

1 **Parallel RNA and DNA analysis after Deep-sequencing (PRDD-seq) reveals cell type-**
2 **specific lineage patterns in human brain**

3
4 August Yue Huang^{a,b,c,1}, Pengpeng Li^{a,b,c,1}, Rachel E. Rodin^{a,b,c,d}, Sonia N. Kim^{a,b,c,d}, Yanmei
5 Dou^e, Connor J. Kenny^{a,b,c}, Shyam K. Akula^{a,b,c,d}, Rebecca D. Hodge^f, Trygve E. Bakken^f,
6 Jeremy A. Miller^f, Ed S. Lein^f, Peter J. Park^e, Eunjung Alice Lee^{a,b,c}, Christopher A. Walsh^{a,b,c,d,2}

7
8 ^aDivision of Genetics and Genomics, Howard Hughes Medical Institute, and Manton Center for
9 Orphan Disease Research, Boston Children's Hospital, Boston, MA, USA.

10 ^bDepartments of Neurology and Pediatrics, Harvard Medical School, Boston, MA, USA.

11 ^cBroad Institute of MIT and Harvard, Cambridge, MA, USA.

12 ^dProgram in Neuroscience and Harvard/MIT MD-PHD Program, Harvard Medical School,
13 Boston, MA, USA.

14 ^eDepartment of Biomedical Informatics and Ludwig Center at Harvard, Harvard Medical School,
15 Boston, MA, USA.

16 ^fAllen Institute for Brain Science, Seattle, WA, USA

17 ¹These authors contributed equally to this work.

18 ²Corresponding Author

19 Christopher A. Walsh, MD. PhD

20 Division of Genetics and Genomics, Boston Children's Hospital

21 Center for Life Sciences 15062

22 3 Blackfan Circle, Boston, MA 02115

23 Phone: 617-919-2923

24 Email: Christopher.walsh@childrens.harvard.edu

25

26 **Classification**

27 Biological Sciences: Neuroscience

28

29 **Keywords**

30 PRDD-seq, single-cell MosaicHunter, birthdating, cortical layer, neurodevelopment

31

32 **Author Contributions**

33 A.Y.H. and P.L. conceived the project and C.A.W. supervised it. A.Y.H. and P.L. developed the

34 scMH algorithm. P.L. developed PRDD-seq and performed experiments. A.Y.H. performed

35 computational and statistical analyses. R.E.R., S.N.K., and Y.D. helped with validation of

36 sSNVs. C.J.K. helped with and provided insight on the comparison of PRDD-seq and scWTA.

37 S.K.A. assisted with interpretation of neurodevelopmental discoveries. R.D.H., T.E.B., J.M. and

38 E.S.L. generated the MTG single-cell RNA sequencing data, and provided it prior to publication.

39 E.A.L. and P.J.P. provided suggestions on computational analyses. A.Y.H. and P.L. wrote the

40 manuscript supervised by C.A.W., with input from all other authors.

41

42 **This PDF file includes:**

43 Main Text

44 Figures 1 to 5

45 *SI Appendix*, Figure S1 to S5 and Tables S1 to S3

46

47 **Abstract**

48 Elucidating the lineage relationships among different cell types is key to understanding human
49 brain development. Here we developed Parallel RNA and DNA analysis after Deep-sequencing
50 (PRDD-seq), which combines RNA analysis of neuronal cell types with analysis of nested
51 spontaneous DNA somatic mutations as cell lineage markers, identified from joint analysis of
52 single cell and bulk DNA sequencing by single-cell MosaicHunter (scMH). PRDD-seq enables
53 the first-ever simultaneous reconstruction of neuronal cell type, cell lineage, and sequential
54 neuronal formation (“birthdate”) in postmortem human cerebral cortex. Analysis of two human
55 brains showed remarkable quantitative details that relate mutation mosaic frequency to clonal
56 patterns, confirming an early divergence of precursors for excitatory and inhibitory neurons, and
57 an “inside-out” layer formation of excitatory neurons as seen in other species. In addition our
58 analysis allows the first estimate of excitatory neuron-restricted precursors (about 10) that
59 generate the excitatory neurons within a cortical column. Inhibitory neurons showed complex,
60 subtype-specific patterns of neurogenesis, including some patterns of development conserved
61 relative to mouse, but also some aspects of primate cortical interneuron development not seen in
62 mouse. PRDD-seq can be broadly applied to characterize cell identity and lineage from diverse
63 archival samples with single-cell resolution and in potentially any developmental or disease
64 condition.

65

66 **Significance Statement**

67 Stem cells and progenitors undergo a series of cell divisions to generate the neurons of the brain,
68 and understanding this sequence is critical to studying the mechanisms that control cell division
69 and migration in developing brain. Mutations that occur as cells divide are known as the basis of
70 cancer, but have more recently been shown to occur with normal cell divisions, creating a
71 permanent, forensic map of the clonal patterns that define the brain. Here we develop new
72 technology to analyze both DNA mutations and RNA gene expression patterns in single cells
73 from human postmortem brain, allowing us to define clonal patterns among different types of
74 human brain neurons, gaining the first direct insight into how they form.

75

76 **Introduction**

77 Although we have learned a great deal about development of the cerebral cortex from
78 animal models, we have remarkably little direct information about how the human brain, which
79 differs vastly in shape, size, and composition from the brains of non-primates, forms the neurons
80 of its cerebral cortex (1-4). Recent studies defining the fundamental cell types of the adult and
81 developing human cortex (5-7) form a foundation for understanding how these cell types develop,
82 how the unique aspects of the human cortex come about, and how developmental brain disorders
83 might alter patterns of cell lineage or cell type in human brain. However, whether individual
84 neural progenitor cells (NPCs) in embryonic stages are restricted to produce certain subtypes of
85 neurons, or multi-potential to generate all neuronal types, is still an open question even in model
86 animal species, since making this distinction requires simultaneous identification of cell lineage
87 and transcriptional analysis of cell type, which remains a technical challenge (8-12).

88 Somatic genetic mutations accumulate with each cell division during early development,
89 when spontaneous DNA damage escapes the DNA repair machinery, with single-nucleotide
90 variants (SNVs) being the most common mutation type (13-15). The timing of somatic mutations
91 can be inferred by either the cell fraction that carries each mutation or the co-occurrence status of
92 multiple mutations, in which early mutations should be shared by a large fraction of cells
93 whereas later mutations should be present in nested subpopulations of cells (16). Previous study
94 has shown the ability to use somatic SNVs as a rich internal lineage map to birthdate the
95 developmental timing of each neurons differentiated from neuronal progenitor cells (14) but has
96 not combined that with direct analysis of the subtypes of neurons, defined by morphology,
97 location, physiology, or RNA transcription pattern.

98 Single-cell transcriptomes provide granular information about cell identity (5-7), but it
99 cannot provide lineage maps as it fails to capture most somatic mutations, since somatic
100 mutations occur throughout the genome, most often in intronic or intergenic regions (16, 17).
101 Similarly, DNA-sequencing alone fails to provide information about cell identity, and so lineage
102 mapping using only somatic mutations from DNA sequencing is unable to address questions
103 about the lineage of specific cell identities in neurodevelopment. Somatic mutations in
104 mitochondrial DNA have been recently suggested as potential lineage marks as well, but the
105 modest target size of the mitochondrial genome, and the multiple diverse mitochondrial genomes
106 in each cell, represent challenges to the use of mitochondrial mutations as a rich source of stable
107 lineage markers (18).

108 To address this challenge, we developed Parallel RNA and DNA analysis after Deep-
109 sequencing (PRDD-seq) that identifies somatic SNVs (sSNVs) from single cell and bulk whole-
110 genome sequencing (WGS) data, with multiplexed detection of sSNVs and multiple RNA
111 marker transcripts from single nuclei. We then benchmarked the performance of the DNA and
112 RNA assays of PRDD-seq against bulk WGS and single-cell RNA sequencing (scRNAseq),
113 respectively. Applying PRDD-seq to two postmortem brains of individuals without neurological
114 disease allowed unprecedented quantitative analysis of cell lineage in the human brain. While
115 revealing the expected patterns of divergence of excitatory and inhibitory lineages and “inside-
116 out” generation of excitatory neurons, our PRDD-seq data also directly suggest complex patterns
117 of interneuron formation in the human brain.

118 **Results**

119 **Simultaneous cell type and lineage analysis of single-cells by PRDD-seq**

120 The workflow of PRDD-seq is illustrated in Figure 1. Single NeuN+ cortical neuronal
121 nuclei from prefrontal cortex (PFC) of postmortem human brain tissue were purified by
122 fluorescence-activated nuclear sorting (FANS) (16) (Fig. 1A), and subjected to one-step RT-
123 qPCR with target-specific primers for 1] cDNA specific for up to 30 marker genes of major
124 neuronal cell types, and 2] specific genomic DNA (gDNA) loci representing identified somatic
125 mutations (see below) as markers of cell lineage (Fig. 1B). Aliquots of the pre-amplified gDNA
126 and cDNA libraries were analyzed for the presence of specific somatic mutations and transcripts
127 by microfluidic genotyping and gene expression profiling, respectively, using the Fluidigm
128 Biomark system (Fig. 1C). The somatic mutations used in PRDD-seq were identified by single-
129 cell MosaicHunter (scMH), described below, a new bioinformatic tool to identify lineage-
130 informative sSNVs, jointly considering WGS data from MDA-amplified single cells and
131 matched deep (>200X) WGS from bulk DNA samples collected from the same brain region (Fig.
132 1D).

133 We first created a map of neuronal cell types by analyzing >25,000 single neuronal nuclei
134 -- FANS-sorted based on NeuN immunoreactivity -- by scRNAseq from two different datasets,
135 to create a cell type landscape onto which PRDD-seq analyzed neurons could be located. We
136 performed 10X Genomics scRNAseq of 10,967 NeuN+ nuclei from the same PFC region of one
137 of the brains from which DNA mutations were identified (Fig. 1E). t-SNE analysis of this dataset
138 defined 21 transcriptionally distinct cell clusters, including 8 excitatory neuron clusters that
139 further clustered into upper, middle, and lower layers, and 13 inhibitory neuron clusters that
140 could be further classified into SST+, PV+, VIP+, and LAMP5+ subtypes (Fig. 1F and *SI*

141 *Appendix*, Fig. S1) (5, 7). A recently published scRNAseq dataset of 15,928 single neuronal
142 nuclei from human middle temporal gyrus (MTG) (5), sorted by NeuN immunoreactivity
143 following microdissection of cerebral cortical layers, provided additional direct information
144 about layer location of neuronal types (Fig. 1G and *SI Appendix*, Fig. S2) and so was used for
145 cell type mapping in parallel. PFC and MTG share relatively generic cerebral cortical
146 architecture as “association” cortex, and clustering analysis of the two datasets (Fig. 1H) shows
147 that they identified similar major cell types, with cells clustering by cell type rather than by
148 platform, although the SMART-seq dataset from MTG defined finer subdivisions of cell type as
149 expected because of its larger sample size and deeper sequence depth.

