

1 **Classifying cell cycle states and a quiescent-like G0 state**
2 **using single-cell transcriptomics**

3 Samantha A. O'Connor¹, Leonor Garcia², Rori Hoover¹, Anoop P. Patel^{3,4}, Benjamin B. Bartelle
4¹, Jean-Philippe Hugnot², Patrick J. Paddison⁵, Christopher L. Plaisier^{1,*}

5
6¹ School of Biological and Health Systems Engineering, Arizona State University, Tempe AZ,
7 USA.

8² Institut de Génomique Fonctionnelle, Université de Montpellier, CNRS, INSERM, 141 rue de la
9 Cardonille, 34091, Montpellier, France.

10³ Brotman-Baty Institute for Precision Medicine, University of Washington, Seattle, WA, USA.

11⁴ Department of Neurosurgery, Preston Robert Tisch Brain Tumor Center, Duke University,
12 Durham, NC, USA.

13⁵ Human Biology Division, Fred Hutchinson Cancer Research Center, Seattle WA, USA

14 * To whom correspondence should be addressed: Christopher Plaisier, E-mail:

15 plaisier@asu.edu, Phone: (480) 965-6832, Address: P.O. Box 879709, Tempe, AZ 87287-9709

16

17 **Abstract**

18 Single-cell transcriptomics has unveiled a vast landscape of cellular heterogeneity in which the
19 cell cycle is a significant component. We trained a high-resolution cell cycle classifier (ccAFv2)
20 using single cell RNA-seq (scRNA-seq) characterized human neural stem cells. The ccAFv2
21 classifies six cell cycle states (G1, Late G1, S, S/G2, G2/M, and M/Early G1) and a quiescent-
22 like G0 state (Neural G0), and it incorporates a tunable parameter to filter out less certain
23 classifications. The ccAFv2 classifier performed better than or equivalent to other state-of-the-
24 art methods even while classifying more cell cycle states, including G0. We demonstrate that
25 the ccAFv2 classifier effectively generalizes the S, S/G2, G2/M, and M/Early G1 states across
26 cell types derived from all three germ layers. While the G0, G1, and Late G1 states perform well
27 in neuroepithelial cell types, their accuracy is lower in other cell types. However,
28 misclassifications are confined to the G0, G1, and Late G1 states. We showcased the versatility
29 of ccAFv2 by successfully applying it to classify cells, nuclei, and spatial transcriptomics data in
30 humans and mice, using various normalization methods and gene identifiers. We provide
31 methods to regress the cell cycle expression patterns out of single cell or nuclei data to uncover
32 underlying biological signals. The classifier can be used either as an R package integrated with
33 Seurat or a PyPI package integrated with SCANPY. We proved that ccAFv2 has enhanced
34 accuracy, flexibility, and adaptability across various experimental conditions, establishing
35 ccAFv2 as a powerful tool for dissecting complex biological systems, unraveling cellular
36 heterogeneity, and deciphering the molecular mechanisms by which proliferation and
37 quiescence affect cellular processes.

38 **Introduction**

39 Single-cell RNA sequencing (scRNA-seq) is a robust method for dissecting the transcriptional
40 states of individual cells obtained from specific conditions. These cellular transcriptional states
41 are influenced by various biological signals, including cell type and the phase of the cell cycle.
42 The cell cycle is a tightly regulated and intricately coordinated biological process that
43 orchestrates the division of a cell into two daughter cells. Adult stem cell populations often
44 reside in a quiescent G0 state outside of the cell cycle, reactivating only upon receiving
45 appropriate signals to divide (Doetsch 2003; Obernier et al. 2018). Current state of the art
46 methods to predict cell cycle states based on scRNA-seq transcriptome profiles lump G0 cells
47 with G1 cells (Hao et al. 2021; Zheng et al. 2022; Schwabe et al. 2020; Liu et al. 2017; Hsiao et
48 al. 2020; Scialdone et al. 2015). The grouping of G0 with G1 fails to recognize the clear
49 differences in expression patterns and quiescent phenotype displayed by G0 cells, making them
50 readily distinguishable from G1 cells (O'Connor et al. 2021). The aim of this research is to
51 develop a cell cycle classifier capable of identifying the G0 state in neuroepithelial cells and to
52 determine whether this state can be generalized to other cell types.
53

54 Developing cell cycle classifiers from scRNA-seq transcriptional profiles is challenging due to
55 the scarcity of datasets with experimentally validated ground truth cell cycle labels. Training a
56 classifier requires having example transcriptome profiles labeled with cell cycle states. Previous
57 studies have used Hoechst (DNA stain; (Buettnner et al. 2015) or FUCCI (Leng et al. 2015) to
58 sort embryonic stem cells into G1, S, and G2M subpopulations. However, there are some
59 caveats to these studies. Firstly, these cells were not fixed, meaning they could continue to
60 cycle after sorting and may not be transcriptionally in the same state as they were when they
61 were sorted. Secondly, the markers used for sorting focused on DNA, protein, and post-
62 translational modification abundances, which may not accurately reflect the transcriptional state
63 of the cells. Thirdly, it has been established that embryonic stem cells do not have well-defined
64 G1 or G0 cell cycle states as they quickly transition through cell cycles to produce many cells in
65 the embryo (Ballabeni et al. 2011; White and Dalton 2005). In preliminary analyses it was found
66 that the cell cycle labels were significantly out of alignment (error rates ≥ 0.7) with the
67 transcription states of the cells as determined by ccSeurat (**Supplemental Table S1**), which is
68 the *de facto* standard in the field.
69

