. Query terms included names of specific cell types (e.g. astrocytes, pyramidal cells) along with blanket terms such as “brain cell expression” and “purified brain cells”. Only samples derived from postnatal (> 14 days), wild type, untreated animals were included. Datasets obtained from cell cultures or cell lines were excluded due to the reported expression differences between cultured cells and primary cells (Cahoy et al., 2008; Halliwell, 2003; Januszyk et al., 2015). We also considered RNA-seq data from pooled cells (2016; Zhang et al., 2014) but because such data sets are not available for many cell types, including it in the merged resource was not technically feasible without introducing biases (though we were able to incorporate a single-cell RNA-seq data set, described in the next section). While we plan to incorporate more pooled cell RNA-seq data in the future, for this study we limited their use to validation of marker selection.

As a first step in the quality control of the data, we manually validated that each sample expressed the gene that was used as a marker for purification (expression greater than median expression among all gene signals in the dataset), along with other well established marker genes for the relevant cell type (e.g. *Pcp2* for Purkinje cells, *Gad1* for GABAergic interneurons). We next excluded contaminated samples, namely, samples expressing established marker genes of non-related cell types in levels comparable to the cell type marker itself (for example neuronal samples

expressing high levels of glial marker genes), which lead to the removal of 21 samples. In total, we have 30 major cell types compiled from 24 studies represented by microarray data (summarized in Table 1); a complete list of all samples including those removed is available from the authors).

Single cell RNA-seq data

The study of cortical single cells by Tasic et al. (2016) includes a supplementary file (Supplementary Table 7 in Tasic et al. (2016)) linking a portion of the molecularly defined cell clusters to known cell types previously described in the literature. Using this file, we matched the cell clusters from Tasic et al. with pooled cortical cell types represented by microarray data (Table 3). For most cell types represented by microarray (e.g. glial cells, Martinotti cells), the matching was based on the correspondence information provided by Tasic et al. (2016). However, for some of the cell clusters from Tasic et al. (2016), the cell types were matched manually, based on the description of the cell type in the original publication (e.g., cortical layer, high expression of a specific gene). For example, Glt25d2⁺ pyramidal cells from Schmidt et al. (2012), described by the authors as “layer 5b pyramidal cells with high Glt25d2 and Fam84b expression” were matched with two cell clusters from Tasic et al. - “L5b Tph2” and “L5b Cdh13”, 2 of the 3 clusters described as Layer 5b glutamatergic cells by Tasic et al., since both of these clusters represented pyramidal cells from cortical layer 5b and exhibited high level of the indicated genes. Cell clusters identified in Tasic et al. that did not match to any of the pooled cell types were integrated into to the combined data if they fulfilled the following criteria: 1) They represented well-characterized cell types and 2) we could determine with high confidence that they did not correspond to more than one cell type represented by microarray data. Table 3 contains information regarding the matching between pooled cell types from microarray data and cell clusters from single cell RNA-seq data from Tasic et al.

In total, the combined database contains expression profiles for 36 major cell types, 10 of which are represented by both pooled cell microarray and single cell RNA-seq data, and five which are represented by single cell RNA-seq only (summarized in Table 3). Due to the substantial

differences between microarray and RNA-seq technologies, we analysed these data separately (see next sections). For visualization only, in neuroexpresso.org we rescaled the data to allow them to be plotted on the same axes. Details are provided on the web site.

Microarray cell type	Tasic et al. cell cluster	Matching method	NeuroExpresso cell type name
Astrocyte	Astro Gja1	Direct match	Astrocyte
Microglia	Micro Ctss	Direct match	Microglia
Oligodendrocyte	Oligo Opalin	Direct match	Oligodendrocyte
FS Basket (G42)	Pvalb Gpx3, Pvalb Rspo2, Pvalb Wt1, Pvalb Obox3, Pvalb Cpne5	Definition: fast spiking pval positive interneurons	FS Basket (G42)
Martinotti (GIN)	Sst Cbln4	Direct match	Martinotti (GIN)
VIPReIn (G30)	Vip Sncg	Unique Vip and Sncg expression, high Sncg expression in microarray cell type	VIPReIn (G30)
Pyramidal Glt25d2	L5b Tph2, L5b Cdh13	Definition: Glt25d2 positive Fam84b positive	Pyramidal Glt25d2
Pyramidal S100a10	L5a Hsd11b1, L5a Batf3, L5a Tcerg1l, L5a Pde1c	Definition: S100a10 expressing cells from layer 5a	Pyramidal S100a10
Pyramidal CrtThalamic	L6a Car12, L6a Syt17	Direct match	Pyramidal Crt-Thalamic
---	Endo Myl9, Endo Tbc1d4	New cell type	Endothelial
---	OPC Pdgfra	New cell type	Oligodendrocyte precursors
---	L4 Ctxn3, L4 Scnn1a, L4 Arf5	New cell type	Layer 4 Pyra
---	L2 Ngb, L2/3 Ptgs2	New cell type	Layer 2 3 Pyra
---	L6a Mgp, L6a Sla	New cell type	Layer 6a Pyra
---	L6b Serpinb11, L6b Rgs12	New cell type	Layer 6b Pyra

Table 3: Matching single cell RNA sequencing data from Tasic to well defined cell types. List of molecular cell types identified by Tasic et al. and their corresponding cell types in NeuroExpresso. Matching method column defines how the matching was performed. Direct matches are one to one matching between the definition provided by Tasic et al. for the molecular cell types and definition provided by microarray samples. For “Definition” matches, description of the cell type in the original source is used to find molecular cell types that fit the definition. VIPReIn – Vip Sncg matching was done based on unique Sncg expression in VIPReIn cells in the microarray data. New cell types are well defined cell types that have no counterpart in microarray data.

Grouping and re-assignment of cell type samples

When possible, samples were assigned to specific cell types based on the descriptions provided in their associated original publications. When expression profiles of closely related cell types were too similar to each other and we could not find sufficient number of differentiating marker genes

meeting our criteria, they were grouped together into a single cell type. For example, A10 and A9 dopaminergic cells had no distinguishing markers (provided the other cell types presented in the midbrain region) and were grouped as “dopaminergic neurons”. In the case of pyramidal cells, while we were able to detect marker genes for pyramidal cell subtypes, we found that they were often few in numbers and most of them were not present in the human microarray chip (Affymetrix Human Exon 1.0 ST Array) used in the downstream analysis. Hence, several pyramidal cell types were pooled into a “pan-pyramidal” set of markers for downstream analysis. Both the “pan-pyramidal” and the pyramidal sub-type marker genes are reported in Table 1.

Since our focus was on finding markers specific to cell types within a given brain region, samples were grouped based on the brain region from which they were isolated, guided by the anatomical hierarchy of brain regions (Figure 1B). Brain sub-regions (e.g. locus coeruleus) were added to the hierarchy if there were multiple cell types represented in the sub-region. An exception to the region assignment process are glial samples. Since these samples were only available from either cortex or cerebellum regions or extracted from whole brain, the following assignments were made: Cerebral cortex-derived astrocyte and oligodendrocyte samples were included in the analysis of other cerebral regions as well as thalamus, brainstem and spinal cord. Bergmann glia and cerebellum-derived oligodendrocytes were used in the analysis of cerebellum. The only microglia samples available were isolated from whole brain homogenates and were included in the analysis of all brain regions.

Selection of cell type markers

Marker gene sets (MGSs) were selected for each cell type in each brain region, based on fold change and clustering quality (see below). For cell types that are represented by both microarray and single cell data (cortical cells), two sets of MGSs were created and later merged as described below.

Marker genes were selected for each brain region based on the following steps:

1. For RNA-seq data, each of the relevant clusters identified in Tasic et al. was considered as a single sample, where the expression of each gene was calculated by taking the mean RPKM values of the individual cells representing the cluster. Table 3 shows which clusters represent which cell types.
2. Expression level of a gene in a cell type was calculated by taking the mean expression of all replicate samples originating from the same study and averaging the resulting values across different studies per cell type.
3. The quality of clustering was determined by “mean silhouette coefficient” and “minimal silhouette coefficient” values (where silhouette coefficient is a measure of group dissimilarity ranged between -1 and 1 (Rousseeuw, 1987)). Mean silhouette coefficient was calculated by assigning the samples representing the cell type of interest to one cluster and samples from the remaining cell types to another, and then calculating the mean silhouette coefficient of all samples. The minimal silhouette coefficient is the minimal value of mean silhouette coefficient when it is calculated for samples representing the cell type of interest in comparison to samples from each of the remaining cell types separately. The two measures were used to ensure that the marker gene robustly differentiates the cell type of interest from other cell types. Silhouette coefficients were calculated with the “silhouette” function from the “cluster” R package version 1.15.3 (Maechler et al., 2016), using the expression difference of the gene between samples as the distance metric.
4. A background value was selected below which the signal cannot be discerned from noise. Different background values are selected for microarray (6 – all values are log₂ transformed) and RNA-seq (0.1) due to the differences in their distribution.

Based on these metrics, the following criteria were used:

1. A threshold expression level was selected to help ensure that the gene’s transcripts will be detectable in bulk tissue. Genes with median expression level below this threshold were excluded from further analyses. For microarrays, this threshold was chosen to be 8. Theoretically, if a gene has an expression level of 8 in a cell type, and the gene is specific

to the cell type, an expression level of 6 would be observed if $1/8^{\text{th}}$ of a bulk tissue is composed of the cell type. As many of the cell types in the database are likely to be as rare as or rarer than $1/8^{\text{th}}$, and 6 is generally close to background for these data, we picked 8 as a lower level of marker gene expression. For RNA-seq data, we selected a threshold of 2.5 RPKM, which in terms of quantiles corresponds to the microarray level of 8.

2. If the expression level in the cell type of interest is higher than 10 times the background threshold, there must be at least a 10-fold difference from the median expression level of the remaining cell types in the region. If the expression level in the cell type is less than 10 times the background, the expression level must be higher than the expression level of every other cell type in that region. This criterion was added because below this expression level, for a 10-fold expression change to occur, the expression median of other cell types needs be lower than the background. Values below the background signal that do not convey meaningful information but can prevent potentially useful marker genes from being selected.

3. The mean silhouette coefficient for the gene must be higher than 0.5 and minimum silhouette coefficient must be greater than zero for the associated cell type.

4. The conditions above must be satisfied only by a single cell type in the region.

To ensure robustness against outlier samples, we used the following randomization procedure, repeated 500 times: One third (rounded) of all samples were removed. For microarray data, in order to prevent large studies from dominating the silhouette coefficient, when studies representing the same cell types did not have an equal number of samples, N samples were picked randomly from each of the studies, where N is the smallest number of samples coming from a single study.

Our next step was combining the MGSs created from the two expression data types. For cell types and genes represented by both microarray and RNA-seq data, we first looked at the intersection between the MGSs. For most of the cell types, the overlap between the two MGSs was about 50%. We reasoned that this could be partially due to numerous “near misses” in both data sources. Namely, since our method for marker gene selection relies on multiple steps with hard thresholds,

it is very likely that some genes were not selected simply because they were just below one of the required thresholds. We thus adopted a soft intersection: A gene was considered as a marker if it fulfilled the marker gene criteria in one data source (pooled cell microarray or single cell RNA-seq), and its expression in the corresponding cell type from the other data source was higher than in any other cell type in that region. For example, Ank1 was originally selected as a marker of FS Basket cells based on microarray data, but not based on RNA-seq data. Since the expression level of Ank1 in the RNA-seq data is higher in FS Basket cells than in any other cell type from this data source, based on the soft intersection criterion, Ank1 is still considered as a marker of FS Basket cells. For genes and cell types that were only represent by one data source, the selection was based on this data source only. Some previously described markers (such as Prox1 for dentate granule cells) are absent from our marker gene lists. This is due in some cases absence from the microarray platforms used, while others failed to meet our stringent selection criteria. Final marker gene lists are available at <http://www.chibi.ubc.ca/supplement-to-mancarci-et-al-neuroexpresso/>.

Human homologues of mouse genes were defined by NCBI HomoloGene (<ftp://ftp.ncbi.nih.gov/pub/HomoloGene/build68/homologene.data>).

Microglia enriched genes

Microglia expression profiles differ significantly between activated and inactivated states and to our knowledge, the samples in our database represent only the inactive state (Holtman et al., 2015). In order to acquire marker genes with stable expression levels regardless of microglia activation state, we removed the genes differentially expressed in activated microglia based on Holtman et al. (2015). This step resulted in removal of 408 out of the original 720 microglial genes in cortex (microarray and RNA-seq lists combined) and 253 of the 493 genes in the context of other brain regions (without genes from single cell data). Microglial marker genes which were differentially expressed in activated microglia are referred to as Microglia_activation and Microglia_deactivation (up or down-regulated, respectively) in the marker gene lists provided.

S100a10⁺ pyramidal cell enriched genes

The paper (Schmidt et al., 2012) describing the cortical S100a10⁺ pyramidal cells emphasizes the existence of non-neuronal cells expressing S100a10⁺. Schmidt et al. therefore limited their analysis to 7,853 genes specifically expressed in neurons and advised third-party users of the data to do so as well. Since a contamination caveat was only concerning microarray samples from Schmidt et al. (the only source of S100a10⁺ pyramidal cells in microarray data), we removed marker genes selected for S100a10⁺ pyramidal cells based on the microarray data if they were not among the 7,853 genes indicated in Schmidt et al. We also removed S100a10 itself since based on the author's description it was not specific to this cell type. In total, 36 of the 47 S100a10 pyramidal genes originally selected based on microarray data were removed in this step. Of note, none of the removed genes were selected as a marker of S100a10 cell based on RNA-seq data.

Dentate granule cell enriched genes

We used data from (Cembrowski et al., 2016) (Hipposeq) for validation and refinement of dentate granule markers (as noted above these data are not currently included in Neuroexpresso for technical reasons). FPKM values were downloaded (GEO accession GSE74985) and log₂ transformed. Based on these values, we removed marker genes we selected for granule cells if their expression in Hipposeq (mean of dorsal and ventral granule cells) was lower than other cell types in Hipposeq. In total, 15 of the 39 genes that were selected were removed in this step.

In situ hybridization

Male C57BL/6J mice aged 13-15 weeks at time of sacrifice were used (n=5). Mice were euthanized by cervical dislocation and then the brain was quickly removed, frozen on dry ice, and stored at -80°C until sectioned via cryostat. Brain sections containing the sensorimotor cortex were cut along the rostral-caudal axis using a block advance of 14 µm, immediately mounted on glass slides and dried at room temperature (RT) for 10 minutes, and then stored at -80°C until processed using multi-label fluorescent in situ hybridization procedures.