150 We jointly analyzed single PRDD-seq cells and scRNAseq cells and mapped each
151 PRDD-seq cell onto the t-SNE maps of scRNAseq based on gene expression similarity (Fig. 1I,
152 see Methods). The cell type and cortical layer information of each PRDD-seq cell was then
153 imputed based on its assigned cluster in scRNAseq datasets. Finally, the combination of
154 genotype and gene expression information of PRDD-seq cells allowed lineage and birthdate
155 analysis of particular cell types (Fig. 1J), as well as analysis of cell type differentiation of
156 particular lineages (Fig. 1K).

157

158 **Discovery of lineage-informative sSNVs from bulk brain and single-neuron WGS data**

159 The resolution of lineage reconstruction is dependent on having a comprehensive list of
160 somatic mutations identified from the specific brain under analysis. Whereas deep WGS (e.g.,
161 200-250X coverage) of “bulk” DNA, isolated from tissue, efficiently identifies sSNVs present in
162 4% or more cells (19), it is insensitive to detecting later-occurring sSNVs that mark late cell
163 lineage events. On the other hand, WGS of DNA amplified from single neuronal nuclei (16)

164 identifies later-occurring sSNVs but is limited by cost and subject to artifacts during single-cell
165 amplification. Therefore, we developed scMH, which incorporates a Bayesian graphic model (20,
166 21) that integrates analysis of bulk WGS and single-cell WGS data to distinguish somatic
167 mutations from germline mutations and technical artifacts (Fig. 2A; see Methods). scMH first
168 calculates the likelihood and mosaic fraction of candidate sSNVs from a bulk DNA sample, and
169 then applies these values as the priors to genotype each candidate SNV across every single cell
170 being analyzed. The shared presence of a given sSNV in bulk DNA and one or more single cells
171 serves as validation of the sSNV. To expand the utility of scMH when a matched bulk sample is
172 unavailable, we further designed a “bulk-free” mode that can utilize a “synthetic” bulk WGS
173 dataset, generated by *in silico* merging of the many WGS datasets of multiple single-cells
174 obtained from the same donor. We benchmarked scMH using 45X single-cell WGS of 24
175 neurons—22 of which were sequenced in previous studies (16, 17)—as well as ~200X bulk
176 WGS of PFC (both from the brain of the same individual, UMB1465, who died at age 17 with no
177 neurological diagnosis), against existing single-cell sSNV callers including Monovar (22),
178 SCcaller (23), LiRA (24), and Conbase (25). Sensitivity and false discovery rate (FDR) were
179 estimated based on experimentally validated mutations and clade annotations identified
180 previously (16). With either PFC bulk or synthetic bulk, scMH outperformed the other tools and
181 achieved ~70% sensitivity to detect lineage-informative mutations with < 5% FDR; combining
182 both the default and “bulk-free” modes improved detection sensitivity to 93% without increasing
183 the FDR, suggesting that the “bulk-free” mode of scMH can detect sSNVs that are present in
184 multiple single-cells but may be undetectable in the bulk 200X WGS samples because of the low
185 mosaic fraction of these late mutations (Fig. 2B).

186 Applying scMH to data from brains of three normal individuals (UMB1465, UMB4638,
187 and UMB4643 (16, 17), identified and validated 42, 19, and 22 sSNVs, respectively (Fig. 2C-E,
188 and *SI Appendix*, Table S1), with an overall validation rate of 74.8% determined by Sanger
189 sequencing of independently sorted neurons from the same brain region. The number and
190 validation rate of lineage-informative sSNVs detected by scMH dramatically increased from
191 previous studies (16, 17). sSNVs identified from all three brains showed an enrichment in C>T
192 mutations, especially in CpG sites (*SI Appendix*, Fig. S3), a pattern observed in other studies of
193 embryonic mutations and cancer mutations (13, 26), since such C>T mutations appear to be
194 caused by cytosine deamination that is replicated into a fixed SNV before it can be repaired (27).
195 Unsupervised clustering analysis grouped the 24 sequenced neurons from UMB1465 into six
196 different clades; no cells harbored mutations of multiple clades, suggesting the high accuracy of
197 scMH for single-cell genotyping of sSNVs (Fig. 2C). In clades C and E, we observed neurons
198 that shared early mutations but harbored different sets of later mutations, suggesting that they
199 were derived from different branches of the same clades (Fig. 2C). Clustering of ten and nine
200 sequenced neurons from UMB4638 and UMB4643—respectively by their sSNVs—
201 demonstrated similar nested patterns forming three primary clades for each individual and also
202 showed evidence for branches of these clades (Fig. 2D, E). The mosaic fraction of each sSNV in
203 “bulk” DNA (Fig. 2C, D, E) was used as an additional indicator of the sequence in which sSNV
204 occurred, since early sSNVs tend to be found in many single cells, as well as at higher mosaic
205 fraction in bulk DNA, whereas later mutations appear in fewer cells and lower mosaic fraction in
206 bulk DNA. These two findings correlated very strongly.

207

208 **Lineage and cell type identity of single-neurons revealed by PRDD-seq**

209 To assess the performance of PRDD-seq in capturing lineage and cell type information
210 from single-cells, we applied PRDD-seq to 1,710 cortical neurons from UMB1465 PFC, using
211 probes to detect 30 out of 42 validated sSNVs in UMB1465, for which we successfully designed
212 highly specific and sensitive probes (*SI Appendix*, Table S1), along with 30 marker genes whose
213 expression levels distinguish major inhibitory and excitatory neuronal subtypes and cortical
214 layers identified in the scRNAseq datasets (5, 7) (*SI Appendix*, Table S2). Overall, PRDD-seq
215 mapped 1,112/1,710 (65%) cortical neurons from UMB1465 PFC into 20 lineage branches and 6
216 major clades (Fig. 3A). For each major clade, birthdate-ordered lineage branches were inferred
217 from the nested sSNVs, where earlier derived neurons contained fewer clonal mutations, and
218 neurons generated later harbored additional mutations from subsequent cell divisions (16). The
219 nested nature of sSNVs in clades allow cells to be placed into clades using multiple sSNVs, so
220 that cells whose genomes were subject to allelic dropout—which is not uncommon when single
221 cell DNA molecules are amplified—could still be placed into clades based on other sSNV from
222 the same clade (Fig. 3A and *SI Appendix*, Table S1). On the other hand, only 71/1710 (4.2%)
223 neurons contained sSNVs from multiple clades, suggesting a low rate of false positive
224 amplification or sorting of multiple nuclei into single wells in the DNA assay of PRDD-seq (Fig.
225 3B, upper panel). 527/1710 (30.8%) neurons showed the absence of any sSNVs from the 6
226 clades; these neurons may be from other clades in which we did not discover sSNV markers (Fig.
227 3B, upper panel). In PRDD-seq cells, mosaic fractions of sSNVs correlated linearly with the
228 fractions calculated from ~200X bulk WGS, indicating generally unbiased sSNV detection (Fig.
229 3B, lower panel and Fig. 3C), and allowing confident inference of the developmental sequence
230 of sSNVs according to the nested pattern.

231 Among the 1,112 PRDD-seq cells that were successfully claded, we ran the RNA assay
232 of PRDD-seq to measure the expression of 30 marker genes for each cell. Our evaluation using
233 simulation data derived from our own and published scRNAseq datasets (see Methods)
234 suggested that these 30 marker genes were sufficiently informative to infer many aspects of cell
235 type and dissected layer annotation (Fig. 3D), with an average accuracy of 84% for cortical layer
236 classification (within +/- one-layer difference) and 83% for inhibitory neuron subtype
237 classification. We then utilized expression of these 30 makers to successfully classify 747/1,112
238 PRDD-seq neurons (67.2%) from UMB1465 into 3 excitatory subgroups—corresponding to
239 upper, middle, or lower cortical layers—and 4 inhibitory subgroups: somatostatin positive
240 (SST+), vasoactive intestinal peptide-positive (VIP+), lysosomal associated membrane protein 5-
241 positive (LAMP5+), and putative parvalbumin-positive (putative PVALB+, or pPVALB+), since
242 probes for PVALB were not always directly assayed (Fig. 3E). PRDD-seq cells assigned to
243 upper, middle, and lower layers by the 10X PFC scRNAseq dataset were also enriched in L2-L3,
244 L4-L5, and L6 markers according to the SMART-seq MTG scRNAseq dataset, respectively,
245 indicating the similarity of the cell type compositions between PFC and MTG, the similarity of
246 the results with the two RNAseq methods, as well as the robustness of the mapping algorithm
247 (Fig. 3E, upper panel). Both our 10X scRNAseq dataset and PRDDseq analysis of UMB1465
248 and UMB4638 showed higher proportions of inhibitory neurons (43-47%) than reported with
249 other methods, however this ratio was very similar between the three experiments, suggesting
250 that the ratio reflects our particular NeuN+ sorting protocol rather than technical aspects of the
251 cell typing methods (Fig. 3F upper panel). We observed remarkably similar layer and subtype
252 distribution between PRDD-seq and scRNAseq cells for excitatory neurons (Chi-square test; Fig.
253 3F, middle panel). Among inhibitory neurons, pPVALB+ inhibitory neurons showed a higher

254 proportional representation in PRDD-seq than in scRNAseq, suggesting that a few neurons in
255 this category might reflect amplification failure of the other inhibitory probes (SST, VIP, and
256 LAMP5). In summary, our analysis suggests that PRDD-seq captures the major aspects of cell
257 types, without systematic loss of any given cell type.