70 Previously, we used scRNA-seq of U5 human Neural Stem Cells (U5-hNSCs; Davis and
71 Temple 1994; Johe et al. 1996) grown *in vitro* to discern seven cell cycle states including a
72 quiescent-like G0 state (O'Connor et al. 2021). An Artificial Neural Network (ANN) (Ma and
73 Pellegrini 2020) classifier named the cell cycle ASU/Fred Hutch (ccAF) was trained to predict
74 these seven cell cycle states in cells from new datasets (O'Connor et al. 2021). In those studies,
75 the ccAF classifier was applied to a host of neuroepithelial derived cells characterized by
76 scRNA-seq, including glioblastoma patient tumor cells. The underlying software packages for
77 constructing ANNs (TensorFlow and Keras) have been significantly improved and we
78 hypothesized that reimplementation of the ccAF classifier would significantly improve classifier
79 performance and provide likelihoods for each classification, a feature not available in the original
80 ccAF implementation.
81

82 In addition to the advancements in ANN methodology, numerous new scRNA-seq studies have
83 been conducted that include actively dividing cells. Particularly valuable for assessing the
84 quality and generalizability of the classifier is an atlas of 245,906 cells from 15 different cell
85 types, spanning all three germ layers, derived from human fetal tissue 3 to 12 weeks post-
86 conception (Zeng et al. 2023). A second atlas of developing human spinal cord (Zhang et al.
87 2021) will be used to evaluate whether the classifier can be applied to both single cell and single

88 nuclei RNA-seq (scRNA-seq and snRNA-seq). An atlas of adult neurogenesis in the ventricular-
89 subventricular zone (V-SVZ) (Cebrian-Silla et al. 2021) will be used to demonstrate that the
90 classifier can be applied to mouse cells. It will also allow comparisons to be made between the
91 cell cycle proportions of cell types from adult mouse neurogenesis in the V-SVZ (Cebrian-Silla
92 et al. 2021) and the developing human telencephalon (Nowakowski et al. 2017). Two studies of
93 quiescent neural stem cells will be crucial for demonstrating the identity of the G0 cell state
94 (Llorens-Bobadilla et al. 2015; Dulken et al. 2017). Additionally, we collected scRNA-seq for two
95 IDH mutant low-grade glioma (LGG) cell lines in conditions with and without growth factors. This
96 will allow us to gain insights into the performance of the classifier when confronted with a higher
97 proportion of non-cycling cells. We will also apply the classifier to *in vivo* glioblastoma tumor
98 cells and *in vitro* glioblastoma tumor derived cancer stem cells that were not included in the
99 previous ccAF classifier studies (Couturier et al. 2020). Finally, application of the classifier to a
100 high-resolution spatial-transcriptomics (ST-seq) study of a mouse embryo at 15.5 weeks post-
101 conception (E15.5) will allow us to resolve canonical biological and morphological phenomena
102 for the developmental stage. These datasets offer a robust foundation for rigorously testing and
103 validating the improved ccAF version 2 (ccAFv2) classifier, showcasing its versatility across
104 species, single cells and nuclei, and generalizability across cell types from all germ layers.
105
106 The goal of this research is to develop an improved cell cycle classifier using current state of the
107 art machine learning technology. We aim to demonstrate that the classifier outperforms existing
108 models and generalizes well across various cell types, library preparation methods (scRNA-seq,
109 snRNA-seq, ST-seq), gene annotations, and normalization techniques. Lastly, we aim to
110 provide a classifier with a more user-friendly interface to facilitate its application in future
111 studies.

112 Results

113 Implementation of neural network classifier for ccAFv2

114 We implemented the core algorithm of ccAFv2 to take advantage of significant improvements in
115 machine learning tools that should improve classifier performance and provide likelihoods for
116 each predicted cell cycle classification. The ccAFv2 core algorithm is broken up into two steps.
117 First, the input data is run through the artificial neural network (ANN) to compute likelihoods for
118 each class (i.e., Neural G0, G1, Late G1, S, S/G2, G2/M or M/Early G1; **Figure 1A-B**). The
119 underlying ANN for ccAFv2 starts with a dense input layer connected to two hidden layers that
120 connect to a softmax output layer (**Figure 1A**). Overfitting in the ANN is mitigated by dropout
121 regularization via two dropout layers. The first dropout layer is positioned between the first and
122 second hidden layers and the second dropout layer is between the second hidden layer and the
123 softmax output (**Figure 1A**; Xie et al. 2019). Second, the likelihoods calculated by the ANN for
124 each cell cycle state are used to determine which state should be assigned for each cell (**Figure**
125 **1B**). The cell cycle state with the maximum likelihood is identified and if the likelihood is greater
126 than or equal to the likelihood threshold then the state is returned. Otherwise, if the maximum
127 likelihood is less than the likelihood threshold a state of “Unknown” is returned. These
128 improvements to the core ANN of ccAFv2 will be rigorously tested in the subsequent sections.