Fluorescent in situ hybridization probes were designed by Advanced Cell Diagnostics, Inc. (Hayward, CA, USA) to detect mRNA encoding Cox6a2, Slc32a1, and Pvalb. Two sections per animal were processed using the RNAscope® 2.5 Assay as previously described (Wang et al., 2012). Briefly, tissue sections were incubated in a protease treatment for 30 minutes at RT and then the probes were hybridized to their target mRNAs for 2 hours at 40°C. The sections were exposed to a series of incubations at 40°C that amplifies the target probes, and then counterstained with NeuroTrace blue-fluorescent Nissl stain (1:50; Molecular Probes) for 20 minutes at RT. Cox6a2, Pvalb, and Slc32a1 were detected with Alexa Fluor® 488, Atto 550 and Atto 647, respectively.

Data were collected on an Olympus IX83 inverted microscope equipped with a Hamamatsu Orca-Flash4.0 V2 digital CMOS camera using a 60x 1.40 NA SC oil immersion objective. The equipment was controlled by cellSens (Olympus). 3D image stacks (2D images successively captured at intervals separated by 0.25 µm in the z-dimension) that are 1434 x 1434 pixels (155.35 µm x 155.35 µm) were acquired over the entire thickness of the tissue section. The stacks were collected using optimal exposure settings (i.e., those that yielded the greatest dynamic range with no saturated pixels), with differences in exposures normalized before analyses.

Laminar boundaries of the sensorimotor cortex were determined by cytoarchitectonic criteria using NeuroTrace labeling. Fifteen image stacks across the gray matter area spanning from layer 2 to 6 were systematic randomly sampled using a sampling grid of 220 x 220 µm², which yielded a total of 30 image stacks per animal. Every NeuroTrace labeled neuron within a 700 x 700 pixels counting frame was included for analyses; the counting frame was placed in the center of each image to ensure that the entire NeuroTrace labeled neuron was in the field of view. The percentage (± standard deviation) of NeuroTrace labeled cells containing Cox6a2 mRNA (Cox6a2+) and that did not contain Slc32a1 mRNA (Slc32a1-), that contained Slc32a1 but not Pvalb mRNA (Slc32a1+/Pvalb-), and that contained both Slc32a1 and Pvalb mRNAs (Slc32a1+/Pvalb+) were manually assessed.

Allen Brain Atlas in situ hybridization (ISH) data

We downloaded in situ hybridization (ISH) images using the Allen Brain Atlas API (<http://help.brain-map.org/display/mousebrain/API>). Assessment of expression patterns was done by visual inspection. If a probe used in an ISH experiment did not show expression in the region, an alternative probe targeting the same gene was sought. If none of the probes showed expression in the region, the gene was considered to be not expressed.

Validation of marker genes using external single cell data

Mouse cortex single cell RNA sequencing (RNA-seq) data were acquired from Zeisel et al. (2015) (available from <http://linnarssonlab.org/cortex/>, GEO accession: GSE60361, 1691 cells) Human single cell RNA sequencing data were acquired from Darmanis et al. (2015) (GEO accession: GSE67835, 466 cells). For both datasets, pre-processed expression data were encoded in a binary matrix with 1 representing any nonzero value. For all marker gene sets, Spearman's ρ was used to quantify internal correlation. A null distribution was estimated by calculating the internal correlation of 1000 randomly-selected prevalence-matched gene groups. Gene prevalence was defined as the total number of cells with a non-zero expression value for the gene. Prevalence matching was done by choosing a random gene with a prevalence of $\pm 2.5\%$ of the prevalence of the marker gene. P-values were calculated by comparing the internal correlation of marker gene set to the internal correlations of random gene groups using Wilcoxon rank-sum test.

Pre-processing of microarray data

All microarray data used in the study were pre-processed and normalized with the "rma" function of the "oligo" (Affymetrix gene arrays) or "affy" (Affymetrix 3'IVT arrays) (Carvalho and Irizarry, 2010) R packages. Probeset to gene annotations were obtained from Gemma (Zoubarev et al., 2012) (<http://gemma.chibi.ubc.ca/>). Probesets with maximal expression level lower than the median among all probeset signals were removed. Of the remaining probesets, whenever several probesets were mapped to the same gene, the one with the highest variance among the samples

was selected for further analysis.

Samples from pooled cell types that make up the NeuroExpresso database were processed by an in-house modified version of the “rma” function that enabled collective processing of data from Mouse Expression 430A Array (GPL339) and Mouse Genome 430 2.0 Array (GPL1261) which share 22690 of their probesets. As part of the rma function, the samples are quantile normalized at the probe level. However, possibly due to differences in the purification steps used by different studies (Okaty et al., 2011), we still observed biases in signal distribution among samples originating from different studies. Thus, to increase the comparability across studies, we performed a second quantile normalization of the samples at a probeset level before selection of probes with the highest variance. After all processing the final data set included 11564 genes.

For comparison of marker gene profiles in white matter and frontal cortex, we acquired expression data from pathologically healthy brain samples from Trabzuni et al. (2013) (GEO accession: GSE60862). For estimation of dopaminergic marker gene profiles in Parkinson’s disease patients and controls, we acquired substantia nigra expression data from (Lesnick et al., 2007) (GSE7621) and (Moran et al., 2006) (GSE8397). Expression data for the Stanley Medical Research Institute (SMRI), which included post-mortem prefrontal cortex samples from bipolar disorder, major depression and schizophrenia patients along with healthy donors, were acquired through <https://www.stanleygenomics.org/>, study identifier 2.

Estimation of marker gene profiles (MGPs)

For each cell type relevant to the brain region analysed, we used the first principal component of the corresponding marker gene set expression as a surrogate for cell type proportions. This method of marker gene profile estimation is similar to the methodology of multiple previous works that aim to estimate relative abundance of cell types in a whole tissue sample (Chikina et al., 2015; Westra et al., 2015; Xu et al., 2013). Principal component analysis was performed using the “prcomp” function from the “stats” R package, using the “scale = TRUE” option. It is plausible that some marker genes will be transcriptionally differentially regulated under some conditions (e.g.

disease state), reducing the correspondence between their expression level with the relative cell proportion. A gene that is thus regulated is expected to have reduced correlation to the other marker genes with expression levels primarily dictated by cell type proportions, which will reduce their loading in the first principal component. To reduce the impact of regulated genes on the estimation process, we removed marker genes from a given analysis if their loadings had the opposite sign to the majority of markers when calculated based on all samples in the dataset and recalculate loadings and components using the remaining genes. This was repeated until all remaining genes had loadings with the same signs. Since the sign of the loadings of the rotation matrix (as produced by prcomp function) is arbitrary, to ease interpretation between the scores and the direction of summarized change in the expression of the relevant genes, we multiplied the scores by -1 whenever the sign of the loadings was negative. For visualization purposes, the scores were normalized to the range 0-1. Two sided Wilcoxon rank-sum test (“wilcox.test” function from the “stats” package in R, default options) was used to compare between the different experimental conditions.

For estimations of cell type MGPs in samples from frontal cortex and white matter from the Trabzuni study (Trabzuni et al., 2013), results were subjected to multiple testing correction by the Benjamini & Hochberg method (Benjamini and Hochberg, 1995). For the Parkinson’s disease datasets from Moran et al. (2006) and Lesnick et al. (2007), we estimated MGPs for dopaminergic neuron markers in control and PD subjects. Moran et al. data included samples from two sub-regions of substantia nigra. Since some of the subjects were sampled in only one of the sub-regions while others in both, the two sub-regions were analysed separately. The list of 22 PD genes, described as the “first PD expression signature”, were taken from Moran et al. (2006). Dopaminergic MGPs were calculated for the samples. Correlations of the PD gene expression as well as correlation of 1st principal component of PD gene expression to dopaminergic MGP was calculated using Spearman’s ρ .

For the SMRI collection of psychiatric patients we estimated oligodendrocytes MGPs based on expression data available through the SMRI website (as indicated above) and compared our

results to experimental cell counts from the same cohort of subjects previously reported by Uranova et al. (2004). Figure 5B representing the oligodendrocyte cell counts in each disease group was adapted from Uranova et al. (2004). The data presented in the figure was extracted from Figure 1A in Uranova et al. (2004) using WebPlotDigitizer (<http://arohatgi.info/WebPlotDigitizer/app/>).

Acknowledgements

This work is supported by a NeuroDevNet grant to PP, the UBC bioinformatics graduate training program (BOM), a CIHR post-doctoral fellowship to ST, by the Campbell Family Mental Health Research Institute of CAMH (ES and BR), NIH grants MH077159 to ES, and MH111099 and GM076990 to PP, and an NSERC Discovery Grant to PP. We thank Ken Sugino for providing access to raw CEL files for Purkinje and TH⁺ cells from locus coeruleus, Chee Yeun Chung for providing access to raw CEL files for dopaminergic cells, Dean Attali for providing insight on the usage of the Shiny platform, the Pavlidis lab for their inputs to the project during its development and Rosemary McCloskey for aid in editing the manuscript.

The authors declare no competing financial interests

References

- Anandasabapathy, N., Victora, G.D., Meredith, M., Feder, R., Dong, B., Kluger, C., Yao, K., Dustin, M.L., Nussenzweig, M.C., Steinman, R.M., Liu, K., 2011. Flt3L controls the development of radiosensitive dendritic cells in the meninges and choroid plexus of the steady-state mouse brain. *J. Exp. Med.* 208, 1695–1705. doi:10.1084/jem.20102657
- Barrett, T., Wilhite, S.E., Ledoux, P., Evangelista, C., Kim, I.F., Tomashevsky, M., Marshall, K.A., Phillippy, K.H., Sherman, P.M., Holko, M., Yefanov, A., Lee, H., Zhang, N., Robertson, C.L., Serova, N., Davis, S., Soboleva, A., 2013. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res.* 41, D991–D995. doi:10.1093/nar/gks1193
- Beckervordersandforth, R., Tripathi, P., Ninkovic, J., Bayam, E., Lepier, A., Stempfhuber, B., Kirchhoff, F., Hirrlinger, J., Haslinger, A., Lie, D.C., Beckers, J., Yoder, B., Irmeler, M., Götz, M., 2010. In Vivo Fate Mapping and Expression Analysis Reveals Molecular Hallmarks of Prospectively Isolated Adult Neural Stem Cells. *Cell Stem Cell* 7, 744–758. doi:10.1016/j.stem.2010.11.017
- Bellesi, M., Pfister-Genskow, M., Maret, S., Keles, S., Tononi, G., Cirelli, C., 2013. Effects of Sleep and Wake on Oligodendrocytes and Their Precursors. *J. Neurosci.* 33, 14288–14300. doi:10.1523/JNEUROSCI.5102-12.2013
- Benjamini, Y., Hochberg, Y., 1995. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* 57, 289–300.
- Bennett, M.L., Bennett, F.C., Liddel, S.A., Ajami, B., Zamanian, J.L., Fernhoff, N.B., Mulinyawe, S.B., Bohlen, C.J., Adil, A., Tucker, A., Weissman, I.L., Chang, E.F., Li, G., Grant, G.A., Hayden Gephart, M.G., Barres, B.A., 2016. New tools for studying microglia in the mouse and human CNS. *Proc. Natl. Acad. Sci. U. S. A.* 113, E1738-1746. doi:10.1073/pnas.1525528113
- Bowling, K., Ramaker, R.C., Lasseigne, B.N., Hagenauer, M., Hardigan, A., Davis, N., Gertz, J., Cartagena, P., Walsh, D., Vawter, M., Schatzberg, A., Barchas, J., Watson, S., Bunney, B., Akil, H., Bunney, W., Li, J., Cooper, S., Myers, R.M., 2016. Post-mortem molecular profiling of three psychiatric disorders reveals widespread dysregulation of cell-type associated transcripts and refined disease-related transcription changes. *bioRxiv* 061416. doi:10.1101/061416
- Cahoy, J.D., Emery, B., Kaushal, A., Foo, L.C., Zamanian, J.L., Christopherson, K.S., Xing, Y., Lubischer, J.L., Krieg, P.A., Krupenko, S.A., Thompson, W.J., Barres, B.A., 2008. A Transcriptome Database for Astrocytes, Neurons, and Oligodendrocytes: A New Resource for Understanding Brain Development and Function. *J. Neurosci.* 28, 264–278. doi:10.1523/JNEUROSCI.4178-07.2008
- Carvalho, B.S., Irizarry, R.A., 2010. A framework for oligonucleotide microarray preprocessing. *Bioinforma. Oxf. Engl.* 26, 2363–2367. doi:10.1093/bioinformatics/btq431
- Celio, M.R., Heizmann, C.W., 1981. Calcium-binding protein parvalbumin as a neuronal marker. *Nature* 293, 300–302. doi:10.1038/293300a0
- Cembrowski, M.S., Wang, L., Sugino, K., Shields, B.C., Spruston, N., 2016. Hipposeq: a comprehensive RNA-seq database of gene expression in hippocampal principal neurons. *eLife* 5, e14997. doi:10.7554/eLife.14997

767 Chikina, M., Zaslavsky, E., Sealfon, S.C., 2015. CellCODE: A robust latent variable approach to
768 differential expression analysis for heterogeneous cell populations. *Bioinformatics* *btv015*.
769 doi:10.1093/bioinformatics/btv015

770 Chung, C.Y., Seo, H., Sonntag, K.C., Brooks, A., Lin, L., Isacson, O., 2005. Cell type-specific gene
771 expression of midbrain dopaminergic neurons reveals molecules involved in their
772 vulnerability and protection. *Hum. Mol. Genet.* *14*, 1709–1725. doi:10.1093/hmg/ddi178

773 Dalal, J., Roh, J.H., Maloney, S.E., Akuffo, A., Shah, S., Yuan, H., Wamsley, B., Jones, W.B.,
774 Strong, C. de G., Gray, P.A., Holtzman, D.M., Heintz, N., Dougherty, J.D., 2013.
775 Translational profiling of hypocretin neurons identifies candidate molecules for sleep
776 regulation. *Genes Dev.* *27*, 565–578. doi:10.1101/gad.207654.112

777 Darmanis, S., Sloan, S.A., Zhang, Y., Enge, M., Caneda, C., Shuer, L.M., Gephart, M.G.H., Barres,
778 B.A., Quake, S.R., 2015. A survey of human brain transcriptome diversity at the single cell
779 level. *Proc. Natl. Acad. Sci.* *112*, 7285–7290. doi:10.1073/pnas.1507125112

780 Dougherty, J.D., Maloney, S.E., Wozniak, D.F., Rieger, M.A., Sonnenblick, L., Coppola, G., Mahieu,
781 N.G., Zhang, J., Cai, J., Patti, G.J., Abrahams, B.S., Geschwind, D.H., Heintz, N., 2013.
782 The Disruption of *Celf6*, a Gene Identified by Translational Profiling of Serotonergic
783 Neurons, Results in Autism-Related Behaviors. *J. Neurosci.* *33*, 2732–2753.
784 doi:10.1523/JNEUROSCI.4762-12.2013