258

259 **Early divergence of progenitors for excitatory and inhibitory neurons**

260 The simultaneous analysis of lineage and gene expression from the same neurons enabled
261 us to study the change of cell type contribution during early neurogenesis. Using PRDD-seq, we
262 profiled >2700 neurons from two brains, UMB1465 and UMB4638, and successfully captured
263 both lineage and cell type information from 747 and 480 neurons, respectively. In both
264 UMB4638 and UMB1465, all lineage clades showed early sSNVs in both excitatory and
265 inhibitory neurons, reflecting mutations occurring during early embryogenesis before the
266 divergence of these cell types, whereas late SNVs show progressive restriction to one or the
267 other cell type (Fig. 4A, B). Among the six major clades in UMB1465, clade C contained seven
268 nested branches with mosaic fractions diminishing from 0.33 to 0.0067 (Figure 3A and *SI*
269 *Appendix*, Table S1), with an increasing percentage of excitatory neurons containing mutations
270 C1 to C5, and only excitatory neurons containing mutations C6 to C7 (Fig. 4A), while clade F
271 showed similar progressive restriction. Similarly, both clade A and B in UMB4638 showed
272 nested mutations that became progressively limited to excitatory neurons (Fig. 4B). Interestingly,
273 the excitatory neurons appeared exclusively in branches with mosaic fraction below ~0.04 (Fig.
274 4A, B, and *SI Appendix*, Table S1), corresponding to a progenitor giving rise to about 4% of the
275 total cells in that cortical sample. Considering that ~40% of cortical cells are excitatory neurons,
276 with the remainder being glial cells or inhibitory neurons (28, 29), this observation suggests that

277 ten or more excitatory neuronal progenitor cells (NPCs) generate excitatory neurons in a given
278 cortical area, or “column”; the fact that 6-7 (including a branched clade) excitatory precursors
279 are explicitly marked by non-overlapping clades, and account for 60-70% of excitatory neurons
280 in our sample, independently supports this estimate. On the other hand, two clades (clade A and
281 B) from UMB1465 are statistically enriched for inhibitory neurons (two-sided one-proportion Z-
282 test’s $P < 0.05$), with the percentage of inhibitory neurons increasing from B1 to B2 (Fig. 4A).
283 These results show that at least some human NPCs demonstrate restricted cell type output,
284 supporting the model first established in mice (30-32) and strongly supported by conserved gene
285 expression patterns in the ganglionic eminence between humans and non-humans (33, 34), that
286 excitatory and inhibitory neurons are generated from distinct progenitor regions.

287

288 **“Inside-out” order of cortical layer formation for excitatory neurons**

289 Further sub-typing of excitatory neurons using laminar markers revealed layer-specific
290 patterns of excitatory neuron neurogenesis. For example, in UMB1465, the percentage of lower
291 layer neurons carrying a mutation decreased from mutations C1 to C4, and no deep-layer
292 neurons were detected carrying C5 to C7, with the percentage of upper layer neurons increasing
293 correspondingly from C1 to C7 (Pearson correlation’s $P = 2.9 \times 10^{-3}$; Fig. 4C, upper panel). To
294 gain more precise layer identities of PRDD-seq cells, we mapped them to the SMART-seq MTG
295 scRNAseq dataset obtained after layer microdissection using the same methods as earlier (5),
296 which generated similar “birthdate” patterns in clade C, with early lineage sSNVs present in all
297 layers, and later sSNVs restricted to middle and upper layers (Pearson correlation’s $P = 1.4 \times 10^{-3}$;
298 Fig. 4C, lower panel). A similar trend was also observed in clades A and B in UMB4638.
299 Mapping PRDD-seq cells of UMB4638 to both 10X PFC and SMART-seq MTG scRNAseq

300 datasets showed that cells with later lineage markers were restricted to middle and upper layers
301 (Fig. 4D). These results together directly indicate that human cortical excitatory neurons are
302 formed in “inside-out” sequence after preplate cells are born, similar to mouse and non-human
303 primates (35-37). Furthermore, it suggests that neurons in lower cortical layers begin becoming
304 postmitotic relatively quickly after progenitors are specialized for excitatory neuron production.

305

306 **Diverse spatiotemporal patterns of development of inhibitory neuron subtypes**

307 Mapping PRDD-seq cells onto two different scRNAseq datasets also allowed analysis of
308 cortical inhibitory neurons, which originate from multiple developmentally transient structures of
309 the ventral telencephalon, including the medial, lateral and caudal ganglionic eminences (MGE,
310 LGE, and CGE), and migrate into dorsal cortex (30, 38). However, the highly dispersed nature of
311 inhibitory neuron clones observed in animal models (39-41) suggests that sSNVs in the
312 inhibitory lineage are likely to be present at exceedingly low allele frequencies in bulk DNA and
313 tiny fractions of single cells, so that only sSNVs occurring relatively early in development have
314 been analyzed so far. Inhibitory neurons derived from MGE and CGE can be distinguished by
315 expression of specific markers (5, 6), and PRDD-seq analysis showed that interneurons with
316 diverse marker genes were generated over the same developmental window (Fig. 5A, B). The
317 analyzed sSNVs were shared by multiple inhibitory subtypes, with hints that late marks might be
318 more limited to cell types, but no differences that reached statistical significance (FDR-adjusted
319 Chi-square test's $P > 0.05$). Previous studies cataloging interneurons in mouse and human have
320 suggested that MGE-derived inhibitory neuron subtypes (SST+ and PVALB+) are enriched in
321 infragranular cortical layers, while CGE-derived interneuron subtypes (LAMP5/PAX6+, VIP+)
322 tend to occupy upper cortical layers preferentially (5, 42, 43) and thus our mapping of PRDD-seq

323 cells onto scRNAseq reflected these patterns. Birthdating analyses in mice and non-human
324 primates have reached contradictory conclusions about whether inhibitory neurons follow inside-
325 out patterns of generation similar to excitatory neurons (44, 45), though recent analyses in mice
326 suggest that previous contradictions may reflect the convolution of multiple patterns of
327 generation that may be subtype specific (46). We found that MGE-derived pPVALB+ subtype
328 neurons, enriched in layer IV-VI, showed if anything a trend for the latest-generated neurons to
329 show markers of deeper layers (Fig. 5C, D). SST+ neurons, widely distributed in layer II-VI,
330 similarly did not show an inside-out pattern detectable with the mutations and cells analyzed (Fig.
331 5C, D). We robustly detected SST+ neurons with expression of layer I markers in human PFC
332 (SST-like subclass) (Fig. 5C, D), consistent with observations in MTG (5, 47) and in mice,
333 where such layer I SST+ expressing cells are rare but present (43, 47). These upper layer, CGE-
334 derived SST-like cells are a subclass of LAMP5+ interneurons that are more transcriptionally
335 related to VIP neurons than MGE derived SST+ interneurons, though they lack VIP expression
336 (5, 47). Our data further confirm that LAMP5+ interneurons express markers suggesting broad
337 laminar location, but also did not reveal a simple inside-out progression of formation (5).
338 Interestingly, we observed a substantial proportion of LAMP5+ inhibitory neurons, particularly
339 the SST-like class, labeled by later mutations, indicating that this subtype may be generated later
340 during development than other inhibitory cell types (Fig. 5C, D). Overall, our findings suggest
341 little evidence of the inside-out patterns of neurogenesis demonstrated by excitatory neurons, but
342 also show that detailed analysis of interneurons will likely require deep datasets of sSNV
343 occurring at late stages of interneuron development, and higher-throughput methods of analysis.

344

345 **Discussion**

346 We have developed scMH and PRDD-seq that allowed us, to our knowledge, the first
347 simultaneous analysis of cell lineage and transcriptional cell type in human brain—and
348 potentially, any mammalian brain—through improved identification of sSNVs in deep bulk and
349 single-cell sequencing data. Our analysis of a single cortical area (PFC) in two individual brains
350 revealed some conserved patterns of cell lineage compared to nonhumans, including that
351 inhibitory and excitatory neurons diverge early in humans, and that excitatory neurons form
352 following a similar “inside-out” order as seen in the animal models. However, PRDD-seq also
353 provides the first quantitative estimate in any species of number of progenitor cells
354 (approximately 10) that generate the excitatory neurons in a given cortical area. Furthermore,
355 PRDD-seq also provided some direct insight into inhibitory neuron development in humans,
356 supporting parallel development of different subtypes of inhibitory neurons, with spatial and
357 temporal associations specific only to some subtypes. Our data show that, as methods improve to
358 capture sSNVs present in small numbers of cells, the natural occurrence of sSNVs with each cell
359 division (13, 14, 17) is likely sufficient to provide a very rich map of cell lineage patterns in any
360 given postmortem human brain.

361 The human cerebral cortex has been thought to contain approximately 80% excitatory
362 glutamatergic neurons and 20% GABAergic interneurons (48), although recent scRNAseq
363 studies have reported a somewhat lower ratio of about 70% excitatory neurons (*SI Appendix*,
364 Table S3) (5, 49, 50). Although our PRDD-seq analysis showed 661 excitatory versus 566
365 inhibitory PRDD-seq cells in total for UMB1465 and UMB4638, which represents 54%
366 excitatory neurons (*SI Appendix*, Table S3), this higher proportion of inhibitory neurons seems to
367 reflect either aspects of the tissue (which was stored for long periods frozen), or our NeuN+-
368 sorting method, since similar ratios are seen in 10X scRNAseq from the one brain analyzed (*SI*

369 *Appendix*, Table S3). On the other hand, PRDD-seq cells are studied as containing at least one
370 sSNV identified from scMH using a small number of deeply sequenced neuronal nuclei isolated
371 from the same region, and so do not represent an unbiased sampling of the human brain region.
372 Nonetheless, the fact that we can assign 60-70% of all excitatory neurons to clades in UMB1465,
373 and that neurons with identified SNVs represent most major neuronal types in scRNAseq (Fig.
374 3E), suggests that our sampling has captured the majority of the lineage of the cortical patch,
375 although rare lineages are likely to be missed without much deeper sequencing. Moreover, the
376 presence of 6-7 explicitly marked clades, and the ability to correlate the allele frequency of a
377 sSNV to the excitatory-restriction of the cells carrying that sSNV, allows two independent
378 quantitative assessments of how many progenitors (approximately 10) contribute to the neurons
379 of the patch of cortex from which neurons were isolated, illustrating the remarkable quantitative
380 potential of this approach.