129 Training the ccAFv2 classifier

130 The training data for ccAFv2 is comprised of scRNA-seq from actively dividing U5 human neural
131 stem cells (U5-hNSCs) cultured *in vitro* (O’Connor et al. 2021). The U5-hNSCs were cultivated
132 from the telencephalon of a human fetus 8 weeks post-conception (Bressan et al. 2017). We
133 previously identified 7 transcriptional states in the U5-hNSCs that were mapped to cell cycle
134 states (i.e., Neural G0, G1, Late G1, S, S/G2, G2/M, and M/Early G1; O’Connor et al. 2021).
135 The U5-hNSC scRNA-seq data were reanalyzed using current quality control and normalization
136 methods which resulted in 2,962 good quality single-cell transcriptome profiles (**Supplemental**
137 **Figure S1A**). The U5-hNSC scRNA-seq profiles, along with the previously established cell cycle
138 labels (O’Connor et al. 2021), represent the most meticulously curated training dataset available
139 for cell cycle classification.

140 We compared the newly implemented ccAFv2 classifier against four distinct classification
141 methods: support vector machine with rejection (SVMrej), random forest (RF), scRNA-seq
142 optimized *k*-nearest neighbor (KNN), and ACTINN (Ma and Pellegrini 2020) which was used to
143 build ccAF (O’Connor et al. 2021). The training dataset for all classifiers consisted of the pre-
144 processed U5-hNSC scRNA-seq subset to the 861 genes upregulated in cell cycle states
145 ($\log_2\text{FC} \geq 0.25$, adjusted p-value ≤ 0.05 ; **Supplemental Table S2**). We applied 10-fold cross-
146 validation (CV) for each classification method (**Supplemental Figure S1A**) and observed that
147 ccAFv2 exhibited significantly improved F1 scores for each cell cycle state compared to other
148 classification methods (p-values $\leq 2.8 \times 10^{-6}$; **Figure 1C**), establishing it as the most accurate
149 cell cycle classifier overall. A benefit of using the F1-score as the performance metric is that it
150 accounts for the imbalance in class label proportions within the training set. We evaluated the
151 impact of balancing label proportions in the training dataset, but this resulted in worse model
152 performance (**Supplemental Figure S1B**). The accuracy of ccAFv2 when applied to U5-hNSCs
153 was 88.4%, and the main difference when compared to ccAF was an improvement in Late G1
154 cell predictions (**Supplemental Figure S1C-D**). The overall error rate for ccAFv2 was 3.3%,
155 which is a considerable improvement from the 18.4% of ccAF (O’Connor et al. 2021). The
156 reimplementation of the ANN for the ccAFv2 classifier has significantly improved its
157 performance across all cell cycle states, providing a robust foundation for further optimization
158 and comprehensive characterization of its capabilities.

159 **Optimizing the number of neurons in hidden layers**

160 A crucial factor in optimizing the parameters of the ccAFv2 ANN was determining the ideal
161 number of neurons in each hidden layer. We conducted a systematic comparison of 18 different
162 combinations for the number of neurons in the two hidden layers (first hidden layer: ranging
163 from 200 to 700 neurons, and second hidden layer: ranging from 100 to 400 neurons) across
164 U5-hNSCs (O'Connor et al. 2021), a low grade glioma stem cell line (LGG275), six glioma stem
165 cell lines (BT322, BT324, BT326, BT333, BT363, and BT368; Couturier et al. 2020), and two
166 glioma tumors (BT363 and BT368; Couturier et al. 2020). The optimal combination was
167 determined by having the highest average F1-score and Adjusted Mutual Information (AMI)
168 score using ccSeurat as the reference (**Figure 1D; Supplemental Table S3**). We chose to
169 employ the ccSeurat classifier (Butler et al. 2018) to predict the reference labels because true
170 cell cycle state labels do not exist for all datasets. The ccSeurat classifier was chosen for three
171 reasons: 1) it is the de facto standard method for cell cycle classification currently, 2) it performs
172 well when applied to many different datasets, and 3) it uses a totally different underlying
173 algorithm to classify cell cycle state than ccAFv2. We found that configuring the ccAFv2 ANN
174 with 600 neurons in the first hidden layer and 200 in the second hidden layer yielded the largest
175 average F1 score and second largest AMI score (**Figure 1D**). This specific parameterization has
176 been assigned for the hidden layers of the ccAFv2 ANN, and all prior and subsequent ccAFv2
177 classifications use this parameterization.