785 Doyle, J.P., Dougherty, J.D., Heiman, M., Schmidt, E.F., Stevens, T.R., Ma, G., Bupp, S., Shrestha,
786 P., Shah, R.D., Doughty, M.L., Gong, S., Greengard, P., Heintz, N., 2008. Application of a
787 Translational Profiling Approach for the Comparative Analysis of CNS Cell Types. *Cell* *135*,
788 749–762. doi:10.1016/j.cell.2008.10.029

789 Edgar, R., Domrachev, M., Lash, A.E., 2002. Gene Expression Omnibus: NCBI gene expression
790 and hybridization array data repository. *Nucleic Acids Res.* *30*, 207–210.
791 doi:10.1093/nar/30.1.207

792 Enjin, A., Rabe, N., Nakanishi, S.T., Vallstedt, A., Gezelius, H., Memic, F., Lind, M., Hjalt, T.,
793 Tourtellotte, W.G., Bruder, C., Eichele, G., Whelan, P.J., Kullander, K., 2010. Identification
794 of novel spinal cholinergic genetic subtypes disclose *Chodl* and *Pitx2* as markers for fast
795 motor neurons and partition cells. *J. Comp. Neurol.* *518*, 2284–2304.
796 doi:10.1002/cne.22332

797 Erny, D., Hrabě de Angelis, A.L., Jaitin, D., Wieghofer, P., Staszewski, O., David, E., Keren-Shaul,
798 H., Mhlahkoi, T., Jakobshagen, K., Buch, T., Schwierzeck, V., Utermöhlen, O., Chun, E.,
799 Garrett, W.S., McCoy, K.D., Diefenbach, A., Staeheli, P., Stecher, B., Amit, I., Prinz, M.,
800 2015. Host microbiota constantly control maturation and function of microglia in the CNS.
801 *Nat. Neurosci.* *18*, 965–977. doi:10.1038/nn.4030

802 Fomchenko, E.I., Dougherty, J.D., Helmy, K.Y., Katz, A.M., Pietras, A., Brennan, C., Huse, J.T.,
803 Milosevic, A., Holland, E.C., 2011. Recruited Cells Can Become Transformed and Overtake
804 PDGF-Induced Murine Gliomas In Vivo during Tumor Progression. *PLoS ONE* *6*, e20605.
805 doi:10.1371/journal.pone.0020605

806 Galloway, J.N., Shaw, C., Yu, P., Parghi, D., Poidevin, M., Jin, P., Nelson, D.L., 2014. CGG repeats
807 in RNA modulate expression of TDP-43 in mouse and fly models of fragile X tremor ataxia
808 syndrome. *Hum. Mol. Genet.* *ddu314*. doi:10.1093/hmg/ddu314

809 Görlich, A., Antolin-Fontes, B., Ables, J.L., Frahm, S., Ślimak, M.A., Dougherty, J.D., Ibañez-Tallon,
810 I., 2013. Reexposure to nicotine during withdrawal increases the pacemaking activity of

811 cholinergic habenular neurons. *Proc. Natl. Acad. Sci.* 110, 17077–17082.
812 doi:10.1073/pnas.1313103110

813 Gudi, V., Moharregh-Khiabani, D., Skripuletz, T., Koutsoudaki, P.N., Kotsiari, A., Skuljec, J., Trebst,
814 C., Stangel, M., 2009. Regional differences between grey and white matter in cuprizone
815 induced demyelination. *Brain Res.* 1283, 127–138. doi:10.1016/j.brainres.2009.06.005

816 Hagenauer, M.H., Li, J.Z., Walsh, D.M., Vawter, M.P., Thompson, R.C., Turner, C.A., Bunney, W.E.,
817 Myers, R.M., Barchas, J.D., Schatzberg, A.F., Watson, S.J., Akil, H., 2016. Inference of cell
818 type composition from human brain transcriptomic datasets illuminates the effects of age,
819 manner death, dissection, and psychiatric diagnosis. *bioRxiv* 089391. doi:10.1101/089391

820 Halliwell, B., 2003. Oxidative stress in cell culture: an under-appreciated problem? *FEBS Lett.* 540,
821 3–6. doi:10.1016/S0014-5793(03)00235-7

822 Handley, A., Schauer, T., Ladurner, A.G., Margulies, C.E., 2015. Designing Cell-Type-Specific
823 Genome-wide Experiments. *Mol. Cell* 58, 621–631. doi:10.1016/j.molcel.2015.04.024

824 Heiman, M., Heilbut, A., Francardo, V., Kulicke, R., Fenster, R.J., Kolaczyk, E.D., Mesirov, J.P.,
825 Surmeier, D.J., Cenci, M.A., Greengard, P., 2014. Molecular adaptations of striatal spiny
826 projection neurons during levodopa-induced dyskinesia. *Proc. Natl. Acad. Sci.* 111, 4578–
827 4583. doi:10.1073/pnas.1401819111

828 Holtman, I.R., Noback, M., Bijlsma, M., Duong, K.N., van der Geest, M.A., Ketelaars, P.T., Brouwer,
829 N., Vainchtein, I.D., Eggen, B.J.L., Boddeke, H.W.G.M., 2015. Glia Open Access Database
830 (GOAD): A comprehensive gene expression encyclopedia of glia cells in health and
831 disease. *Glia* n/a-n/a. doi:10.1002/glia.22810

832 Hu, H., Gan, J., Jonas, P., 2014. Fast-spiking, parvalbumin+ GABAergic interneurons: From
833 cellular design to microcircuit function. *Science* 345, 1255263.
834 doi:10.1126/science.1255263

835 Januszyk, M., Rennert, R.C., Sorkin, M., Maan, Z.N., Wong, L.K., Whittam, A.J., Whitmore, A.,
836 Duscher, D., Gurtner, G.C., 2015. Evaluating the Effect of Cell Culture on Gene Expression
837 in Primary Tissue Samples Using Microfluidic-Based Single Cell Transcriptional Analysis.
838 *Microarrays* 4, 540–550. doi:10.3390/microarrays4040540

839 Jones, M.J., Islam, S.A., Edgar, R.D., Kobor, M.S., 2017. Adjusting for Cell Type Composition in
840 DNA Methylation Data Using a Regression-Based Approach. *Methods Mol. Biol.* Clifton NJ
841 1589, 99–106. doi:10.1007/7651_2015_262

842 Kawaguchi, Y., Katsumaru, H., Kosaka, T., Heizmann, C.W., Hama, K., 1987. Fast spiking cells in
843 rat hippocampus (CA1 region) contain the calcium-binding protein parvalbumin. *Brain Res.*
844 416, 369–374. doi:10.1016/0006-8993(87)90921-8

845 Kuhn, A., Thu, D., Waldvogel, H.J., Faull, R.L.M., Luthi-Carter, R., 2011. Population-specific
846 expression analysis (PSEA) reveals molecular changes in diseased brain. *Nat. Methods* 8,
847 945–947. doi:10.1038/nmeth.1710

848 Lesnick, T.G., Papapetropoulos, S., Mash, D.C., French-Mullen, J., Shehadeh, L., de Andrade, M.,
849 Henley, J.R., Rocca, W.A., Ahlskog, J.E., Maraganore, D.M., 2007. A genomic pathway
850 approach to a complex disease: axon guidance and Parkinson disease. *PLoS Genet.* 3,
851 e98. doi:10.1371/journal.pgen.0030098

852 Lobo, M.K., Karsten, S.L., Gray, M., Geschwind, D.H., Yang, X.W., 2006. FACS-array profiling of

853 striatal projection neuron subtypes in juvenile and adult mouse brains. *Nat. Neurosci.* 9,
854 443–452. doi:10.1038/nn1654

855 Maechler, M., original), P.R. (Fortran, original), A.S. (S, original), M.H. (S, maintenance(1999-
856 2000)), K.H. (port to R., Studer, M., Roudier, P., 2016. cluster: “Finding Groups in Data”:
857 Cluster Analysis Extended Rousseeuw et al.

858 Margolis, E.B., Lock, H., Hjelmstad, G.O., Fields, H.L., 2006. The ventral tegmental area revisited:
859 is there an electrophysiological marker for dopaminergic neurons? *J. Physiol.* 577, 907–
860 924. doi:10.1113/jphysiol.2006.117069

861 Maze, I., Chaudhury, D., Dietz, D.M., Von Schimmelmann, M., Kennedy, P.J., Lobo, M.K., Sullivan,
862 S.E., Miller, M.L., Bagot, R.C., Sun, H., Turecki, G., Neve, R.L., Hurd, Y.L., Shen, L., Han,
863 M.-H., Schaefer, A., Nestler, E.J., 2014. G9a influences neuronal subtype specification in
864 striatum. *Nat. Neurosci.* 17, 533–539. doi:10.1038/nn.3670

865 Moran, L.B., Duke, D.C., Deprez, M., Dexter, D.T., Pearce, R.K.B., Graeber, M.B., 2006. Whole
866 genome expression profiling of the medial and lateral substantia nigra in Parkinson's
867 disease. *Neurogenetics* 7, 1–11. doi:10.1007/s10048-005-0020-2

868 Newman, A.M., Liu, C.L., Green, M.R., Gentles, A.J., Feng, W., Xu, Y., Hoang, C.D., Diehn, M.,
869 Alizadeh, A.A., 2015. Robust enumeration of cell subsets from tissue expression profiles.
870 *Nat. Methods* 12, 453–457. doi:10.1038/nmeth.3337

871 Ogura, K., Ogawa, M., Yoshida, M., 1994. Effects of ageing on microglia in the normal rat brain:
872 immunohistochemical observations. *Neuroreport* 5, 1224–1226.

873 Okaty, B.W., Miller, M.N., Sugino, K., Hempel, C.M., Nelson, S.B., 2009. Transcriptional and
874 electrophysiological maturation of neocortical fastspiking GABAergic interneurons. *J.*
875 *Neurosci. Off. J. Soc. Neurosci.* 29, 7040–7052. doi:10.1523/JNEUROSCI.0105-09.2009

876 Okaty, B.W., Sugino, K., Nelson, S.B., 2011. A Quantitative Comparison of Cell-Type-Specific
877 Microarray Gene Expression Profiling Methods in the Mouse Brain. *PLoS ONE* 6, e16493.
878 doi:10.1371/journal.pone.0016493

879 Pantazatos, S.P., Huang, Y.-Y., Rosoklija, G.B., Dwork, A.J., Arango, V., Mann, J.J., 2016. Whole-
880 transcriptome brain expression and exon-usage profiling in major depression and suicide:
881 evidence for altered glial, endothelial and ATPase activity. *Mol. Psychiatry*.
882 doi:10.1038/mp.2016.130

883 Pattyn, A., Morin, X., Cremer, H., Goridis, C., Brunet, J.F., 1997. Expression and interactions of the
884 two closely related homeobox genes *Phox2a* and *Phox2b* during neurogenesis. *Dev.*
885 *Camb. Engl.* 124, 4065–4075.

886 Paul, A., Cai, Y., Atwal, G.S., Huang, Z.J., 2012. Developmental coordination of gene expression
887 between synaptic partners during GABAergic circuit assembly in cerebellar cortex. *Front.*
888 *Neural Circuits* 6, 37. doi:10.3389/fncir.2012.00037

889 Perrone-Bizzozero, N.I., Tanner, D.C., Mounce, J., Bolognani, F., 2011. Increased Expression of
890 Axogenesis-Related Genes and Mossy Fibre Length in Dentate Granule Cells from Adult
891 HuD Overexpressor Mice. *ASN Neuro* 3, AN20110015. doi:10.1042/AN20110015

892 Phani, S., Gonye, G., Iacovitti, L., 2010. VTA neurons show a potentially protective transcriptional
893 response to MPTP. *Brain Res.* 1343, 1–13. doi:10.1016/j.brainres.2010.04.061

- 894 Pickel, V.M., Joh, T.H., Reis, D.J., 1976. Monoamine-synthesizing enzymes in central
895 dopaminergic, noradrenergic and serotonergic neurons. Immunocytochemical localization
896 by light and electron microscopy. *J. Histochem. Cytochem.* 24, 792–792.
897 doi:10.1177/24.7.8567
- 898 Poulin, J.-F., Tasic, B., Hjerling-Leffler, J., Trimarchi, J.M., Awatramani, R., 2016. Disentangling
899 neural cell diversity using single-cell transcriptomics. *Nat. Neurosci.* 19, 1131–1141.
900 doi:10.1038/nn.4366
- 901 Ren, J., Qin, C., Hu, F., Tan, J., Qiu, L., Zhao, S., Feng, G., Luo, M., 2011. Habenula “Cholinergic”
902 Neurons Corelease Glutamate and Acetylcholine and Activate Postsynaptic Neurons via
903 Distinct Transmission Modes. *Neuron* 69, 445–452. doi:10.1016/j.neuron.2010.12.038
- 904 Ren, L., Wienecke, J., Hultborn, H., Zhang, M., 2016. Production of dopamine by aromatic L-amino
905 acid decarboxylase cells after spinal cord injury. *J. Neurotrauma.*
906 doi:10.1089/neu.2015.4037
- 907 Rong, Y., Wang, T., Morgan, J.I., 2004. Identification of candidate Purkinje cell-specific markers by
908 gene expression profiling in wild-type and *pcd(3J)* mice. *Brain Res. Mol. Brain Res.* 132,
909 128–145. doi:10.1016/j.molbrainres.2004.10.015
- 910 Rossner, M.J., Hirrlinger, J., Wichert, S.P., Boehm, C., Newrzella, D., Hiemisch, H., Eisenhardt, G.,
911 Stuenkel, C., Ahsen, O. von, Nave, K.-A., 2006. Global Transcriptome Analysis of
912 Genetically Identified Neurons in the Adult Cortex. *J. Neurosci.* 26, 9956–9966.
913 doi:10.1523/JNEUROSCI.0468-06.2006
- 914 Rousseeuw, P.J., 1987. Silhouettes: A graphical aid to the interpretation and validation of cluster
915 analysis. *J. Comput. Appl. Math.* 20, 53–65. doi:10.1016/0377-0427(87)90125-7
- 916 Satoh, J.-I., Kino, Y., Asahina, N., Takitani, M., Miyoshi, J., Ishida, T., Saito, Y., 2016. TMEM119
917 marks a subset of microglia in the human brain. *Neuropathol. Off. J. Jpn. Soc. Neuropathol.*
918 36, 39–49. doi:10.1111/neup.12235
- 919 Saunders, A., Granger, A.J., Sabatini, B.L., 2015. Corelease of acetylcholine and GABA from
920 cholinergic forebrain neurons. *eLife* 4. doi:10.7554/eLife.06412
- 921 Schmidt, E.F., Warner-Schmidt, J.L., Otopalik, B.G., Pickett, S.B., Greengard, P., Heintz, N., 2012.
922 Identification of the Cortical Neurons that Mediate Antidepressant Responses. *Cell* 149,
923 1152–1163. doi:10.1016/j.cell.2012.03.038
- 924 Shannon, C.P., Balshaw, R., Chen, V., Hollander, Z., Toma, M., McManus, B.M., FitzGerald, J.M.,
925 Sin, D.D., Ng, R.T., Tebbutt, S.J., 2017. Enumerateblood – an R package to estimate the
926 cellular composition of whole blood from Affymetrix Gene ST gene expression profiles.
927 *BMC Genomics* 18. doi:10.1186/s12864-016-3460-1
- 928 Shay, T., Jojic, V., Zuk, O., Rothamel, K., Puyraimond-Zemmour, D., Feng, T., Wakamatsu, E.,
929 Benoist, C., Koller, D., Regev, A., ImmGen Consortium, 2013. Conservation and divergence
930 in the transcriptional programs of the human and mouse immune systems. *Proc. Natl.*
931 *Acad. Sci. U. S. A.* 110, 2946–2951. doi:10.1073/pnas.1222738110
- 932 Shrestha, P., Mousa, A., Heintz, N., 2015. Layer 2/3 pyramidal cells in the medial prefrontal cortex
933 moderate stress induced depressive behaviors. *eLife* 4. doi:10.7554/eLife.08752
- 934 Sibille, E., Arango, V., Joeyen-Waldorf, J., Wang, Y., Leman, S., Surget, A., Belzung, C., Mann,
935 J.J., Lewis, D.A., 2008. Large-scale estimates of cellular origins of mRNAs: enhancing the