381 Since occasional dropout of DNA marks and RNA markers in PRDD-seq is unavoidable,
382 limited by the quality of isolated nuclei, we emphasize that our results are most robust when
383 analyzing cells positive for both. The quality of postmortem brain tissues can influence the
384 integrity of both genomic DNA and mRNA. Regarding DNA, since no whole-genome
385 amplification is performed prior to targeted pre-amplification, only a single molecular copy of
386 each allele is available for genotyping of each sSNV, so occasional dropout is inevitable.
387 However, our lineage strategy is based not only on the presence of clade-specific sSNVs but also
388 the absence of many sSNVs from other clades (Fig. 3A), so the chance for mis-assigning cells
389 should be relatively small. Nevertheless, mapping our sSNVs onto our scRNAseq dataset
390 suggests that lineage marks are present in the major neuronal subtypes, although rare neuronal
391 types are likely to be missed given our modest sample size. Regarding RNA, single nuclei from

392 postmortem human brain contains only a small amount of mRNAs. Fluidigm Biomark assays are
393 microfluidics-based qPCR assays that are sensitive to subtle changes of the input or environment.
394 As a result, we observed a 30.4% dropout rate of DNA markers and similar level of dropout of
395 RNA marker dropout. However, since PRDD-seq analyses excluded these dropout events, and
396 were completely based on the relative cell type proportions across different stages within one
397 lineage, we have no reason to think that the dropouts are systematic with respect to cell type with
398 one exception: the relatively larger proportion of pPVALB⁺ neurons in PRDD-seq than
399 scRNAseq, likely reflecting the failure of some probes for SST, VIP, and LAMP5. Better and
400 richer probe sets are likely to be able to resolve this in the future.

401 There are limitations to our analysis, since we are analyzing a small sample of the vast
402 size of the human brain, and PRDD-seq is relatively low-throughput and expensive, so our initial
403 analysis only can make conclusions about relatively common cell types. The present analysis is
404 somewhat limited in the analysis of late mutations present in 1% of cells, especially interneurons,
405 since it is challenging to detect those mutations with great sensitivity, but will await single-cell
406 studies on subtypes of neurons in the future. On the other hand, the combined analysis of sSNVs
407 and cell types is archival and progressive. The vast size of the human brain means that each
408 subsequent round of DNA sequencing—whether of bulk tissue or of single or pooled cells—adds
409 to the total depth of sequence data, and provides progressively richer information about late
410 sSNVs. Indeed, the likely dispersed nature of inhibitory clones suggests that analyzing one
411 cortical region could provide sequence data useful in the analysis of a completely different
412 cortical region for these cell types.

413 Overall, PRDD-seq has many advantages even beyond the quantitative analysis of
414 lineages and mosaic fractions that we begin to illustrate here. Since the method uses sSNVs as

415 lineage marks, it is inherently genomic and so allows correlation not only of normal
416 developmental patterns, but would immediately capture alterations to lineage patterns caused by
417 function-altering germline or somatic mutations. In addition, since sSNVs serve as *in vivo*
418 cellular markers for drawing a developmental lineage map without any transgenic manipulation
419 as demonstrated in this study, the method promises to be applicable in principle to any species or
420 human disease condition for which post-mortem brain is available.

421

422 **Materials and Methods**

423 **Human tissues whole-genome sequencing**

424 Frozen post-mortem tissues from three neurologically normal individuals, UMB1465 (a 17-year-
425 old male), UMB4638 (a 15-year-old female), and UMB4643 (a 42-year-old female), were
426 obtained from the NIH NeuroBioBank at the University of Maryland, and prepared according to
427 a standardized protocol (<http://medschool.umaryland.edu/btbank/method2.asp>) under the
428 supervision of the NIH NeuroBioBank ethical guidelines. UMB1465 and UMB4638 died of
429 injuries sustained in motor vehicle accidents, while UMB4643 died of cardiovascular disease.
430 Bulk DNA samples and single neuronal nuclei amplified by multiple displacement amplification
431 (MDA) were prepared and whole-genome sequenced by Illumina HiSeq platforms as part of
432 previous studies in our lab (16). The average sequencing depth was about 40X for single neurons
433 and about 200X for bulk brain samples.

434

435 **Estimation of cell-specific dropout rate and error rate**

436 Germline heterozygous mutations were called by GATK HaplotypeCaller (51) from the whole-
437 genome sequencing data from bulk brain DNA samples, and only common SNPs annotated in

438 the 1000 Genome Project (52) were considered to reduce false positive calls. To estimate cell-
439 specific allele dropout rate, we calculate the proportion of germline heterozygous sites that were
440 genotyped as reference- or alternative-homozygous in single-cell sequencing data. One neuron of
441 UMB4643 with significantly lower allele dropout rate (Z -score < -2) was excluded from
442 subsequent analyses, since it likely represented a doublet from FANS sorting. Similarly, we also
443 extracted the reference-homozygous sites at the 3' adjacent position of each germline
444 heterozygous mutation and calculate the proportion of heterozygous and alternative-homozygous
445 genotypes to estimate the genome-wide error rate in each single-cell.

446

447 **Framework of single-cell MosaicHunter**

448 The overall framework of single-cell MosaicHunter (scMH) was illustrated in Fig. 2A. sSNV
449 candidates were first called from the bulk sequencing data using a Bayesian graphical model (20,
450 21), in which the likelihoods of somatic mutation and three genotypes of inherited mutation were
451 calculated with the consideration of binomial sampling variation and base-calling errors (Fig. 2A,
452 left panel). The presence or absence of somatic mutation in each single-cell was then inferred by
453 adapting the likelihood and allele fraction (f) of somatic mutation estimated from bulk sample as
454 prior probability, after controlling the cell-specific allele dropout rate (d) and error rate (e) (Fig.
455 2A, right panel). Specifically, the transition matrix between bulk and single-cell genotypes was
456 developed as below,

$$457 \quad P(G_{sc} | G_{bulk}) = \begin{pmatrix} 1 & 0 & 0 & 1-f \\ 0 & 1 & 0 & f \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

458 where each column denotes reference-homozygous, heterozygous, alternative-homozygous, and
459 mosaic genotype for bulk sequencing, and each row denotes reference-homozygous,

460 heterozygous, and alternative-homozygous genotype for single-cell sequencing. The genotype
461 likelihoods in single-cells were further adjusted for allele dropout rate (d) and error rate (e) as
462 below,

$$463 \quad P(G_{post} | G_{pre}) = \begin{pmatrix} 1-2e+e^2 & d(1-d) & e^2 \\ 2e(1-e) & 1-2d+2d^2 & 2e(1-e) \\ e^2 & d(1-d) & 1-2e+e^2 \end{pmatrix}$$

464 where each column and row denotes reference-homozygous, heterozygous, and alternative-
465 homozygous genotype before and after adjustment for single-cell sequencing. Single-cell
466 genotypes were binarized as mutant or wildtype by comparing the posterior probability of
467 heterozygous genotype to an empirical threshold. For each candidate site, the proportion of
468 mutant cells was calculated to further filter out germline mutations. Candidate sites with >50%
469 cells showing aberrant single-cell allele fractions were also removed to exclude hotspots of
470 technical artifacts. In “bulk-free” mode with synthetic bulk generated from in silico merging
471 sequencing data from multiple single-cells, scMH would only consider sSNVs which were
472 shared by at least two single-cells.

473

474 **Somatic SNV calling and performance comparison**

475 Paired-end reads from bulk and single-cell whole-genome sequencing data were aligned to the
476 GRCh37 human reference genome by BWA (53), and then processed by GATK (51) and Picard
477 (<http://broadinstitute.github.io/picard/>) for the removal of duplicated and error-prone reads, indel
478 realignments, and base-quality recalibrations. sSNVs in neurons of UMB1465, UMB4638, and
479 UMB4643 were called by scMH and four other tools including Monovar (22), SCcaller (23),
480 LiRA (24), and Conbase (25). Sensitivity was estimated as the detected proportion of lineage-
481 informative sSNVs that had been previously identified and validated in these three brain samples

482 (16). False discovery rate (FDR) was measured as the proportion of lineage-informative
483 mutations that were shared by cells from conflicting clades (16).

484

485 **Validation of somatic SNVs**

486 Validation of somatic SNVs called from scMH was performed using PCR of 200-500 bp
487 amplicons including the mutated base, followed by Sanger sequencing. All variants were
488 validated in independently sorted single neuronal nuclei amplified by MDA.

489

490 **Generation of simulated single-cell whole-genome sequencing data**

491 To estimate the sensitivity of scMH to detect lineage-informative sSNVs, bulk and single-cell
492 sequencing data with varied somatic mutation rates was generated *in silico* (SI Appendix, Fig.
493 S4A). First, we developed a simplistic model to mimic the process of early embryogenesis: 1)
494 ten rounds of symmetric cell division was applied to generate 1024 (2^{10}) daughter cells derived
495 from a single zygote, in which somatic mutations was randomly introduced at a rate of 1, 2, 5, or
496 10 mutations per round; 2) each daughter cells accumulated cell-specific somatic mutations for
497 another ten rounds with the same mutation rate. Then, for each daughter cell, sequencing reads of
498 chromosome 1 was generated at 40X by ART (54) with default parameters for Illumina
499 platforms, and then germline mutations identified from NA12878 and somatic mutations
500 generated by our model was introduced to the sequencing read using BAMSurgeon (55), with an
501 allele dropout rate of 1×10^{-2} per base and MDA amplification rate of 1×10^{-7} per base that were
502 estimated from real single-cell sequencing data. Finally, we randomly selected 80 cells
503 (consistent with the detection threshold of scMH in real brain bulk samples) from the 1024
504 daughter cells and merged their sequencing data with a down-sampling of 200X to generate the

505 bulk sequencing data, and another 16 cells was randomly selected for benchmarking the
506 performance of scMH. Our simulation data suggested that scMH was able to detect, on average,
507 67% and 86% of cell-shared sSNVs with PFC bulk or synthetic bulk, respectively (*SI Appendix*,
508 Fig. S4B).

509

510 **Design and selection of Taqman genotyping and gene expression probes**

511 Taqman genotyping probes for all validated sSNVs were designed using custom Taqman assay
512 design tool provided Thermo Fisher Scientific. Off-the-shelf Taqman gene expression probes
513 were ordered from Thermo Fisher Scientific. All designed probes were tested by ddPCR using
514 human genomic DNA (Human male, Promega) as a negative control. Gene expression probes
515 were further tested by isolated bulk brain RNA as a positive control. Genotyping probes were
516 also tested by comparing the detected mosaic fractions and the fractions calculated from bulk
517 sequencing (Fig. 3C, D).

518

519 **Parallel RNA and DNA analysis after Deep-sequencing (PRDD-seq)**

520 Single nuclei from postmortem brain samples were isolated using fluorescence-activated nuclear
521 sorting (FANS) for NeuN as described previously (56). Isolated single neuronal nuclei were
522 directly sorted into CellsDirect One-Step qRT-PCR (Thermo Fisher Scientific) pre-amplification
523 buffers containing 0.14x Taqman gene expression assays and SNP genotyping assays. Pre-
524 amplification of all cDNA and genomic DNA amplicons were performed directly after the FANS
525 sorting. Following pre-amplification, samples were diluted 10-fold and loaded onto 96.96
526 genotyping or 192.24 gene expression dynamic assay integrated fluidic circuits for standard
527 amplification per manufacturer's instructions (Biomark, Fluidigm). Genotype and gene

528 expression were further determined by Biomark machine and analyzed by Biomark & EP1
529 software (Fluidigm).