178 **Most important features for classifying ccAFv2 states**

179 After optimizing the training of the ccAFv2 ANN, it is sensible to determine which features are
180 most essential for classifying each of the seven states. We computed feature importance by
181 permuting one of the 861 genes in the U5-hNSCs dataset and asking what impact that had on
182 the likelihoods for each of the seven states. Randomizing the expression of an important feature
183 for classifying a ccAFv2 state would lead to reductions in the states likelihood for cells known to
184 be of this state. Thus, it is crucial that the dataset used for feature importance have cell cycle
185 labels, which is why the U5-hNSCs were used for feature importance analyses (**Figure 1E**). We
186 report the top 15 most important genes for each of seven ccAFv2 states (**Figure 1F-L**).

187 Eleven of the most important genes for classifying the Neural G0 state (**Figure 1F**) were also
188 marker genes of Neural G0 in the U5-hNSCs. The first most important gene for classifying the
189 G1 state (**Figure 1G**) was *HMGN2*, and in prior studies over-expression of *HMGN2* in
190 osteosarcoma cells led to significantly higher number of cells in G0/G1 (Liang et al. 2015). The
191 top two most important genes for the classifying the Late G1 state (**Figure 1H**) include two
192 Immediate-Early Genes (IEGs) *CCN1* and *CCN2* which are known to be induced rapidly after
193 initiation of cell cycle progression by many factors (Tullai et al. 2007). The top four most
194 important genes for classifying the S state (**Figure 1I**) include three genes required for DNA
195 replication during S phase (*CLSPN*, *GINS2*, and *PCNA*) and the cyclin associated with S phase
196 (*CCNE2*). The top four most important genes for classifying the S/G2 state (**Figure 1J**) are all
197 histones, specifically one *H4* histone and multiple *H1* histones isoforms that enable the
198 condensation of nucleosomes into chromatin. The top five most important genes for classifying
199 the G2/M state (**Figure 1K**) include a gene involved in keeping sister chromatids from
200 separating (*PTTG1*), and two genes involved in kinetochore and centromere maintenance and
201 function (*CENPA*, *HMMR*; Maxwell et al. 2005). Additionally, the ninth most important gene for
202 classifying G2/M is *CCNB1* the cyclin that peaks in mitosis, and *MKI67* which is an established
203 marker of cell proliferation (Scholzen and Gerdes 2000). Finally, the top three most important
204 genes for classifying the M/Early G1 state (**Figure 1L**) are a microtubule component protein
205 *TUBA1B*, a microtubule associated protein *STMN1*, and a component of the chromosome
206 passage protein complex (CPC) which is essential for sister chromatid alignment and
207 segregation during mitosis and cytokinesis (Vong et al. 2005). The functions of the key genes

208 for classifying each state align well with the molecular processes of each cell cycle state,
209 supporting the conclusion that the identified classes in U5-hNSCs reflect the underlying biology
210 of the cell cycle.

211 Next, we evaluated the expression of important genes for each ccAFv2 state in an independent
212 dataset, the *in vivo* hNSCs collected from whole fetal brain at 9 weeks post-conception (PCW 9
213 R1; Zeng et al. 2023). This allows us to assess the generalizability of these genes as key
214 markers across novel datasets, providing insight into their broader applicability and robustness.
215 The expression of important genes for all ccAFv2 states were expressed strongly in the state
216 they marked, except for the important genes for the G1 state (**Supplemental Figure S1E**). In
217 our prior study it was difficult to identify markers for G1 phase cells, and so the lack of
218 translation for important genes for the G1 state is not surprising. The successful translation of
219 key genes to an independent dataset supports the hypothesis that ccAFv2 and the marker
220 genes identified in U5-hNSCs are broadly applicable to *in vivo* hNSCs.