936 yield of transcriptome analyses. *J. Neurosci. Methods* 167, 198–206.
937 doi:10.1016/j.jneumeth.2007.08.009

938 Skene, N.G., Grant, S.G.N., 2016. Identification of Vulnerable Cell Types in Major Brain Disorders
939 Using Single Cell Transcriptomes and Expression Weighted Cell Type Enrichment. *Front.*
940 *Neurosci.* 10. doi:10.3389/fnins.2016.00016

941 Sugino, K., Hempel, C.M., Miller, M.N., Hattox, A.M., Shapiro, P., Wu, C., Huang, Z.J., Nelson,
942 S.B., 2006. Molecular taxonomy of major neuronal classes in the adult mouse forebrain.
943 *Nat. Neurosci.* 9, 99–107. doi:10.1038/nn1618

944 Sugino, K., Hempel, C.M., Okaty, B.W., Arnson, H.A., Kato, S., Dani, V.S., Nelson, S.B., 2014. Cell-
945 Type-Specific Repression by Methyl-CpG-Binding Protein 2 Is Biased toward Long Genes.
946 *J. Neurosci.* 34, 12877–12883. doi:10.1523/JNEUROSCI.2674-14.2014

947 Sunkin, S.M., Ng, L., Lau, C., Dolbeare, T., Gilbert, T.L., Thompson, C.L., Hawrylycz, M., Dang, C.,
948 2013. Allen Brain Atlas: an integrated spatio-temporal portal for exploring the central
949 nervous system. *Nucleic Acids Res.* 41, D996–D1008. doi:10.1093/nar/gks1042

950 Tan, C.L., Plotkin, J.L., Venø, M.T., Schimmelmann, M. von, Feinberg, P., Mann, S., Handler, A.,
951 Kjems, J., Surmeier, D.J., O'Carroll, D., Greengard, P., Schaefer, A., 2013. MicroRNA-128
952 governs neuronal excitability and motor behavior in mice. *Science* 342, 1254–1258.
953 doi:10.1126/science.1244193

954 Tan, P.P.C., French, L., Pavlidis, P., 2013. Neuron-Enriched Gene Expression Patterns are
955 Regionally Anti-Correlated with Oligodendrocyte-Enriched Patterns in the Adult Mouse and
956 Human Brain. *Front. Neurosci.* 7, 5. doi:10.3389/fnins.2013.00005

957 Tasic, B., Menon, V., Nguyen, T.N., Kim, T.K., Jarsky, T., Yao, Z., Levi, B., Gray, L.T., Sorensen,
958 S.A., Dolbeare, T., Bertagnolli, D., Goldy, J., Shapovalova, N., Parry, S., Lee, C., Smith, K.,
959 Bernard, A., Madisen, L., Sunkin, S.M., Hawrylycz, M., Koch, C., Zeng, H., 2016. Adult
960 mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nat. Neurosci.* 19,
961 335–346. doi:10.1038/nn.4216

962 Tomomura, M., Rice, D.S., Morgan, J.I., Yuzaki, M., 2001. Purification of Purkinje cells by
963 fluorescence-activated cell sorting from transgenic mice that express green fluorescent
964 protein. *Eur. J. Neurosci.* 14, 57–63.

965 Trabzuni, D., Ramasamy, A., Imran, S., Walker, R., Smith, C., Weale, M.E., Hardy, J., Ryten, M.,
966 North American Brain Expression Consortium, 2013. Widespread sex differences in gene
967 expression and splicing in the adult human brain. *Nat. Commun.* 4.
968 doi:10.1038/ncomms3771

969 Ugrumov, M.V., 2013. Chapter Four - Brain Neurons Partly Expressing Dopaminergic Phenotype:
970 Location, Development, Functional Significance, and Regulation, in: Eiden, L.E. (Ed.),
971 *Advances in Pharmacology, A New Era of Catecholamines in the Laboratory and Clinic.*
972 Academic Press, pp. 37–91.

973 Uranova, N.A., Vostrikov, V.M., Orlovskaya, D.D., Rachmanova, V.I., 2004. Oligodendroglial
974 density in the prefrontal cortex in schizophrenia and mood disorders: a study from the
975 Stanley Neuropathology Consortium. *Schizophr. Res.* 67, 269–275. doi:10.1016/S0920-
976 9964(03)00181-6

977 Wang, Y., Winters, J., Subramaniam, S., 2012. Functional classification of skeletal muscle
978 networks. II. Applications to pathophysiology. *J. Appl. Physiol. Bethesda Md* 1985 113,

979 1902–1920. doi:10.1152/japplphysiol.01515.2011

980 Westra, H.-J., Arends, D., Esko, T., Peters, M.J., Schurmann, C., Schramm, K., Kettunen, J.,
981 Yaghootkar, H., Fairfax, B.P., Andiappan, A.K., Li, Y., Fu, J., Karjalainen, J., Platteel, M.,
982 Visschedijk, M., Weersma, R.K., Kasela, S., Milani, L., Tserel, L., Peterson, P., Reinmaa, E.,
983 Hofman, A., Uitterlinden, A.G., Rivadeneira, F., Homuth, G., Petersmann, A., Lohrbeier, R.,
984 Prokisch, H., Meitinger, T., Herder, C., Roden, M., Grallert, H., Ripatti, S., Perola, M., Wood, A.R.,
985 Melzer, D., Ferrucci, L., Singleton, A.B., Hernandez, D.G., Knight, J.C., Melchior, R.,
986 Lee, B., Poidinger, M., Zozzoli, F., Larbi, A., Wang, D.Y., van den Berg, L.H., Veldink, J.H.,
987 Rotzschke, O., Makino, S., Salomaa, V., Strauch, K., Völker, U., van Meurs, J.B.J.,
988 Metspalu, A., Wijmenga, C., Jansen, R.C., Franke, L., 2015. Cell Specific eQTL Analysis
989 without Sorting Cells. *PLoS Genet* 11, e1005223. doi:10.1371/journal.pgen.1005223

990 Williams, M.R., Hampton, T., Pearce, R.K.B., Hirsch, S.R., Ansorge, O., Thom, M., Maier, M., 2013.
991 Astrocyte decrease in the subgenual cingulate and callosal genu in schizophrenia. *Eur.*
992 *Arch. Psychiatry Clin. Neurosci.* 263, 41–52. doi:10.1007/s00406-012-0328-5

993 Xu, X., Nehorai, A., Dougherty, J.D., 2013. Cell type-specific analysis of human brain transcriptome
994 data to predict alterations in cellular composition. *Syst. Biomed.* 1, 151–160.
995 doi:10.4161/sysb.25630

996 Zamanian, J.L., Xu, L., Foo, L.C., Nouri, N., Zhou, L., Giffard, R.G., Barres, B.A., 2012. Genomic
997 Analysis of Reactive Astrogliosis. *J. Neurosci.* 32, 6391–6410.
998 doi:10.1523/JNEUROSCI.6221-11.2012

999 Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., Manno, G.L., Juréus, A.,
1000 Marques, S., Munguba, H., He, L., Betsholtz, C., Rolny, C., Castelo-Branco, G., Hjerling-
1001 Leffler, J., Linnarsson, S., 2015. Cell types in the mouse cortex and hippocampus revealed
1002 by single-cell RNA-seq. *Science* 347, 1138–1142. doi:10.1126/science.aaa1934

1003 Zhang, Y., Chen, K., Sloan, S.A., Bennett, M.L., Scholze, A.R., O’Keeffe, S., Phatnani, H.P.,
1004 Guarnieri, P., Caneda, C., Ruderisch, N., Deng, S., Liddelow, S.A., Zhang, C., Daneman, R.,
1005 Maniatis, T., Barres, B.A., Wu, J.Q., 2014. An RNA-Sequencing Transcriptome and
1006 Splicing Database of Glia, Neurons, and Vascular Cells of the Cerebral Cortex. *J. Neurosci.*
1007 34, 11929–11947. doi:10.1523/JNEUROSCI.1860-14.2014

1008 Zhang, Y., James, M., Middleton, F.A., Davis, R.L., 2005. Transcriptional analysis of multiple brain
1009 regions in Parkinson’s disease supports the involvement of specific protein processing,
1010 energy metabolism, and signaling pathways, and suggests novel disease mechanisms. *Am.*
1011 *J. Med. Genet. Part B Neuropsychiatr. Genet. Off. Publ. Int. Soc. Psychiatr. Genet.* 137B,
1012 5–16. doi:10.1002/ajmg.b.30195

1013 Zhang, Y., Sloan, S.A., Clarke, L.E., Caneda, C., Plaza, C.A., Blumenthal, P.D., Vogel, H.,
1014 Steinberg, G.K., Edwards, M.S.B., Li, G., Duncan, J.A., Cheshier, S.H., Shuer, L.M., Chang,
1015 E.F., Grant, G.A., Gephart, M.G.H., Barres, B.A., 2016. Purification and Characterization of
1016 Progenitor and Mature Human Astrocytes Reveals Transcriptional and Functional
1017 Differences with Mouse. *Neuron* 89, 37–53. doi:10.1016/j.neuron.2015.11.013

1018 Zoubarev, A., Hamer, K.M., Keshav, K.D., McCarthy, E.L., Santos, J.R.C., Van Rossum, T.,
1019 McDonald, C., Hall, A., Wan, X., Lim, R., Gillis, J., Pavlidis, P., 2012. Gemma: a resource
1020 for the reuse, sharing and meta-analysis of expression profiling data. *Bioinforma. Oxf. Engl.*
1021 28, 2272–2273. doi:10.1093/bioinformatics/bts430

1022

Supporting information

Supplementary file 1:

Figure S1: Expression of dentate granule cell markers discovered in the study in Allen Brain

Atlas mouse brain in situ hybridization database.

Figure S2: Expression of Purkinje markers discovered in the study in Allen Brain Atlas

mouse brain in situ hybridization database.

Table S1: Validation status of dentate granule cell markers.

Table S2: Validation status of Purkinje cell markers.

Table S3: Intersection of Purkinje markers from NeuroExpresso markers and Rong et al. 2004

study. Rong et al. genes are taken from Table 2 of the paper. Probe to gene annotations are

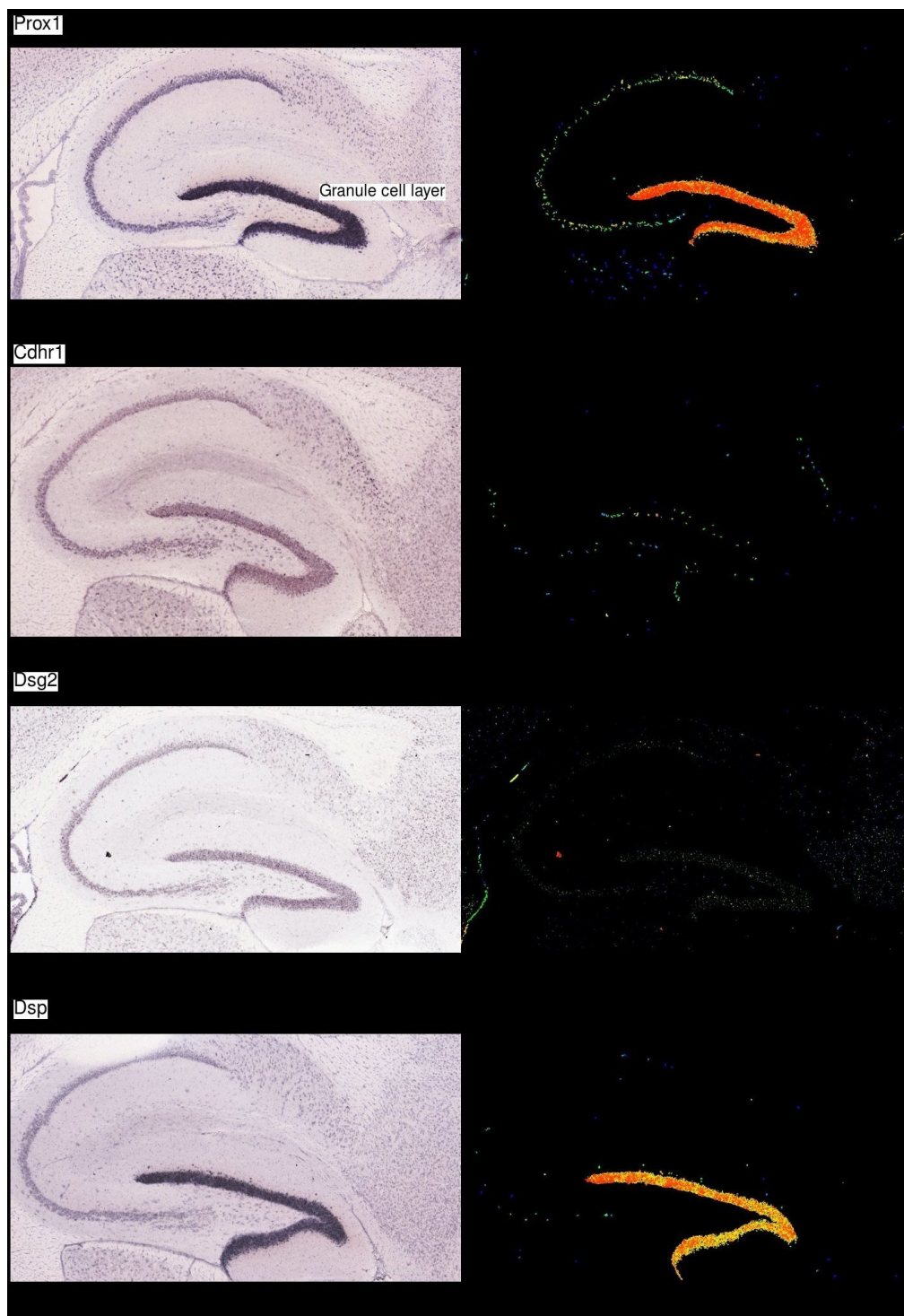
repeated using annotations from Gemma. Fold change column shows the difference of expression

between wild type and pcd^{3J} mouse which lacks purkinje cells according to Rong et al.