530

531 **10X Genomics preparation and sequencing**

532 Standard 10X Genomics Chromium 3' (v2 chemistry) was carried out according to the
533 manufacturer's recommendation. Single nuclei from postmortem brain samples were isolated
534 using FANS for NeuN, and were loaded onto a 10X Genomics Chromium chip. Reverse
535 transcription and library preparation was performed using the 10X Genomics Single Cell v2 kit
536 following the 10X Genomics protocol. The library was then sequenced on one lane of Illumina
537 NextSeq-500 with a high-output kit.

538

539 **Single-cell RNA sequencing analysis**

540 The expression matrix of 10X Genomic single-cell RNA sequencing (scRNAseq) was generated
541 by Cell Ranger following the recommended protocols. The expression matrix and cell
542 annotations of SMART-seq-based scRNAseq for human MTG (5) was downloaded from the
543 website (<https://celltypes.brain-map.org/rnaseq/>). Variance normalization, clustering and
544 visualization were performed by Pagoda2 (57) using a similar protocol to Lake et al (7). Cell
545 clusters containing more than 50 cells were plotted on the t-SNE map, and the annotation of
546 cortical layer (upper, middle, lower) for excitatory neurons and subtypes for inhibitory neurons
547 was manually curated for each cluster according to the expression level of marker genes (*SI*
548 *Appendix*, Fig. S1 and S2). Considering that Layer 1 dissections of MTG nuclei included the
549 upper part of Layer 2 and the absence of excitatory neurons in the Layer 1 of MTG based on in
550 situ labeling (5), all the MTG Layer 1 excitatory neurons were re-annotated as Layer 2. To

551 further compare the expression profile of cells clusters between two scRNAseq datasets, we
552 calculated the cosine similarity of average expression level for marker genes (S) between any
553 two cell clusters. Cell clusters were then hierarchically clustered using the Ward's method with a
554 distance of $1 - S$.

555

556 **Joint analysis of PRDD-seq and scRNAseq cells**

557 To understand the cell type and cell origin of PRDD-seq cells, we utilized their gene expression
558 profiles to map them onto the t-SNE maps of scRNAseq. PRDD-seq cells were firstly separated
559 into excitatory or inhibitory neurons according to the expression of excitatory or inhibitory
560 marker genes (*SI Appendix*, Table S2), and cells with no or conflicting expression of these
561 marker genes were excluded. For excitatory neurons, missing expression status for layer marker
562 genes (*SI Appendix*, Table S2) were inferred if any layer-specific genes for a given layer were
563 expressed. The cosine similarity matrix was then generated by comparing PRDD-seq cells
564 against scRNAseq cells. For each PRDD-seq cell, its cell cluster was determined by the majority
565 voting among its 25-nearest scRNAseq cells in cosine similarity (Fig. 1I), and the cell type and
566 cortical layer information of PRDD-seq cell was further annotated based on their assigned cell
567 cluster in scRNAseq datasets. To benchmark how accurately we could infer cell type and layer
568 annotation from the 30 marker genes profiled in PRDD-seq cells, we randomly sampled 200
569 scRNAseq cells from each of the seven cell types (upper, middle, lower layer excitatory neurons
570 and VIP+, SST+, LAMP5+, pPVALB+ inhibitory neurons) from 10X Genomic dataset and each
571 of the six dissected cortical layers from SMART-seq dataset, and only extracted the expression
572 profiles of 30 marker genes from each scRNAseq cell. Using the same majority voting strategy,
573 we assigned them back to cell clusters on the t-SNE map. As shown in Fig. 3D, the majority of

574 the randomly sampled cells can be correctly assigned to their original cell type and layer
575 annotation, suggesting the accuracy of our mapping strategy in PRDD-seq.

576

577 **Quantification and statistical analysis**

578 All data are reported as mean \pm 95% confident interval (CI) unless mentioned otherwise. All of
579 the statistical details can be found in the figure legends, figures, and Results. Significance was
580 defined for p values smaller than 0.05. All tests were performed using the R software package
581 (version 3.5.0).

582

583 **Data and code availability**

584 Sequencing data was deposited in the NCBI SRA with accession numbers SRP041470 and
585 SRP061939. MosaicHunter is publicly available at <http://mosaichunter.cbi.pku.edu.cn/>. Config
586 files of single-cell MosaicHunter (scMH) and other scripts about PRDD-seq can be accessed at
587 <https://github.com/AugustHuang/PRDD-seq>.

588

589 References

- 590 1. P. Rakic, Evolution of the neocortex: a perspective from developmental biology. *Nat Rev*
591 *Neurosci* **10**, 724-735 (2009).
- 592 2. D. H. Geschwind, P. Rakic, Cortical evolution: judge the brain by its cover. *Neuron* **80**,
593 633-647 (2013).
- 594 3. M. Heide, K. R. Long, W. B. Huttner, Novel gene function and regulation in neocortex
595 expansion. *Curr Opin Cell Biol* **49**, 22-30 (2017).
- 596 4. C. S. Raju *et al.*, Secretagogin is Expressed by Developing Neocortical GABAergic
597 Neurons in Humans but not Mice and Increases Neurite Arbor Size and Complexity.
598 *Cereb Cortex* **28**, 1946-1958 (2018).
- 599 5. R. D. Hodge *et al.*, Conserved cell types with divergent features in human versus mouse
600 cortex. *Nature* **573**, 61-68 (2019).
- 601 6. S. Zhong *et al.*, A single-cell RNA-seq survey of the developmental landscape of the
602 human prefrontal cortex. *Nature* **555**, 524-528 (2018).
- 603 7. B. B. Lake *et al.*, Integrative single-cell analysis of transcriptional and epigenetic states in
604 the human adult brain. *Nat Biotechnol* **36**, 70-80 (2018).
- 605 8. A. McKenna *et al.*, Whole-organism lineage tracing by combinatorial and cumulative
606 genome editing. *Science* **353**, aaf7907 (2016).
- 607 9. K. L. Frieda *et al.*, Synthetic recording and in situ readout of lineage information in
608 single cells. *Nature* **541**, 107-111 (2017).
- 609 10. B. Raj *et al.*, Simultaneous single-cell profiling of lineages and cell types in the
610 vertebrate brain. *Nat Biotechnol* **36**, 442-450 (2018).
- 611 11. B. Spanjaard *et al.*, Simultaneous lineage tracing and cell-type identification using
612 CRISPR-Cas9-induced genetic scars. *Nat Biotechnol* **36**, 469-473 (2018).
- 613 12. A. Rodriguez-Meira *et al.*, Unravelling Intratumoral Heterogeneity through High-
614 Sensitivity Single-Cell Mutational Analysis and Parallel RNA Sequencing. *Mol Cell* **73**,
615 1292-1305 e1298 (2019).
- 616 13. Y. S. Ju *et al.*, Somatic mutations reveal asymmetric cellular dynamics in the early
617 human embryo. *Nature* **543**, 714-718 (2017).
- 618 14. T. Bae *et al.*, Different mutational rates and mechanisms in human cells at pregastrulation
619 and neurogenesis. *Science* **359**, 550-555 (2018).
- 620 15. S. De, Somatic mosaicism in healthy human tissues. *Trends Genet* **27**, 217-223 (2011).
- 621 16. M. A. Lodato *et al.*, Somatic mutation in single human neurons tracks developmental and
622 transcriptional history. *Science* **350**, 94-98 (2015).
- 623 17. M. A. Lodato *et al.*, Aging and neurodegeneration are associated with increased
624 mutations in single human neurons. *Science* **359**, 555-559 (2018).
- 625 18. L. S. Ludwig *et al.*, Lineage Tracing in Humans Enabled by Mitochondrial Mutations and
626 Single-Cell Genomics. *Cell* **176**, 1325-1339 e1322 (2019).
- 627 19. Y. Dou *et al.*, Accurate detection of mosaic variants in sequencing data without matched
628 controls. *Nat Biotechnol* **38**, 314-319 (2020).
- 629 20. A. Y. Huang *et al.*, MosaicHunter: accurate detection of postzygotic single-nucleotide
630 mosaicism through next-generation sequencing of unpaired, trio, and paired samples.
631 *Nucleic Acids Res* **45**, e76 (2017).
- 632 21. A. Y. Huang *et al.*, Postzygotic single-nucleotide mosaicism in whole-genome
633 sequences of clinically unremarkable individuals. *Cell Res* **24**, 1311-1327 (2014).