221 Comparison with existing cell cycle classifiers

222 An important means to test the performance of ccAFv2 is to compare it to existing state-of-the-
223 art methods for cell cycle state classification. We evaluated the following methods: ccAF
224 (O'Connor et al. 2021), ccSeurat (Hao et al. 2021), tricycle (Zheng et al. 2022),
225 Revelio/SchwabeCC (Schwabe et al. 2020), reCAT (Liu et al. 2017), peco (Hsiao et al. 2020),
226 and cyclone (Scialdone et al. 2015). We also evaluated the incorporation of ccAFv2 marker
227 genes into the ccSeurat classification algorithm. However, it had significantly reduced
228 performance compared to ccSeurat and ccAFv2 (**Supplemental Figure S2**). Each tool predicts
229 a different subset of cell cycle phases, uses a different classification algorithm, was trained on
230 different data, and requires different input genes and data formats (**Supplemental Table S4**).
231 We applied ccAFv2 alongside the other state-of-the-art cell cycle classification methods to *in*
232 *vivo* hNSCs collected from whole human fetal brain at PCW 9 R1 (**Figure 2A-B**; Zeng et al.
233 2023). These cells represent an independent dataset for an unbiased comparison of the cell
234 cycle prediction algorithms. The hNSCs from Zeng et al., 2023 were also chosen for their
235 similarity to the U5-hNSCs and their added real-world relevance, as they were collected *in vivo*.
236 We chose to employ the ccSeurat classifier (Butler et al. 2018) to predict the reference labels for
237 classifier comparison for the reasons described above. The AMI score is impacted by the
238 number of cell cycle states in the reference (i.e., three cell cycle states in ccSeurat), and the
239 number of states predicted by each algorithm (e.g., seven cell cycle states in ccAFv2). We used
240 simulation studies to define the expected range of AMI scores that correspond to specific levels
241 of similarity to the reference given the number of cell cycle states in the reference and the
242 classifier being tested. The highest AMI was observed for tricycle, showing an 80% similarity to
243 the reference (**Figure 2A**). This result aligns with the UMAP colorization, indicating a strong
244 match within classifiers that predicted a comparable number of classes to ccSeurat (**Figure 2B**).
245 reCAT and ccAFv2, predicting six and seven cell cycle states, respectively, achieved the next
246 highest AMI scores, both demonstrating over 70% similarity to the reference (**Figure 2A**).
247 Notably, ccAFv2 identified an S/G2 cluster of cells positioned between the S and G2/M cells
248 classified by ccSeurat and tricycle, which is biologically plausible (**Figure 2B**). Additionally,
249 while Neural G0 cells are intermixed with G1 and Late G1 cells within the proliferating cell
250 population on the left side of the UMAP, the right side reveals a distinct cluster of Neural G0
251 cells (**Figure 2B**). This suggests the presence of a quiescent population in these normal human
252 neural stem cells that is not detectable by the ccSeurat, tricycle, or reCAT classifiers.
253 We also applied ccAFv2 alongside the other cell cycle classification methods to cells derived
254 from a glioblastoma (GBM) patient tumor (BT322; Couturier et al. 2020; **Figure 2C-D**). GBM
255 patient tumors are characterized by both quiescent and proliferating subpopulations (Tejero et
256 al. 2019) making them ideal datasets for evaluating and comparing different cell cycle

257 classification methods. We used the ccSeurat labels as the reference because true cell cycle
258 state labels do not exist for this dataset. Like the *in vivo* PCW 9 R1 hNSCs, the largest AMIs
259 were observed for tricycle, reCAT, and ccAFv2; all of which correspond to just below 90%
260 similarity to the reference (**Figure 2C**). These results demonstrate that ccAFv2 delivers at least
261 equivalent performance when compared to contemporary state-of-the-art cell cycle classifiers,
262 while providing the highest resolution of cell cycle state predictions including a quiescent-like G0
263 state (**Figure 2D**).

264 ***In vivo* cyclin expression and marker genes validate ccAFv2 cell cycle states**

265 We explored the distribution of cell cycle states in 94,297 hNSCs collected from human fetal
266 tissue at 3-12 weeks post-conception (Zeng et al. 2023; **Figure 3**). Application of ccAFv2 to the
267 *in vivo* fetal hNSCs was found to differ by week stage (**Figure 3B**). The amount of Neural G0
268 cells from the *in vitro* U5-hNSCs, derived from fetal brain tissue at 8 weeks post-conception,
269 (**Figure 3A**) matches closely to the *in vivo* hNSCs at eight weeks post-conception (**Figure 3B**).
270 Moreover, the expression patterns of cyclins between the *in vitro* (O'Connor et al, 2021) and
271 5,575 *in vivo* hNSCs from PCW 9 R1 were similar (**Figure 3C-G**). In both hNSC populations,
272 CCNE2 exhibited its peak expression during the S phase, while CCNA2 showed highest
273 expression levels during the S/G2 and G2/M phases, and CCNB1 displayed elevated
274 expression in G2/M phase cells (**Figure 3G**). Notably, the highest expression of the key
275 regulator of cell cycle progression, CCND1, was observed in the Late G1 state (**Figure 3G**).

276 Additionally, we identified ccAFv2 marker genes that corresponded to cell cycle state markers in
277 the PCW 9 R1 hNSCs. Differentially expressed genes for each cell cycle state were identified,
278 and only those overlapping with the ccAFv2 marker gene lists were reported (**Figure 3H**).
279 Genes important to the ccAFv2 classifier were enriched among the translatable marker genes
280 for PCW 9 R1 (**Figure 3H**). The expression patterns of these translatable marker genes were
281 consistent with those observed in O'Connor et al., 2021. The exclusive or semi-exclusive
282 expression of these markers in adjacent cell cycle states strongly supports the presence of high-
283 resolution ccAFv2 clusters in hNSCs *in vivo*. Furthermore, the biological function of the
284 translatable marker genes for each ccAFv2 cell cycle state validates the biological basis of the
285 ccAFv2 clusters, providing further evidence of their relevance.

286 **Defining an appropriate classification likelihood threshold**

287 The improved ccAFv2 classifier calculates likelihoods for each cell cycle state which can be
288 used to determine the most likely state and to assess the quality of the classification. We
289 hypothesized that applying a likelihood threshold to ccAFv2 classifications would ensure
290 reliability and confidence in predicted cell cycle states by setting classifications for cells with less
291 certainty to an "Unknown" state. We explored the range of possible likelihood thresholds on the
292 94,297 hNSCs collected by Zeng et al., 2023.