Additional material

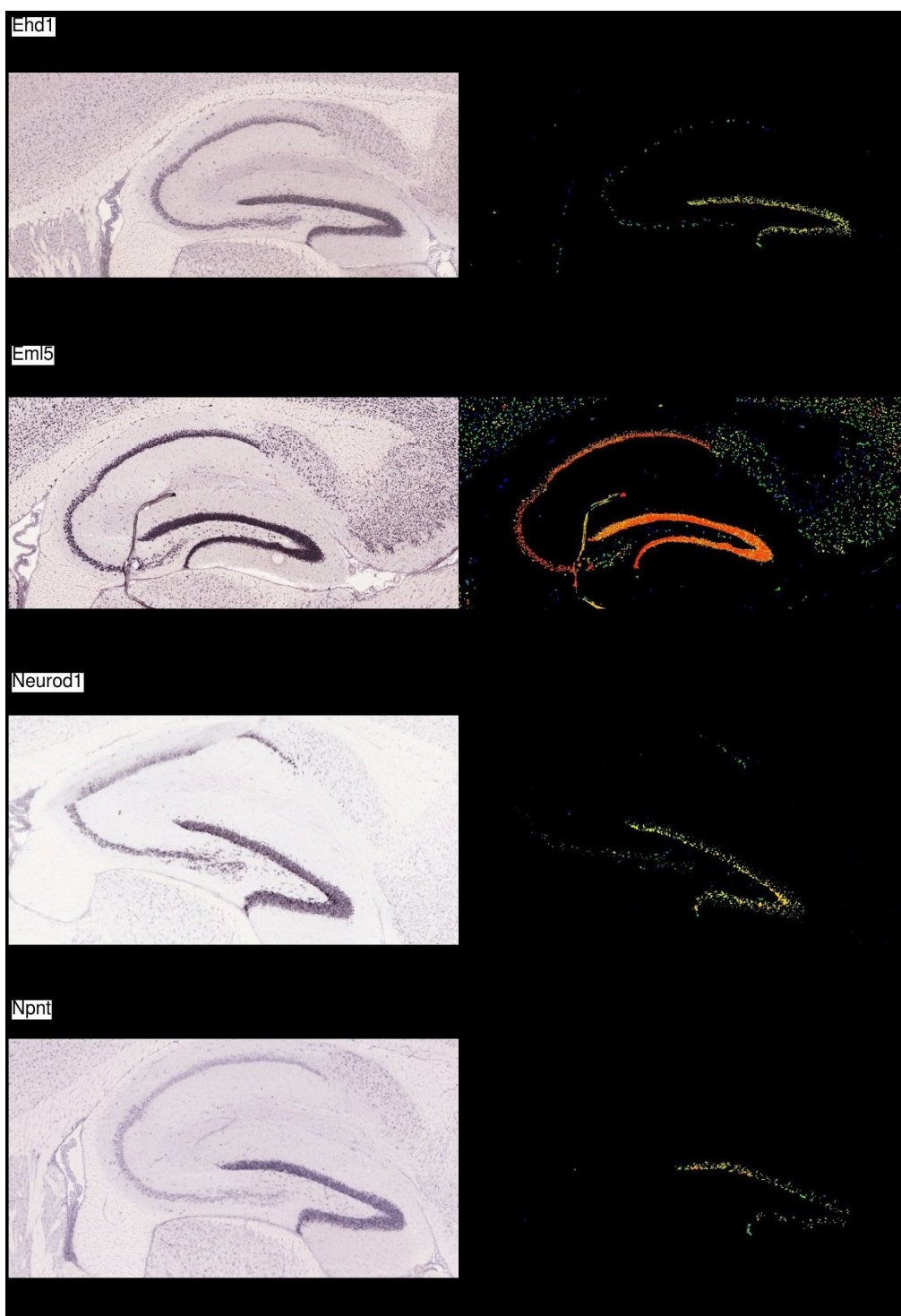
Code for analysis at github.com/oganm/neuroExpressoAnalysis

Web application at neuroexpresso.org



1039

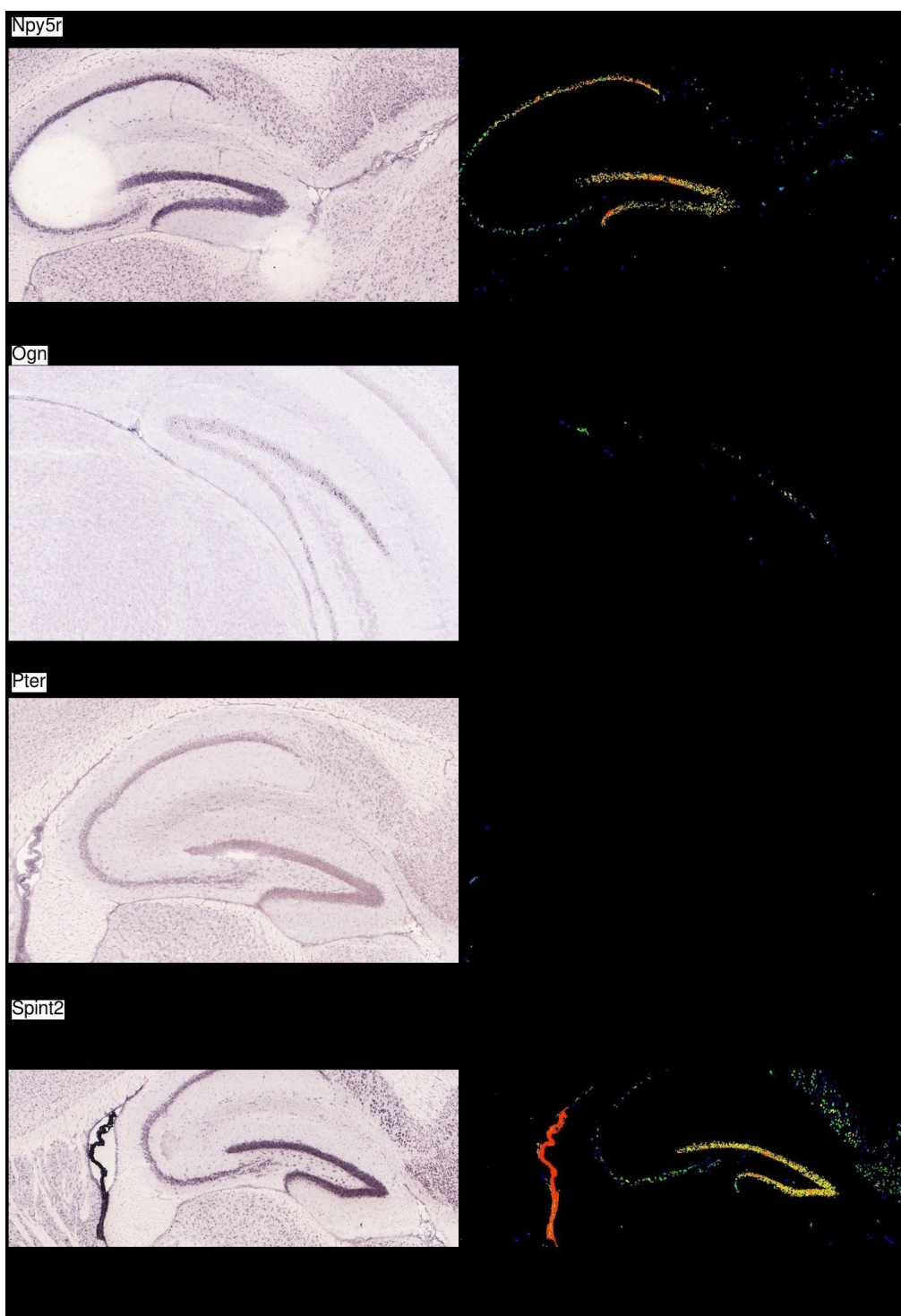
1040 Figure S1: Expression of dentate granule cell markers discovered in the study in Allen Brain Atlas mouse
 1041 brain in situ hybridization database. The first gene is Prox1, a known marker of dentate granule cells. The
 1042 intensity is color-coded to range from blue (low expression intensity), through green (medium intensity) to red
 1043 (high intensity). All images except Ogn is taken from the sagittal view. Ogn is taken from the coronal view.



1044

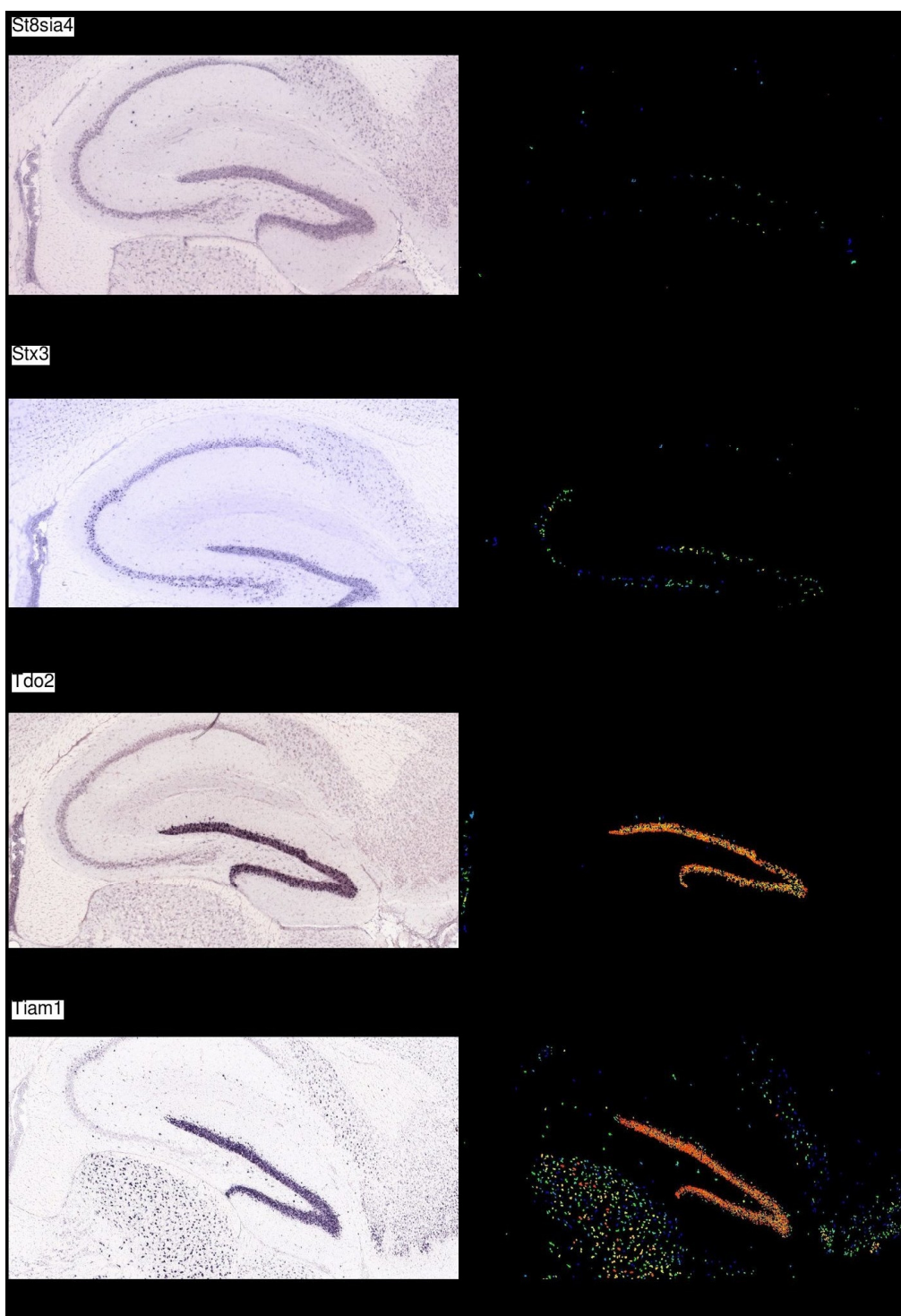
1045 Figure S1 continued.