- 634 22. H. Zafar, Y. Wang, L. Nakhleh, N. Navin, K. Chen, Monovar: single-nucleotide variant
635 detection in single cells. *Nat Methods* **13**, 505-507 (2016).
- 636 23. X. Dong *et al.*, Accurate identification of single-nucleotide variants in whole-genome-
637 amplified single cells. *Nat Methods* **14**, 491-493 (2017).
- 638 24. C. L. Bohrsen *et al.*, Linked-read analysis identifies mutations in single-cell DNA-
639 sequencing data. *Nat Genet* **51**, 749-754 (2019).
- 640 25. J. Hard *et al.*, Conbase: a software for unsupervised discovery of clonal somatic
641 mutations in single cells through read phasing. *Genome Biol* **20**, 68 (2019).
- 642 26. A. Y. Huang *et al.*, Distinctive types of postzygotic single-nucleotide mosaicisms in
643 healthy individuals revealed by genome-wide profiling of multiple organs. *PLoS Genet*
644 **14**, e1007395 (2018).
- 645 27. T. Helleday, S. Eshtad, S. Nik-Zainal, Mechanisms underlying mutational signatures in
646 human cancers. *Nat Rev Genet* **15**, 585-598 (2014).
- 647 28. C. S. von Bartheld, J. Bahney, S. Herculano-Houzel, The search for true numbers of
648 neurons and glial cells in the human brain: A review of 150 years of cell counting. *J*
649 *Comp Neurol* **524**, 3865-3895 (2016).
- 650 29. H. Markram *et al.*, Interneurons of the neocortical inhibitory system. *Nat Rev Neurosci* **5**,
651 793-807 (2004).
- 652 30. S. A. Anderson, D. D. Eisenstat, L. Shi, J. L. Rubenstein, Interneuron migration from
653 basal forebrain to neocortex: dependence on *Dlx* genes. *Science* **278**, 474-476 (1997).
- 654 31. G. Fishell, C. A. Mason, M. E. Hatten, Dispersion of neural progenitors within the
655 germinal zones of the forebrain. *Nature* **362**, 636-638 (1993).
- 656 32. S. Anderson, M. Mione, K. Yun, J. L. Rubenstein, Differential origins of neocortical
657 projection and local circuit neurons: role of *Dlx* genes in neocortical interneuronogenesis.
658 *Cereb Cortex* **9**, 646-654 (1999).
- 659 33. T. Ma *et al.*, Subcortical origins of human and monkey neocortical interneurons. *Nat*
660 *Neurosci* **16**, 1588-1597 (2013).
- 661 34. D. V. Hansen *et al.*, Non-epithelial stem cells and cortical interneuron production in the
662 human ganglionic eminences. *Nat Neurosci* **16**, 1576-1587 (2013).
- 663 35. P. Rakic, Neurons in rhesus monkey visual cortex: systematic relation between time of
664 origin and eventual disposition. *Science* **183**, 425-427 (1974).
- 665 36. J. B. Angevine, Jr., R. L. Sidman, Autoradiographic study of cell migration during
666 histogenesis of cerebral cortex in the mouse. *Nature* **192**, 766-768 (1961).
- 667 37. P. Gao *et al.*, Deterministic progenitor behavior and unitary production of neurons in the
668 neocortex. *Cell* **159**, 775-788 (2014).
- 669 38. C. Mayer *et al.*, Developmental diversification of cortical inhibitory interneurons. *Nature*
670 **555**, 457-462 (2018).
- 671 39. M. Turrero Garcia, E. Mazzola, C. C. Harwell, Lineage Relationships Do Not Drive
672 MGE/PoA-Derived Interneuron Clustering in the Brain. *Neuron* **92**, 52-58 (2016).
- 673 40. K. N. Brown *et al.*, Clonal production and organization of inhibitory interneurons in the
674 neocortex. *Science* **334**, 480-486 (2011).
- 675 41. C. B. Reid, S. F. Tavazoie, C. A. Walsh, Clonal dispersion and evidence for asymmetric
676 cell division in ferret cortex. *Development* **124**, 2441-2450 (1997).
- 677 42. M. J. Nigro, Y. Hashikawa-Yamasaki, B. Rudy, Diversity and Connectivity of Layer 5
678 Somatostatin-Expressing Interneurons in the Mouse Barrel Cortex. *J Neurosci* **38**, 1622-
679 1633 (2018).

- 680 43. B. Rudy, G. Fishell, S. Lee, J. Hjerling-Leffler, Three groups of interneurons account for
681 nearly 100% of neocortical GABAergic neurons. *Dev Neurobiol* **71**, 45-61 (2011).
- 682 44. E. S. Ang, Jr., T. F. Haydar, V. Gluncic, P. Rakic, Four-dimensional migratory
683 coordinates of GABAergic interneurons in the developing mouse cortex. *J Neurosci* **23**,
684 5805-5815 (2003).
- 685 45. V. V. Rymar, A. F. Sadikot, Laminar fate of cortical GABAergic interneurons is
686 dependent on both birthdate and phenotype. *J Comp Neurol* **501**, 369-380 (2007).
- 687 46. S. M. Kelly, R. Raudales, M. Moissidis, G. Kim, Z. J. Huang, Multipotent radial glia
688 progenitors and fate-restricted intermediate progenitors sequentially generate diverse
689 cortical interneuron types. *bioRxiv* 10.1101/735019 (2019).
- 690 47. E. Boldog *et al.*, Transcriptomic and morphophysiological evidence for a specialized
691 human cortical GABAergic cell type. *Nat Neurosci* **21**, 1185-1195 (2018).
- 692 48. C. P. Wonders, S. A. Anderson, The origin and specification of cortical interneurons. *Nat*
693 *Rev Neurosci* **7**, 687-696 (2006).
- 694 49. N. Habib *et al.*, Massively parallel single-nucleus RNA-seq with DroNc-seq. *Nat*
695 *Methods* **14**, 955-958 (2017).
- 696 50. B. B. Lake *et al.*, Neuronal subtypes and diversity revealed by single-nucleus RNA
697 sequencing of the human brain. *Science* **352**, 1586-1590 (2016).
- 698 51. M. A. DePristo *et al.*, A framework for variation discovery and genotyping using next-
699 generation DNA sequencing data. *Nat Genet* **43**, 491-498 (2011).
- 700 52. C. Genomes Project *et al.*, A global reference for human genetic variation. *Nature* **526**,
701 68-74 (2015).
- 702 53. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler
703 transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 704 54. W. Huang, L. Li, J. R. Myers, G. T. Marth, ART: a next-generation sequencing read
705 simulator. *Bioinformatics* **28**, 593-594 (2012).
- 706 55. A. D. Ewing *et al.*, Combining tumor genome simulation with crowdsourcing to
707 benchmark somatic single-nucleotide-variant detection. *Nat Methods* **12**, 623-630 (2015).
- 708 56. G. D. Evrony *et al.*, Cell lineage analysis in human brain using endogenous retroelements.
709 *Neuron* **85**, 49-59 (2015).
- 710 57. J. Fan *et al.*, Characterizing transcriptional heterogeneity through pathway and gene set
711 overdispersion analysis. *Nat Methods* **13**, 241-244 (2016).
- 712

713 **Acknowledgements**

714 We thank R. Mattieu, K. Brownstein, J. Li, Flow Cytometry Facility in Boston Children's
715 Hospital, BCH IDDRC Molecular Genetics Core Facility, and the Research Computing group at
716 Harvard Medical School for assistance. We thank G. Fishell, C. Harwell, and F. Vaccarino for
717 comments on the manuscript. Human tissue was obtained from the NIH NeuroBioBank at the
718 University of Maryland and Autism BrainNet, and we thank the donors and their families for
719 their invaluable donations for the advancement of science. P.L is a Howard Hughes Medical
720 Institute – Helen Hay Whitney Foundation Fellow. C.A.W. is supported by the Manton Center
721 for Orphan Disease Research, the Allen Discovery Center program through The Paul G. Allen
722 Frontiers Group, grants from the NINDS (R01NS032457 and U01MH106883), and grant
723 U01MH106883 from the NIMH. C.A.W. is an Investigator of the Howard Hughes Medical
724 Institute.

725

726

727 **Figure 1. PRDD-seq enables simultaneous assessment of cell identity and lineage in single**

728 **cells.** A. Neuronal nuclei from postmortem human brain were based on NeuN+ immunoreactivity.

729 B. Target-specific one-step RT-qPCR amplification of cDNA and gDNA fragments of interest. C.

730 Single-cell MosaicHunter co-analysis of single-cell and bulk deep sequencing data to identify

731 lineage-informative somatic SNVs. D. Multiplex analysis of the amplified cDNA and gDNA

732 fragments to genotype the somatic SNVs and profile 30 cell type-specific markers of gene

733 expression. E. 10X Genomics scRNAseq was performed on NeuN+ nuclei isolated from the

734 same PFC region. F. 21 cell clusters were identified based on 10X Genomics gene expression

735 data, and then divided into upper, middle, and lower layer of excitatory neurons and four

736 subtypes of inhibitory neurons. G. A second scRNAseq dataset (5) performed on nuclei isolated

737 from the MTG region of another post-mortem healthy human brain was also analyzed where

738 layer information was identified based on layer micro-dissection. Cell types were identified

739 based on gene expression data. H. Transcriptional clustering revealed similar single-cell

740 expression profiles between 10X Genomics PFC and SMART-seq MTG scRNAseq datasets.

741 Cell clusters were color-coded to denote different cell type annotation, and clusters derived from

742 10X Genomics PFC (triangle) and SMART-seq MTG (circle) in general clustered by cell type

743 but not by platform. I. Each PRDD-seq cell was mapped to the t-SNE maps by the cosine

744 similarity of gene expression to scRNAseq cells, and then assigned cell type and dissected layer

745 accordingly by majority voting of 25 nearest neighbors. J-K. A combination of genotype and

746 gene expression information of PRDD-seq cells allowed lineage and birthdate analysis of

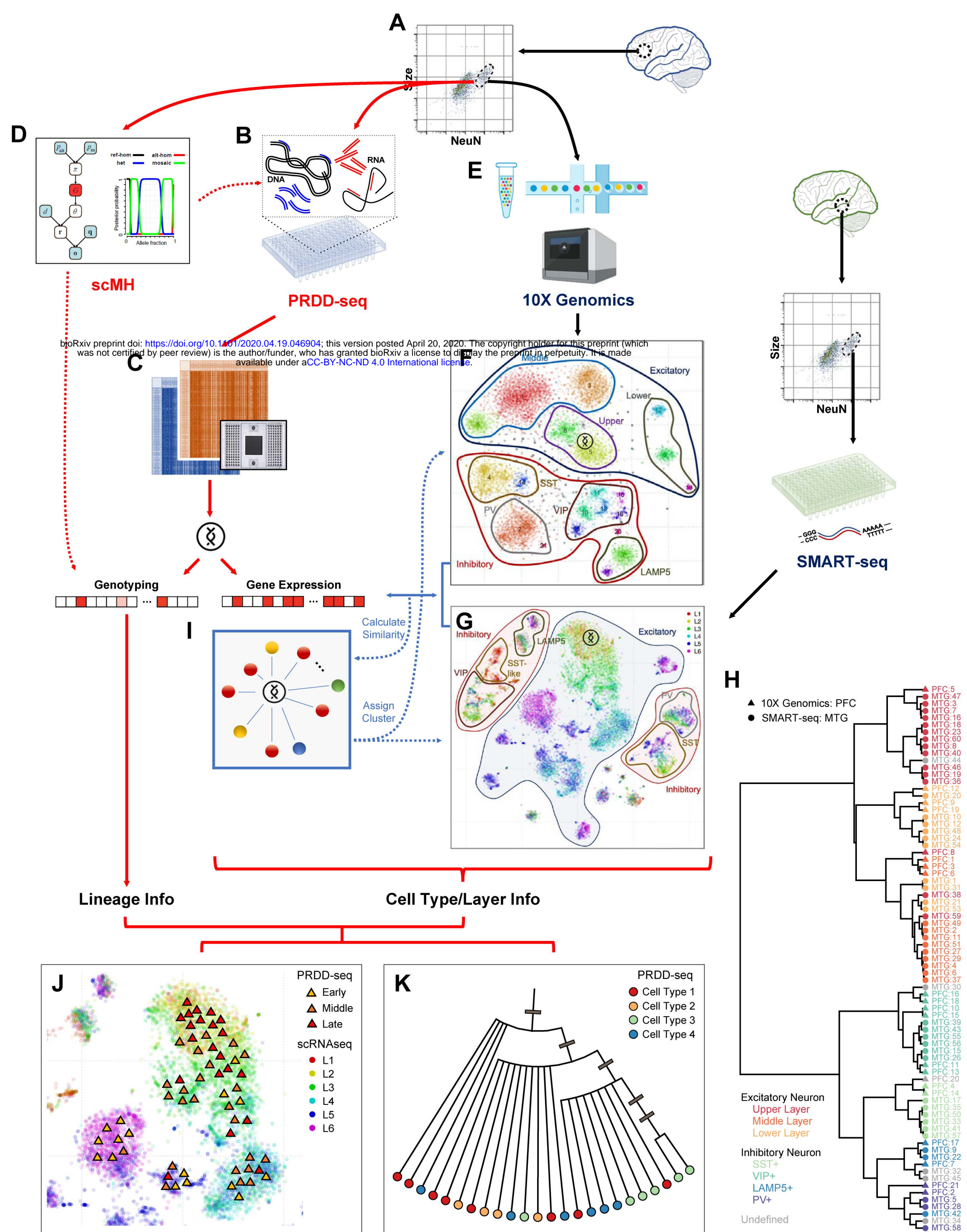
747 particular cell types/layers (J), and cell type differentiation analysis of particular lineage