293 We tested ccAFv2 likelihood thresholds ranging from 0.0 to 0.9 in increments of 0.1
294 (**Supplemental Figure S3**). The calculated cell cycle state likelihood was required to be greater
295 than or equal to the threshold, otherwise an "Unknown" state was returned (**Figure 1B**). Each
296 likelihood threshold was assessed using the percentage of cells predicted and an AMI score
297 with ccSeurat cell cycle states as a reference. As the likelihood threshold increases the number
298 of cells predicted decreases and the AMI scores increase (**Figure 3I; Supplemental Figure S3**;
299 **Supplemental Table S5**). In other words, the removal of less certain classifications improves
300 the accuracy of the overall classifications (**Figure 3I**). Next, we further demonstrated that the
301 increase in AMI resulted from the specific removal of cells which had low classification
302 likelihoods, by comparing it to the random removal of an equivalent number of cells
303 (representative analysis for 9 weeks post-conception is shown in **Figure 3J**). The randomly
304 removed cells do not increase the AMI (**Figure 3J**), only the selected removal of cells with low
305 likelihoods were able to increase the AMI. We found that the median AMI scores calculated with

306 likelihood thresholds of 0.4 to 0.9 were significantly higher than the median AMI scores of the
307 randomly removed cells (**Figure 3J; Supplemental Table S5**), which indicates that the
308 likelihood cutoffs of greater than or equal to 0.4 improve classification accuracy. We selected
309 the likelihood threshold of greater than or equal to 0.5 because it signifies a minimum of 50%
310 certainty in the classified cell cycle state. Additionally, greater than 90% of *in vivo* hNSCs could
311 be assigned a cell cycle state with a likelihood threshold of 0.5 (**Figure 3I**). Thus, the threshold
312 of 0.5 was set as the default for ccAFv2 and used in subsequent analyses, except where noted.
313 We also provide users with the flexibility to adjust the likelihood threshold parameter in ccAFv2,
314 allowing them to adapt the classifier's operation to suit the unique characteristics of their
315 dataset.

316 **Effect of missing gene expression values on ccAFv2**

317 A known limitation of scRNA-seq is that dropouts are common. A dropout occurs when lowly to
318 moderately expressed transcripts are detected in one cell but are not detected in another cell of
319 the same cell type (Qiu 2020). Factors affecting dropouts include the number of sequencing
320 reads from each cell and the complexity of the cell's transcriptome. The ccAFv2 classifier uses
321 the expression of 861 genes to predict cell cycle states. We hypothesized that dropouts could
322 be simulated by randomly setting the expression of a defined percentage of genes to zero and
323 that this would provide a reasonable approximation of the influence of missing genes on the
324 accuracy of ccAFv2's cell cycle state classifications. We evaluated the consequences of these
325 simulated gene dropouts on the classifier error rate, AMI, and the number of cells predicted
326 (**Supplemental Figure S4-5; Supplemental Table S6-7**). As described earlier, the median
327 error rate of applying ccAFv2 to U5-hNSCs was 3.3% with 99% of the input genes (99% is used
328 to allow for cross-validation). Missing information for 20% of the ccAFv2 input genes yielded a
329 smaller median error rate (12.2%) than the original ccAF error rate with all the input genes
330 (18.4%), underscoring the improved performance of the new model. Introducing missing
331 information for 40% of ccAFv2 input genes led to a 29% median error rate, and 96% of cells
332 were predicted (**Supplemental Table S6**). The error rate was the most affected by the
333 introduction of missing information (**Supplemental Figure S4A**) and the median percentage of
334 cells predicted remained above 80% even when 70% of the input gene list was set to missing
335 (**Supplemental Table S6**). When breaking down the error rate by cell cycle state, we observed
336 that S and M/Early G1 had the highest error rates as missing information increased
337 (**Supplemental Figure S4B; Supplemental Table S7**). However, the number of cells predicted
338 remained relatively consistent across all states despite the increasing in missing data
339 (**Supplemental Figure S4C; Supplemental Table S7**). The increase in error rate without a
340 concomitant decrease in the number of cells predicted suggests that raising ccAFv2's likelihood
341 threshold (>0.5) might be required to ensure the quality of predictions for datasets with greater
342 than 20% missing ccAFv2 input genes. Indeed, the error rate for introducing 20% missing
343 information decreased from 12.2% median error rate at 0.5 likelihood threshold to 9.9% with a
344 0.7 likelihood threshold (**Supplemental Figure S4D; Supplemental Table S6**) and 6.2% with a
345 0.9 likelihood threshold (**Supplemental Figure S4E; Supplemental Table S6**). Thus,
346 introducing 20% missing information led to four times the error rate, and the increased error rate
347 can be mitigated in part by increasing the likelihood threshold. Increasing the likelihood
348 threshold decreases the error rate by removing classifications for cells where the missing
349 information has degraded the confidence in the prediction. By removing predictions with less
350 confidence, the error rate decreases, but the overall number of cells classified with cell cycle
351 states decreases. Testing the impact of increasing the likelihood threshold on the number of
352 predicted cells can be quite insightful for choosing an appropriate likelihood threshold
353 (**Supplemental Figure S4F**). Careful consideration of the balance between minimizing errors
354 and retaining enough cells for downstream studies is essential.