1046



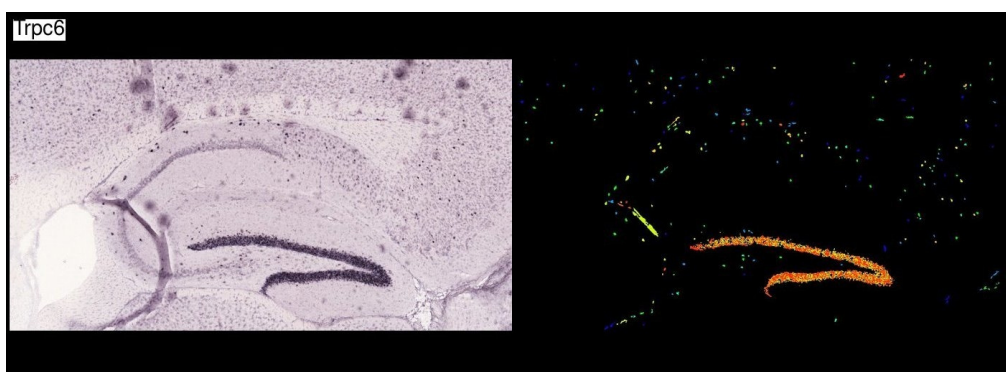
1047

1048 Figure S1 continued.



1049

1050 Figure S1 continued.



1051

1052 Figure S1 continued.

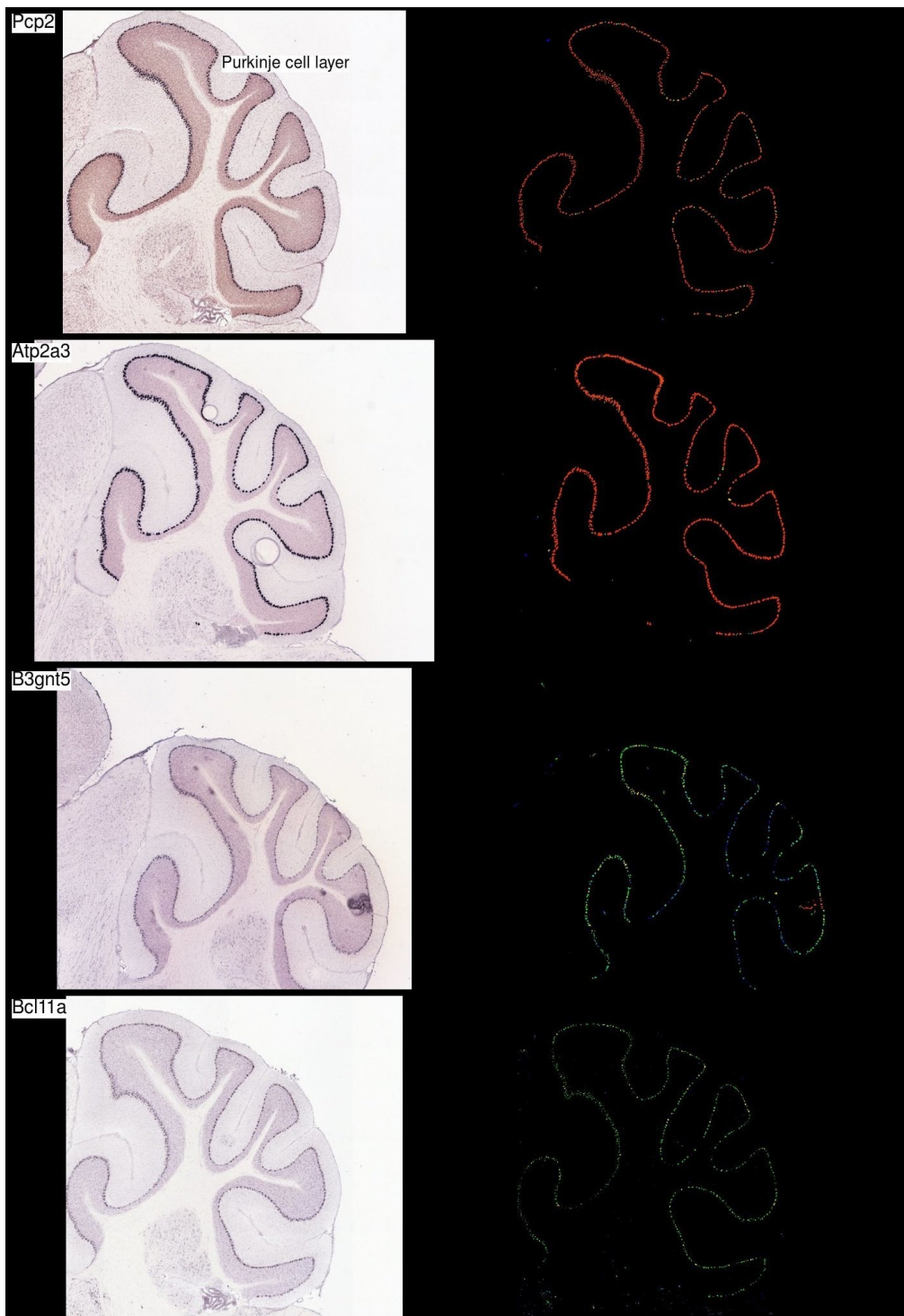
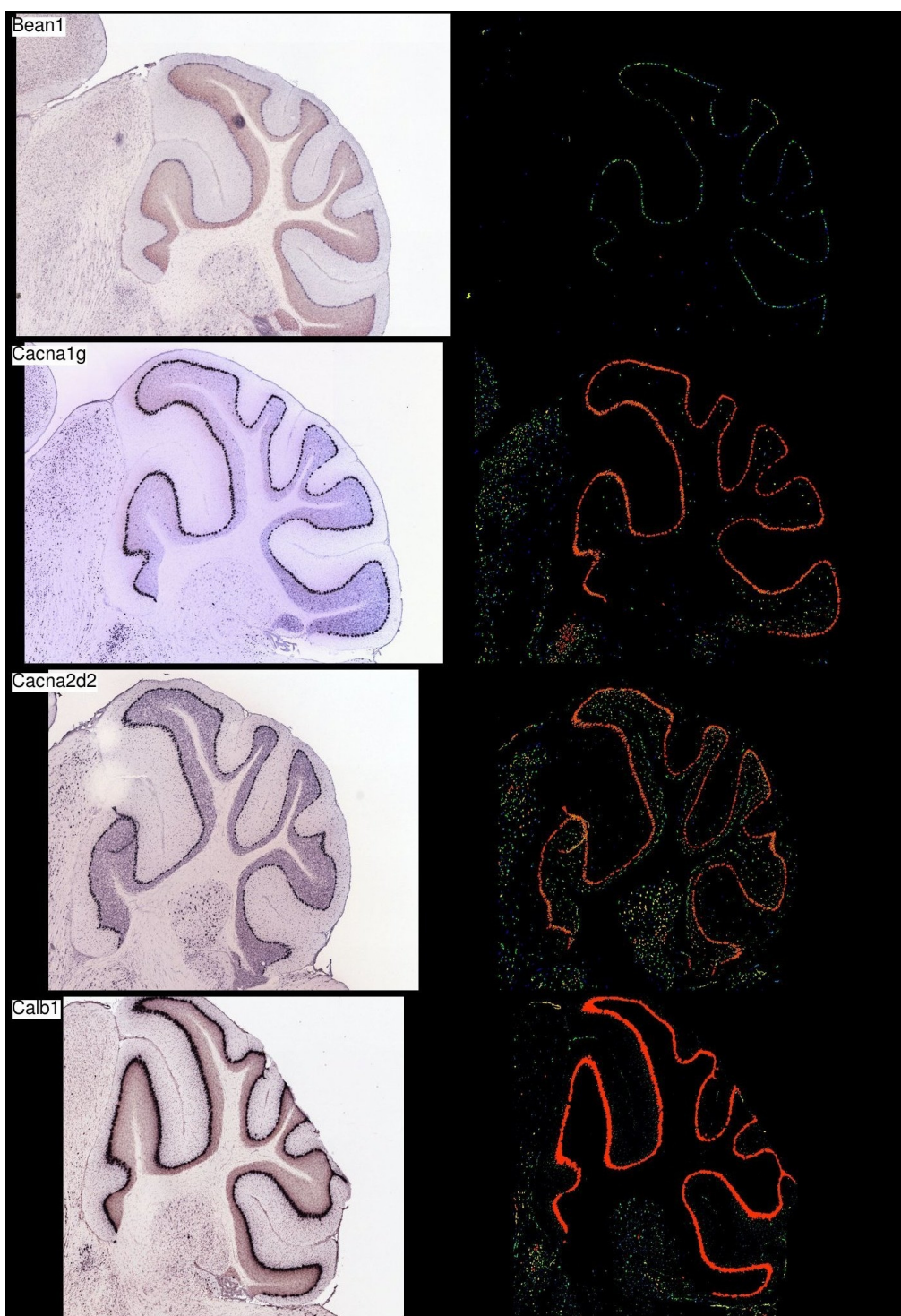


Figure S2: Expression of Purkinje markers discovered in the study in Allen Brain Atlas mouse brain in situ hybridization database. The first gene is Pcp2, a known marker of Purkinje cells. The intensity is color-coded to range from blue (low expression intensity), through green (medium intensity) to red (high intensity).} All images are taken from the sagittal view.



1059

1060 Figure S2 continued.

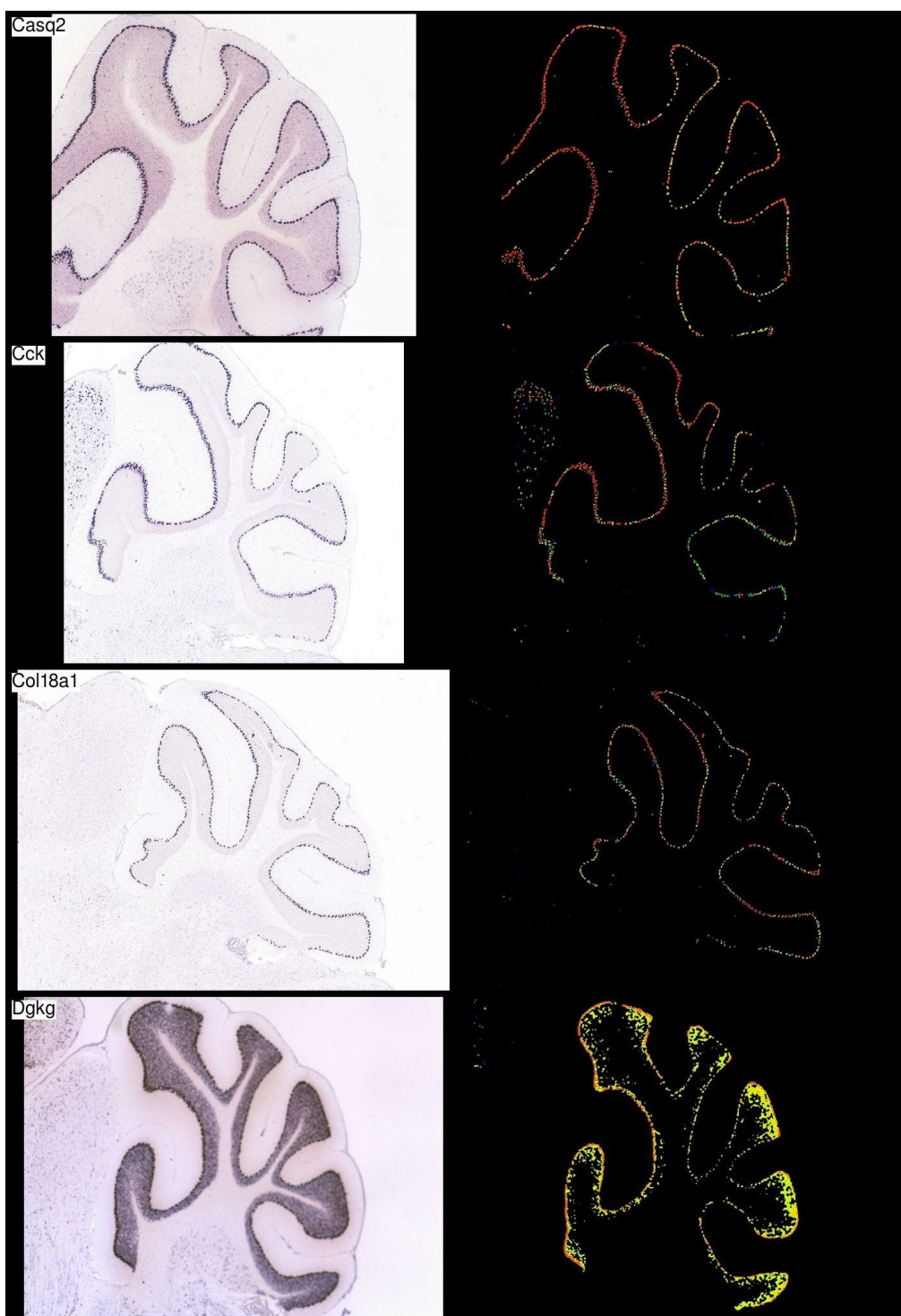


Figure S2 continued.

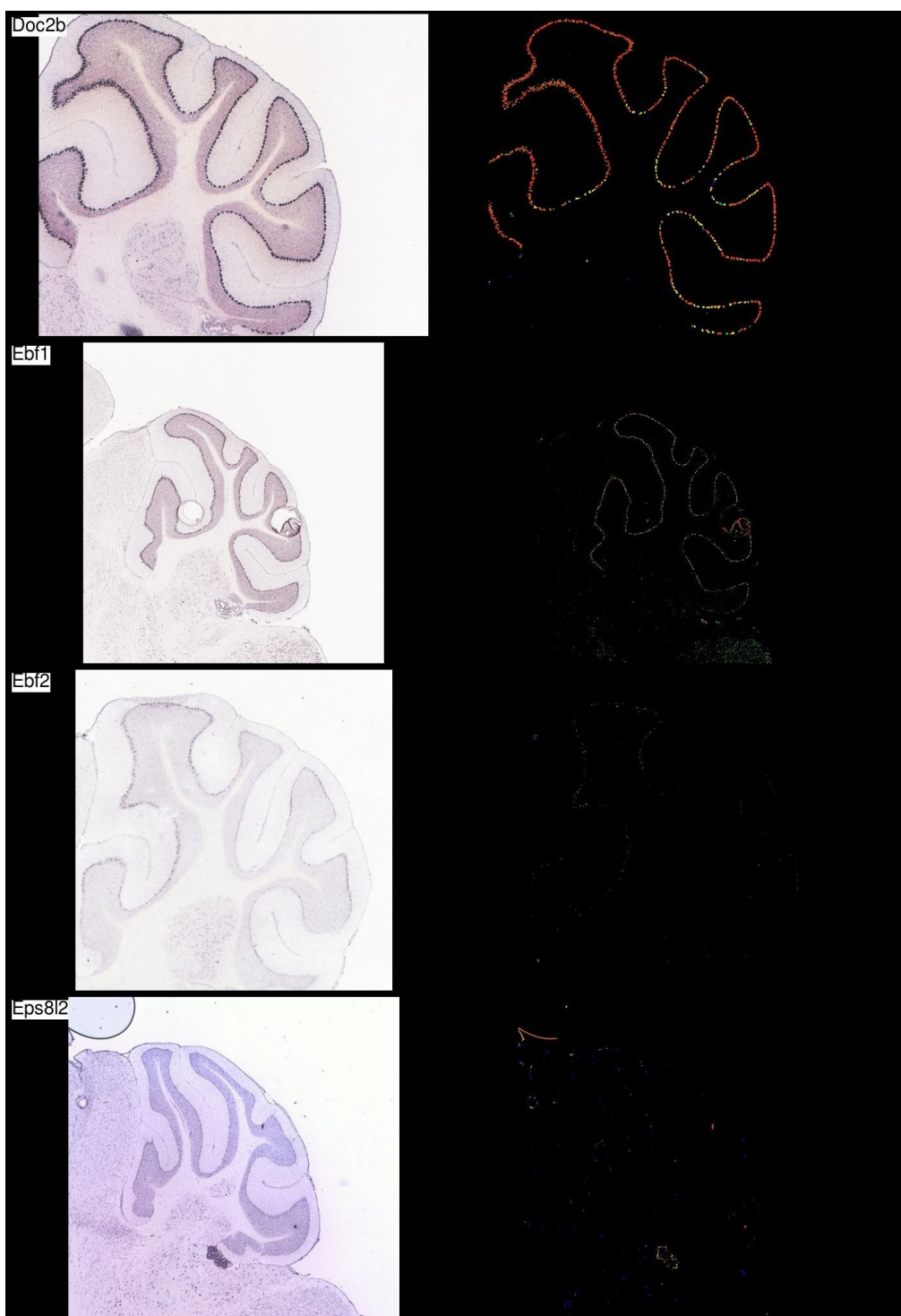


Figure S2 continued.