748 reconstructed by somatic mutations (K). Colored triangles in (I) indicate PRDD-seq cells. Gray

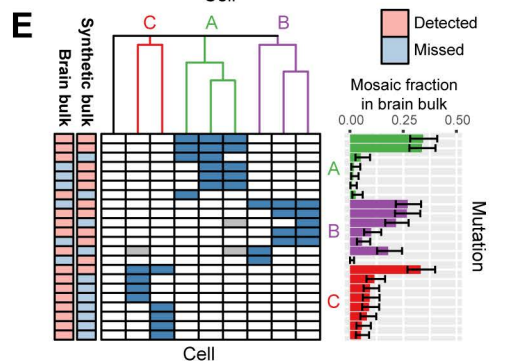
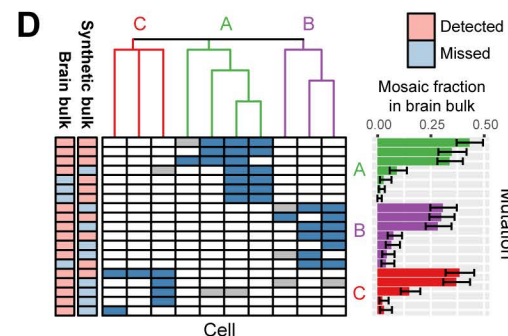
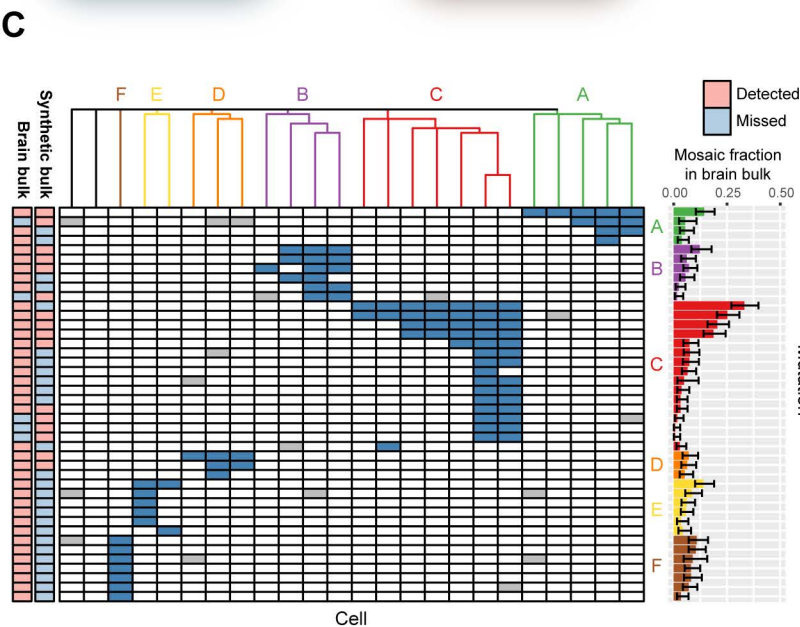
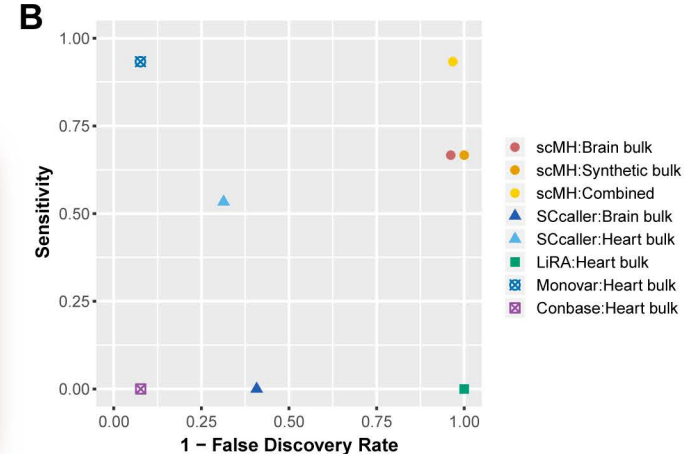
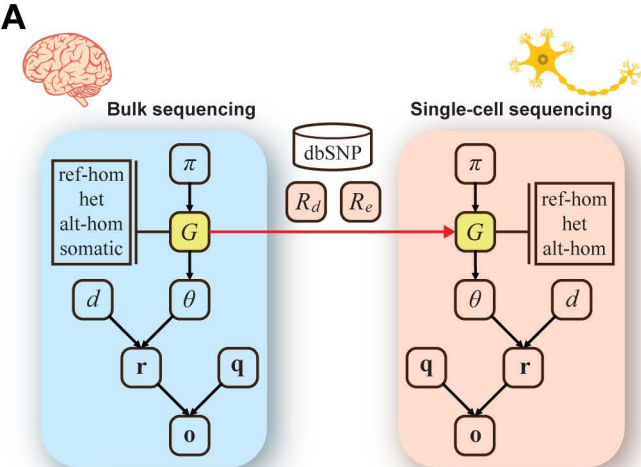
749 bars in (K) indicate occurrences of somatic mutations, whereas all cells in one corresponding

750 sub-clade share the same somatic mutation.