355 **Neural G0 state is enriched in mesenchymal G0 cells**

356 We developed an experimental method to isolate fixed G0 cells using fluorescence-activated
357 cell sorting (FACS) with established markers (Gookin et al. 2017). This approach ensures the
358 selected cells are diploid, non-replicating, and have unphosphorylated RB (pRB⁻), all while
359 preserving RNA integrity. The experimental approach was applied to identify G0 cells from
360 human skeletal muscle satellite cells (hSkMSCs), which are derived from the mesodermal germ
361 layer (**Figure 4A**). Subsequent RNA-seq of the sorted G0 cells captures the characteristic
362 expression pattern of the G0 state for that cell type, enabling direct comparison with the
363 expression profiles of single cells from scRNA-seq data of unsorted cells. We characterized the
364 bulk RNA-seq signatures of 400,000 G0 hSkMSCs for two biological replicates. These G0
365 signatures were mapped onto 6,921 asynchronous, unstained, and unsorted hSkMSCs
366 collected using scRNA-seq. Next, the ccAFv2 cell cycle classifier was applied to the hSkMSC
367 scRNA-seq data. Cells in S, S/G2, G2/M and M/Early G1 states formed distinct clusters ordered
368 in the canonical cell cycle patterning, highlighting the generalizability of these ccAFv2 states to
369 *in vitro* hSkMSCs (**Figure 4B**). The Late G1 cells did form a cluster that was positioned in front
370 of the S phase cells, however, additional Late G1 cells can be seen dispersed within the left half
371 of the cells. This dispersion may be partly attributed to the higher *CCND1* expression in
372 hSkMSCs G0 cells compared to hNSC Neural G0 cells (**Supplemental Figure S6**). The cells
373 inside the dashed region contain almost exclusively Neural G0, G1, and Late G1 cells and do
374 not coalesce into defined clusters (**Figure 4B-C**), suggesting that the ccAFv2 classifier was
375 struggling to accurately discriminate between Neural G0, G1, and Late G1. Correlation of the
376 experimentally determined G0 signature to the cells from the scRNA-seq revealed significant
377 enrichment within the ccAFv2-labeled Neural G0 and G1 cells (**Figure 4D-F**). The majority of
378 G0 cells were classified as Neural G0, while misclassified G0 cells were predominantly labeled
379 G1 or Late G1, and very few were classified as the cycling states (S, S/G2, G2/M, M/Early G1;
380 **Figure 4G**). These findings confirm that ccAFv2 is having difficulty discriminating between the
381 Neural G0, G1, and Late G1 states in hSkMSCs derived from the mesoderm dermal layer.
382 However, the misclassifications are systematic rather than random, and a straightforward
383 solution of merging Neural G0, G1, and Late G1 classifications effectively resolves the
384 misclassifications. To accommodate this, we implemented a switch in the ccAFv2 classifier,
385 enabling users to choose whether to combine Neural G0, G1, and Late G1 or to keep them
386 separate. This feature provides a more cautious and flexible approach for classifying cell types
387 beyond neuroepithelial cells (**Figure 4H**). This provides users with the flexibility to use ccAFv2
388 higher resolution cell cycle classification for non-neuroepithelial cell types.
389

390 Additionally, it should be noted that while the Neural G0 state does not accurately capture the
391 G0 state of hSkMSCs the experimental data demonstrates that it is possible to identify a
392 subpopulation of G0 cells. The marker genes discovered for the hSkMSC G0 cells were
393 significantly overlapping with the Neural G0 marker genes (n = 11; p-value = 4.4 x 10⁻⁴; *FTL*,
394 *IFITM3*, *TIMP4*, *SAT1*, *C1orf21*, *CLU*, *VGLL4*, *SPRY1*, *COL9A3*, *NOVA1*, *NUDTA*; **Figure 4I**;
395 **Supplemental Table S8**. Which strongly suggests that it may be possible to train a classifier
396 with a more generalizable G0 state in future studies using our experimental approach to
397 characterize new training datasets.