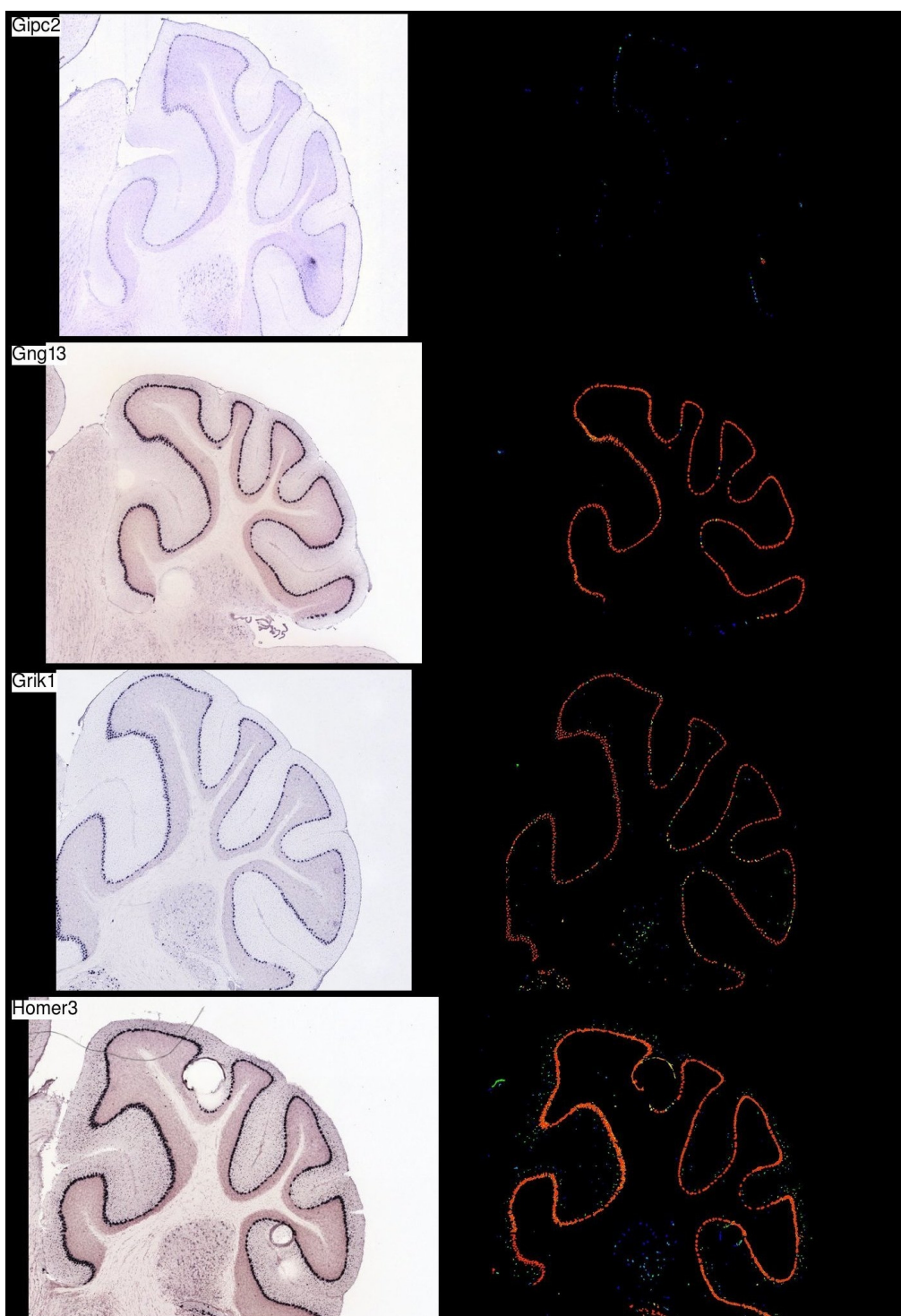
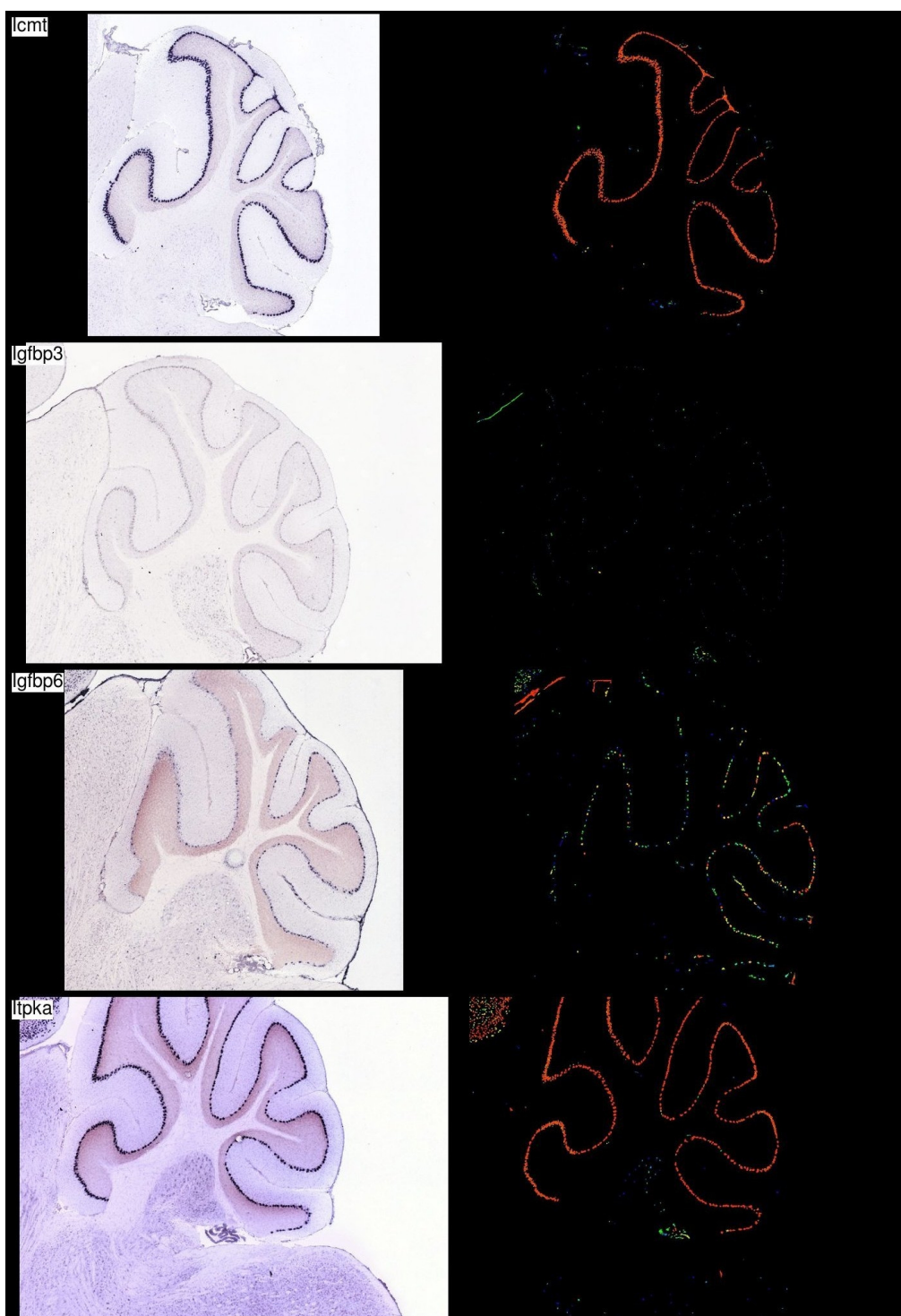


Figure S2 continued.



1067

1068 Figure S2 continued.

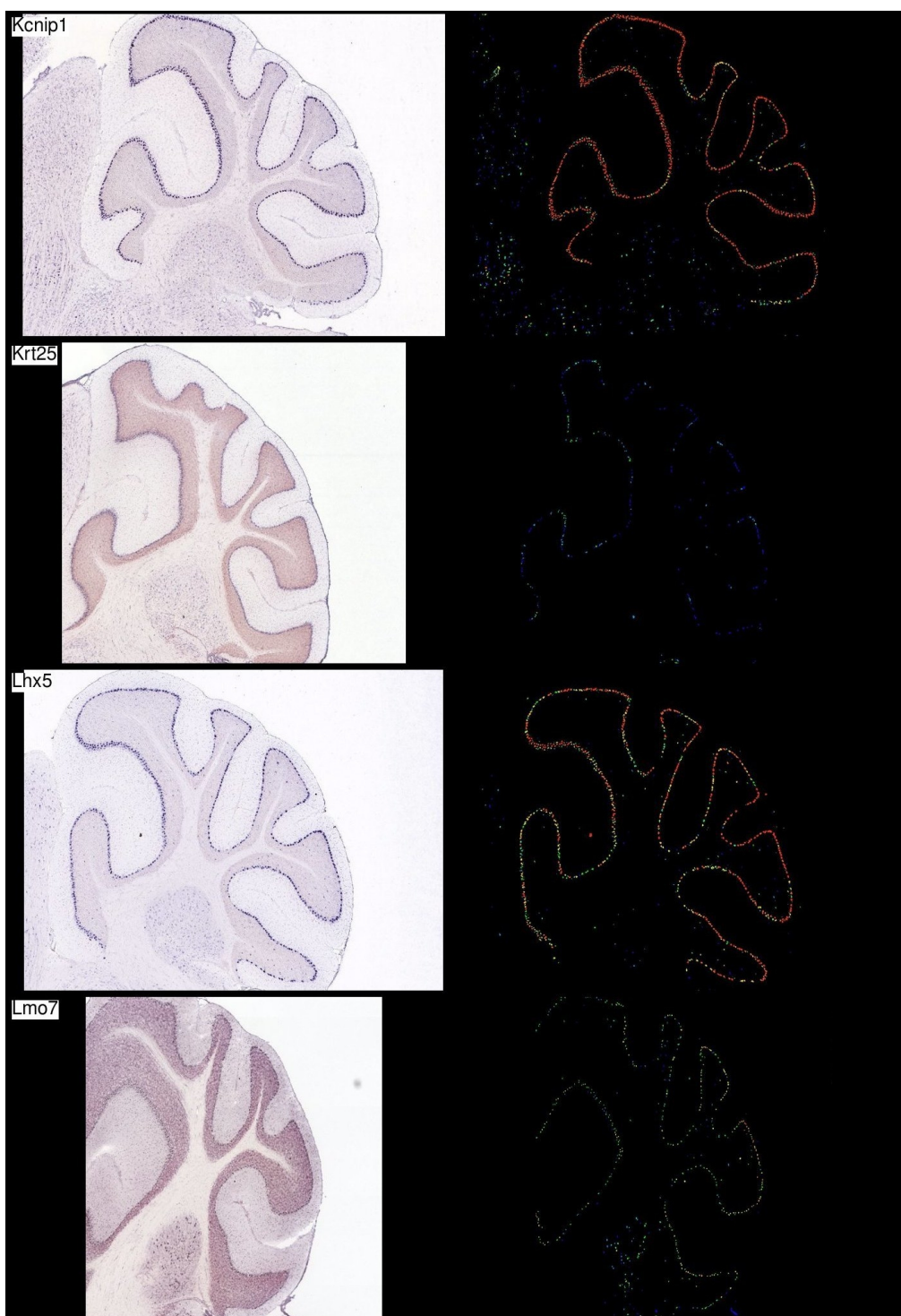


Figure S2 continued.

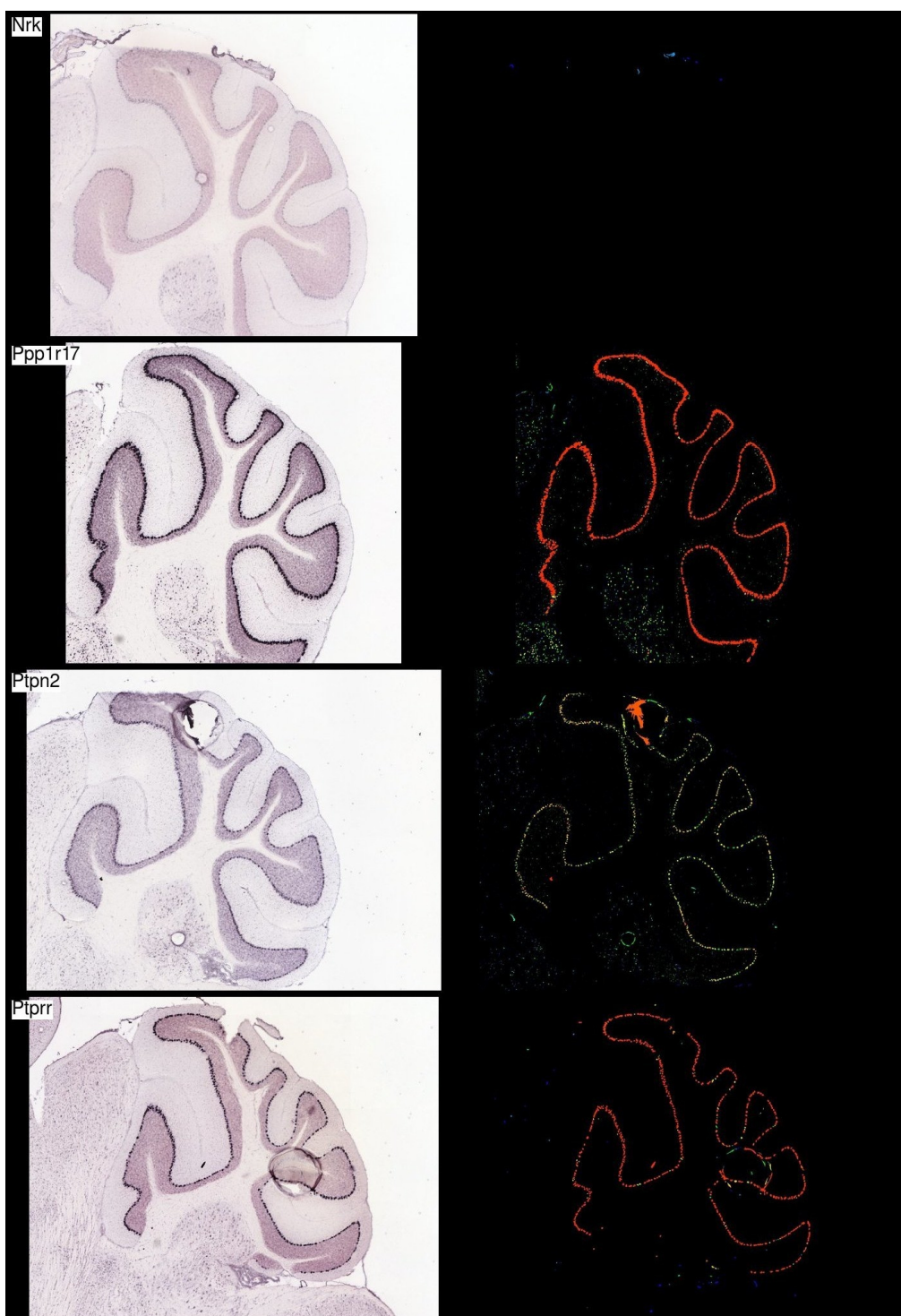


Figure S2 continued.