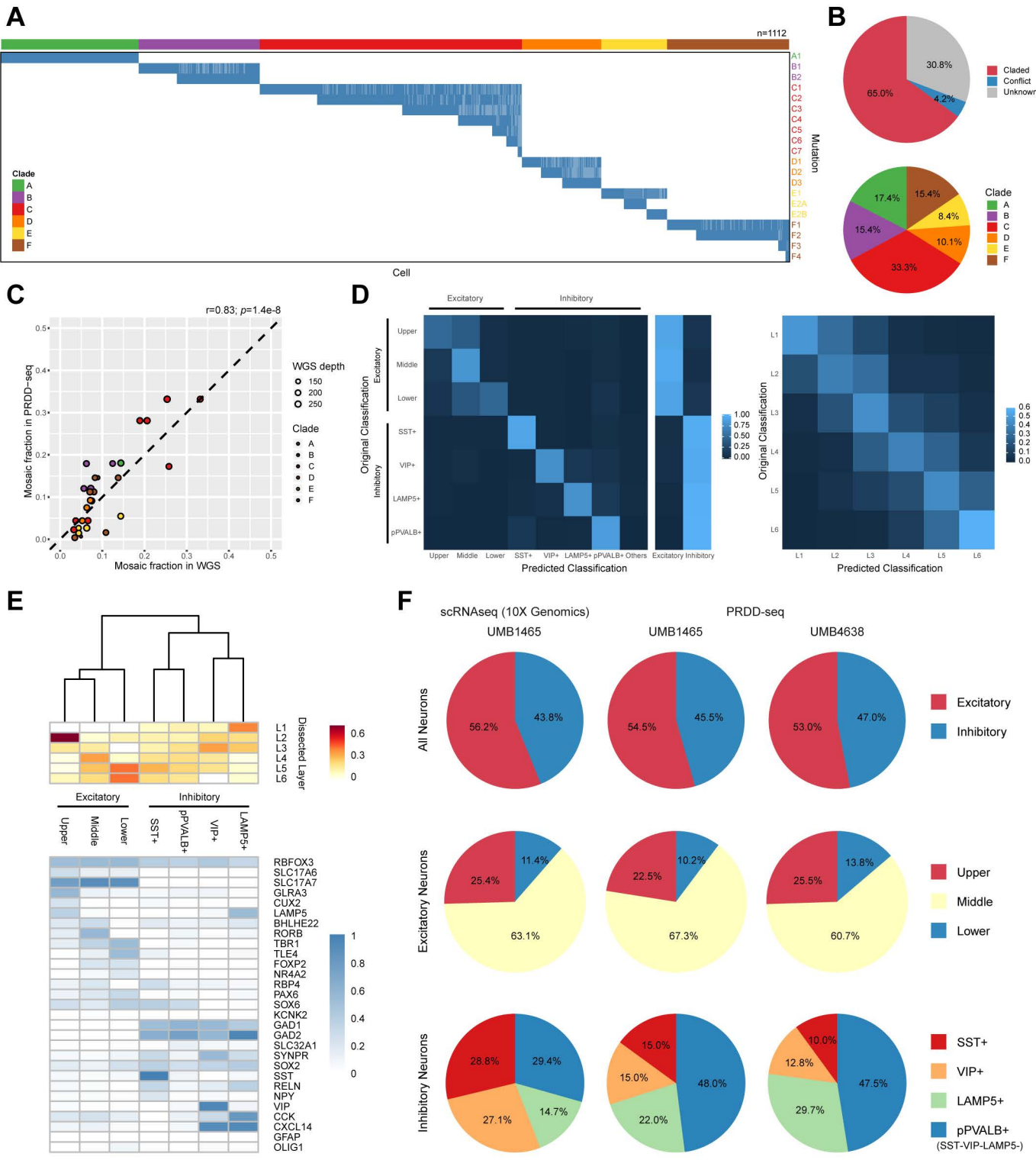
751



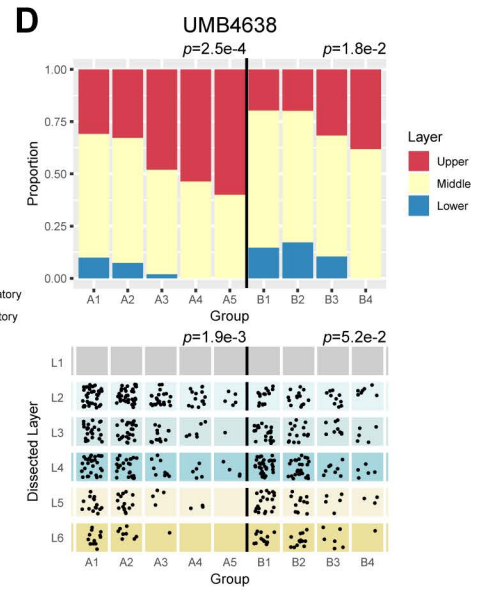
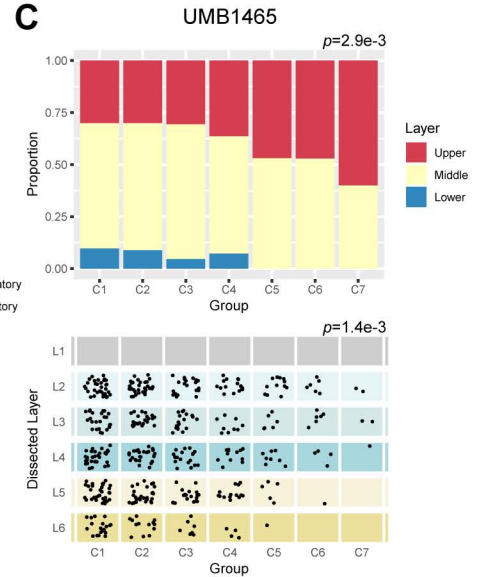
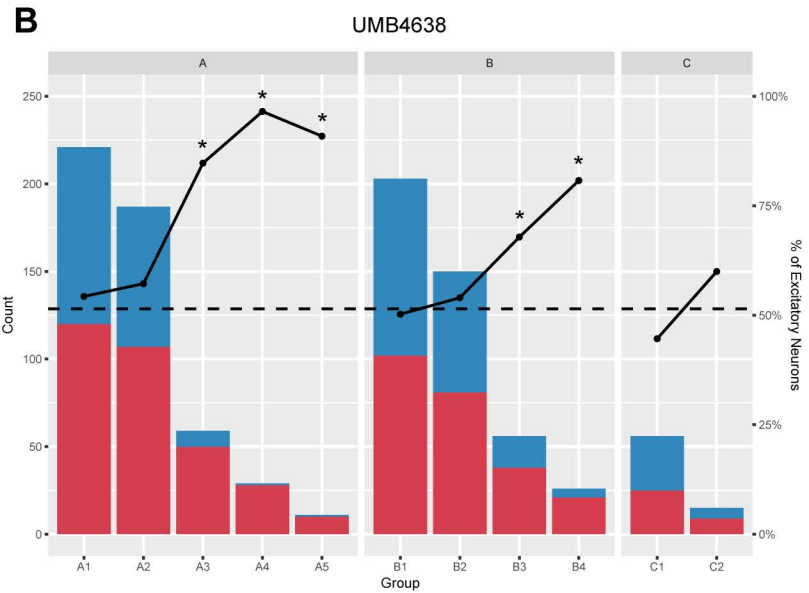
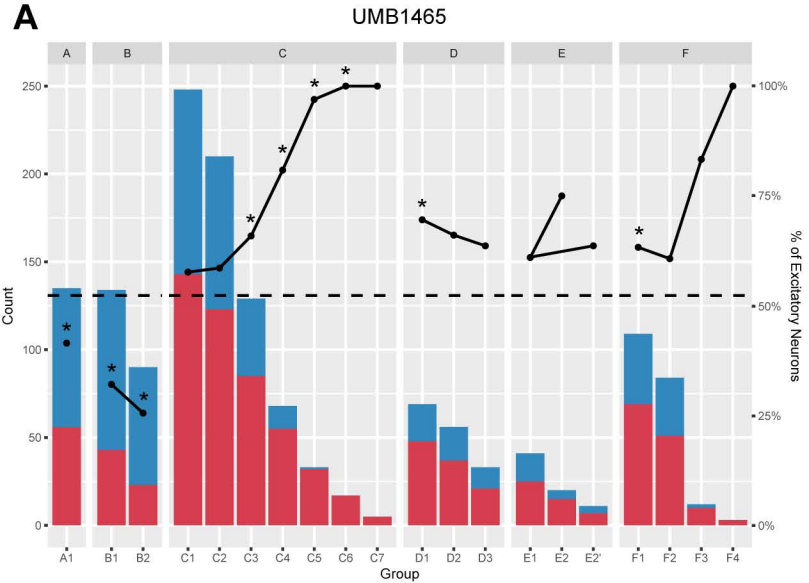
752 **Figure 2. scMH identifies lineage-informative sSNVs from the joint analysis of bulk brain**
753 **and single-neurons.** A. Overview of the extended Bayesian model of scMH to use bulk
754 sequencing data to facilitate sSNV calling in single cells. G denotes the genotype state, π denotes
755 the prior probabilities of genotype, and d , \mathbf{o} , \mathbf{q} denote the depth, observed bases, and their base
756 qualities in bulk or single-cell sequencing data. B. Specificity and precision of identifying sSNVs
757 using scMH and other published callers. scMH outcompeted other callers in both precision and
758 sensitivity. C-E. Validated lineage-informative sSNVs identified by scMH in UMB1465 (C),
759 UMB4638 (D), and UMB4643 (E). Heatmaps demonstrate the genotyping status of sSNVs; dark
760 blue and white squares denote the presence or absence of sSNVs in a given cell, whereas grey
761 squares denote unknown genotype due to locus dropout in single-cell WGS. Bar graphs show the
762 mosaic fraction of each sSNV in WGS of bulk brain sample. Clade E in (C), and clade C in (E),
763 represent likely branching clades where early shared mutations are present, while later sSNVs
764 mark two branches with distinct mutations. Error bars reflect 95% confidence intervals.
765



766 **Figure 3. PRDD-seq profiles single-neurons with varied lineage markers and distinct cell**
767 **type identity.** A. Genotyping results of 30 sSNVs (by rows) from 20 lineages across PRDD-seq
768 cells (by columns) from UMB1465. Blue and white squares represent the presence or absence of
769 sSNV respectively, whereas light blue squares represent the sSNVs that were dropouts in PRDD-
770 seq assay but inferred by the presence of deeper mutations from the same clade. B. Clade
771 classification of PRDD-seq cells profiled in UMB1465. In upper panel, PRDD-seq cells which
772 contained sSNVs from multiple or no clades are labeled as “conflict” and “unknown”
773 respectively. C. Correlation of mosaic fractions from WGS and PRDD-seq (calculated as % of
774 assayed cells carrying a given sSNV) in UMB1465. Both methods showed significantly
775 concordant mosaic fractions (Pearson correlation’s $P < 0.001$). D. Accuracy of cell type (left
776 panel) and cortical layer (right panel) classification based on the expression profile of 30 marker
777 genes used in PRDD-seq. scRNAseq cells from each cell type (10X Genomics) and cortical layer
778 (SMART-seq) were randomly sampled and then re-assigned to clusters of t-SNE map using 30
779 marker genes under PRDD-seq mapping strategy. E. Taxonomy of 3 excitatory layers and 4
780 inhibitory subtypes based on average expression of 30 marker genes in PRDD-seq cells. Relative
781 density of cortical layers for each subgroup is also shown. pPVALB+ denotes PVALB+/SST-
782 VIP- LAMP5- subtype of inhibitory neurons. F. Relative ratio across different cell types of
783 excitatory and inhibitory neurons between PRDD-seq and 10X Genomics scRNAseq.
784



785 **Figure 4. PRDD-seq reveals distinct developmental sequence of excitatory neurons in**
786 **different cortical layers.** A-B. The total number (bar plot) and ratio (dot plot) of excitatory and
787 inhibitory neurons in different lineage clades defined by one or more sSNVs in UMB1465 (A)
788 and UMB4638 (B). Percentage of excitatory neurons increased in later lineage timepoints in
789 clades C and F in UMB1465 and clades A and B in UMB4638. In Clade E of UMB1465, E1
790 branches into two subclades E2A and E2B. Dashed line: average excitatory neuron percentage.
791 Asterisk denotes significantly different excitatory-inhibitory ratio from the average (two-sided
792 one-proportion Z-test's $P < 0.05$). In clades C and F from UMB1465, and clades A and B from
793 UMB4638, later mutations become progressively limited to excitatory neurons. C-D. Layer
794 distributions of excitatory neurons in representative excitatory lineages in UMB1465 (C) and
795 UMB4638 (D), respectively. Layers are determined by mapping PRDD-seq cells onto human
796 PFC scRNAseq (upper panels) or human MTG scRNAseq (lower panels) based on the
797 expression profile similarity of marker genes. In all three illustrated clades, the percentage of
798 upper layer neurons increased while that of lower layer neurons decreased in cells containing
799 sSNVs present at lower mosaic fraction. P -value was calculated by Pearson correlation with
800 ordinal variables.
801



802 **Figure 5. PRDD-seq reveals heterogeneous developmental process for inhibitory neurons.**

803 A-B. Distribution of different subtypes of inhibitory neurons in different lineages in UMB1465
804 (A) and UMB4638 (B), respectively. Major subtypes of inhibitory neurons are widely distributed
805 in different lineages. C-D. Layer distributions of inhibitory subtypes in representative lineages in
806 UMB1465 (C) and UMB4638 (D), respectively. Bar graphs show the proportion of each subtype
807 of neurons in different layers. MGE derived (SST+ and pPVALB+) and CGE derived (VIP+,
808 LAMP5/PAX6+, and SST-like) interneurons showed similar mutation profiles, suggesting that
809 the groups are produced simultaneously. pPVALB+ subtype neurons were enriched in layer IV-
810 VI, while MGE-derived SST+ interneurons showed a similar laminar distribution as pPVALB+
811 interneurons, with no clear evidence of an “inside-out” birth dating pattern. CGE-derived
812 interneurons were broadly distributed across cortical layers, with SST-like cells heavily favoring
813 supragranular layers; LAMP5+, including SST-like cells, were enriched for later lineage marks,
814 suggesting they may be produced later in development than other subtypes.

815