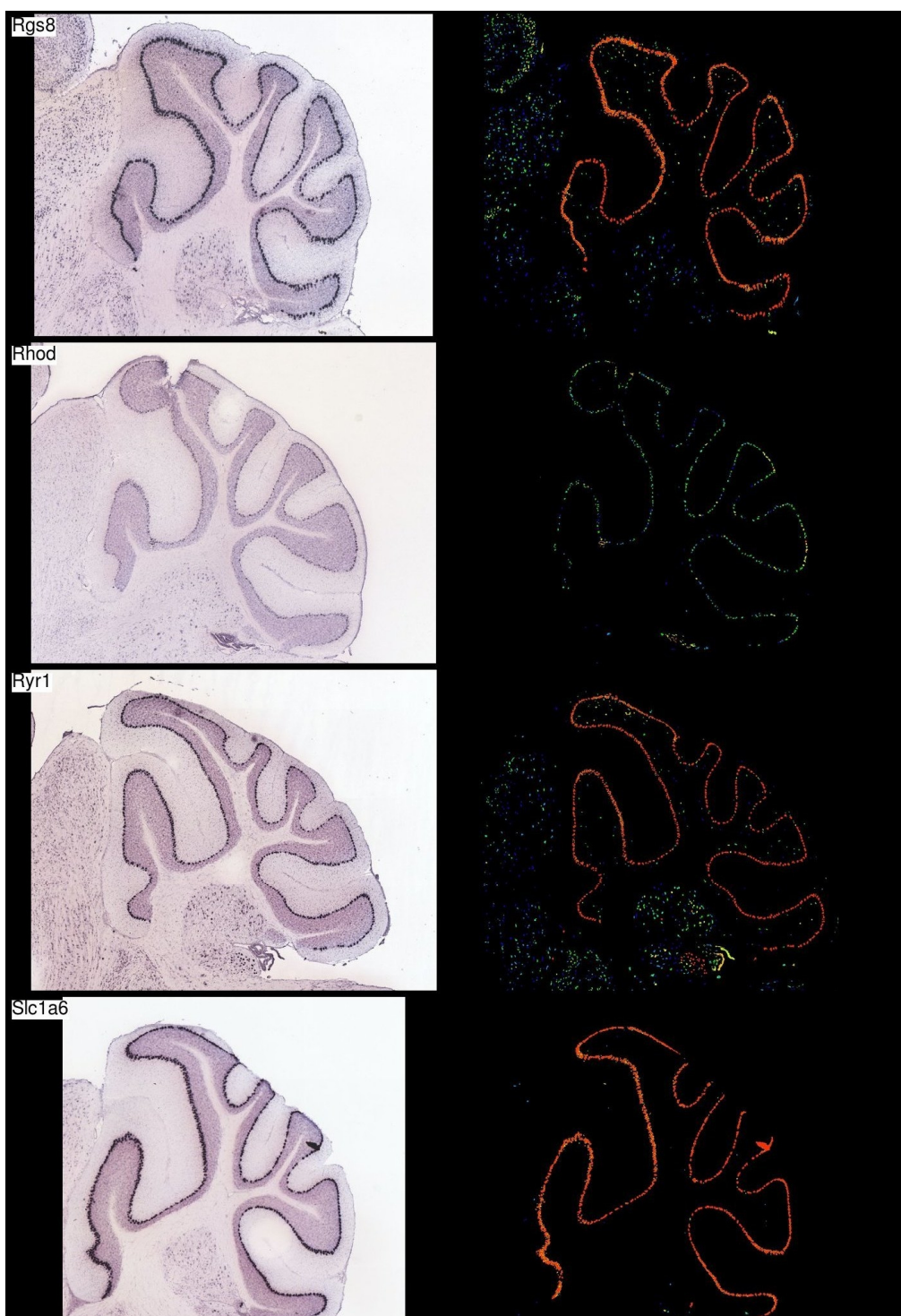


Figure S2 continued.

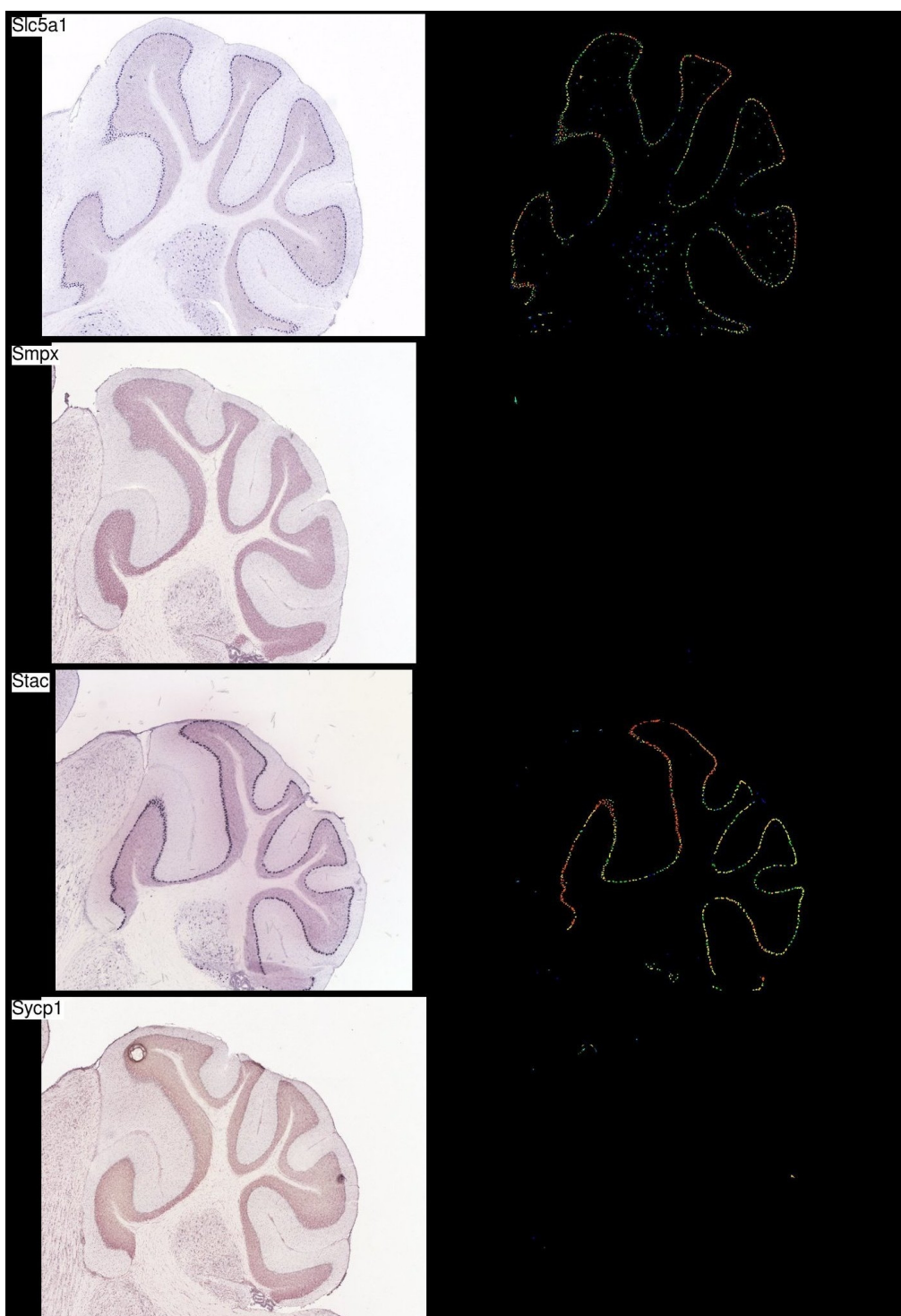
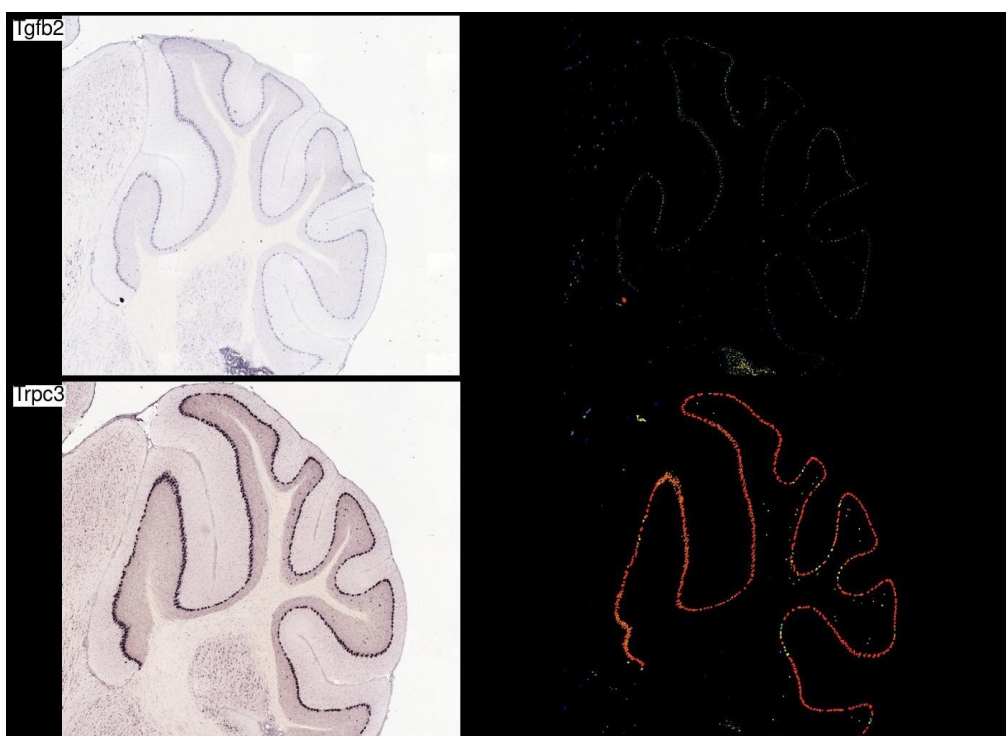


Figure S2 continued.



1077

1078 Figure S2 continued.

Gene	Status	Notes
Cdhr1	Ok	
Dsg2	Ok	
Dsp	Ok	
Ehd1	Ok	
Emi5	Not specific	
Neurod1	Ok	
Npnt	Ok	
Npy5r	Ok	
Ogn	Ok	Low expression
Pter	No signal	Not expressed anywhere else in the rest of the brain
Spint2	Ok	
St8sia4	Ok	Low expression
Stx3	Ok	
Tdo2	Ok	
Tiam1	Ok	
Trpc6	Ok	

1079 Table S1: Validation status of dentate granule cell markers.

Gene	Status	Notes
Atp2a3	Ok	
B3gnt5	Ok	
Bcl11a	Ok	2nd probeset
Bean1	Ok	
Cacna1g	Ok	
Cacna2d2	Not specific	
Calb1	Ok	
Casq2	Ok	
Cck	Ok	
Col18a1	Ok	
Dgkg	Ok	
Doc2b	Ok	
Ebf1	Ok	
Ebf2	Ok	Low signal
Eps8l2	Inconclusive	
Fam174b	Not in ABA	
Gipc2	Ok	Low signal
Gng13	Ok	
Grik1	Ok	
Homer3	Ok	
Icmt	Ok	
Igfbp3	Ok	Low signal
Igfbp6	Ok	
Itpka	Ok	
Kcnip1	Ok	
Krt25	Ok	
Lhx5	Ok	
Lmo7	Ok	
Nrk	Ok	Low signal
Ppp1r17	Ok	
Ptpn2	Ok	
Ptprr	Ok	
Rgs8	Ok	
Rhod	Ok	
Ryr1	Ok	
Slc1a6	Ok	
Slc5a1	Ok	
Smpx	No signal	Not expressed anywhere else in the rest of the brain
Stac	Ok	
Sycp1	No signal	Not expressed anywhere else in the rest of the brain
Tgfb2	Ok	
Trpc3	Ok	
Tuba8	Not in ABA	

1080 Table S2: Validation status of Purkinje cell markers.

Gene Symbol	Probeset	Fold change
Slc1a6	1418933_at	2.03
Cck	1419473_a_at	3.53
Gng13	1419414_at	3.88
Ppp1r17	1449240_at	4.12
Calb1	1417504_at	5.63
Calb1	1448738_at	7.42
Itpka	1424037_at	7.76
Cacna1g	1423365_at	12.51
Doc2b	1420667_at	14.08
Icmt	1426500_at	14.11
Rgs8	1453060_at	15.17
Nrk	1450079_at	15.83
Casq2	1422529_s_at	16.36
Trpc3	1417577_at	21.03
Kcnip1	1448459_at	21.41
Homer3	1424859_at	27.8
Atp2a3	1421129_a_at	28.17
Kcnip1	1416785_at	28.5
Atp2a3	1450124_a_at	31.57
Sycp1	1427291_at	33.16
Doc2b	1420666_at	36.04
Ptprr	1426047_a_at	36.41
Fam174b	1434273_at	43.73
Slc5a1	1419057_at	46.8

Table S3: Intersection of Purkinje markers from NeuroExpresso markers and Rong et al. 2004 study. Rong et al. genes are taken from Table 2 of the paper. Probe to gene annotations are repeated using annotations from Gemma. Fold change column shows the difference of expression between wild type and pcd^{3J} mouse which lacks Purkinje cells according to Rong et al.