398 **ccAFv2 cell cycle states are generalizable across germ layers**

399 Another key consideration when using ccAFv2 is its ability to accurately predict cell cycle states
400 (S, S/G2, G2/M, and M/Early G1) in cell types beyond neuroepithelial cells. In Zeng et al., 2023,
401 they profiled single cells from human fetal tissues, representing all three germ layers
402 (endoderm, mesoderm, and ectoderm; **Figure 5A**). We applied ccAFv2 to 245,906 cells from
403 the atlas first including Neural G0, G1 and Late G1 predictions (**Figure 5B**), and then by
404 collapsing these three predictions into a G0/G1 class (**Figure 5C**). These representations of the

405 Zeng et al., 2023 dataset aggregate across time (PCW 3–12) and tissue collection strategies
406 (whole embryo, whole head, and brain). The proportions of cell cycle states across time for each
407 cell type show strong concordance (**Supplemental Figure S7**), highlighting the consistency of
408 ccAFv2 predictions across biologically similar independent scRNA-seq datasets. Next, we
409 analyzed cyclin expression patterns across the 15 distinct cell types and calculated the average
410 expression for each germ layer. Each germ layer exhibited the expected cyclin expression
411 pattern, with *CCND1* driving entry into the cell cycle (**Figure 5D**), *CCNE2* peaking during S
412 phase (**Figure 5E**), and *CCNA2* (**Figure 5F**) and *CCNB1* coordinating progression and
413 regulation of cell division (**Figure 5G**). These results provide strong evidence that the ccAFv2
414 predicted S, S/G2, G2/M, and M/Early G1 states are accurate across cell types derived from all
415 germ layers.

416
417 We also observed that ccAFv2 proportions align closely with expected developmental patterns.
418 For instance, during PCW 9–12, the fetal brain undergoes rapid development, marked by
419 significant cell division as major structures like the cerebrum, cerebellum, and brainstem
420 become more defined (Belmonte-Mateos and Pujades 2021; Martínez-Cerdeño et al. 2006). We
421 found that intermediate progenitor cells (IPCs) detected at PCW 9 and 12 in whole brains were
422 predominantly in the S, S/G2, and G2/M phases, consistent with the active proliferation
423 necessary for forming these brain structures (**Figure 5**; **Supplemental Figure S7**).
424 Furthermore, these PCW 9–12 IPCs exhibited high expression of EOMES, a critical factor that
425 drives the expansion of the IPC pool (Arnold et al. 2008). The much-reduced proportions of
426 Neural G0, G1, and Late G1 in IPCs validates that ccAFv2 predictions are consistent with
427 known biology. We also observed reduced Neural G0, G1, and Late G1 proportions in the non-
428 neuroepithelial proliferating mesoderm (Prolif. meso.; **Figure 5**) defined by high expression of
429 the proliferation marker *MKI67* ($\log_2(\text{FC}) \geq 1.24$) and mesoderm marker *CDH11* ($\log_2(\text{FC})$
430 >2.27) (Hoffmann and Balling 1995). The reduced number of non-cycling cells in IPCs and
431 proliferating mesoderm cell types is well documented and demonstrates that while Neural G0,
432 G1 and Late G1 misclassifications may occur that the relative proportions of non-cycling cells to
433 cycling cells is accurately determined by ccAFv2.

434 **Capturing the effect of growth factors on cellular proliferation**

435 Growth factors are used to increase cellular proliferation *in vitro*, and we characterized the
436 transcriptomes of LGG cells (grade 2 astrocytoma and grade 3 oligodendrogloma) with and
437 without the application of growth factors (**Supplemental Figure S8**). For this analysis we tested
438 the impact of adjusting the ccAFv2 likelihood threshold across a range of values 0 to 0.9
439 (**Supplemental Figure S9A-B**). Increasing the likelihood threshold values from 0.4 to 0.9 led to
440 an increased proportion of “Unknown” classifications in the samples without growth factors,
441 which is consistent with the known effect of growth factors to stimulate proliferation and the cell
442 cycle. The increased proportion of “Unknown” cells may correspond to new growth factor
443 starvation state(s) not included in ccAFv2 classification states. Additionally, the S, S/G2, and
444 G2/M cell cycle states were disproportionately removed as the likelihood threshold increased
445 (**Supplemental Figure S9A-B**). We then set the likelihood threshold to 0.9 and observed that
446 the cells grown with growth factors form clusters of cell cycle state labels, outlining the expected
447 progression of cell cycle phases ($\text{G1} \rightarrow \text{S} \rightarrow \text{S/G2} \rightarrow \text{G2/M} \rightarrow \text{M/Early G1}$; **Supplemental**
448 **Figure S8A, C, D, & F**). Conversely, cells grown without growth factors exhibit a more
449 dispersed distribution of cell cycle state labels (**Supplemental Figure S8B, C, E, & F**). The
450 ability to change the likelihood threshold of ccAFv2 allows us to observe the biological impact of
451 adding growth factors to LGG cells and demonstrates what to expect when the cell cycle is not
452 the main transcriptional signal in cells.

453 454 **Removing cell cycle expression signatures**

455 The cell cycle generates a strong transcriptional signature that can obscure other less robust
456 transcriptional signatures of interest. Previous studies have shown that statistical methods can
457 effectively remove cell cycle transcriptional signatures, and that the residual transcriptional
458 variance can be used to study less robust transcriptional signatures of interest (Luecken and
459 Theis 2019). We showcase successful removal of the cell cycle transcriptional signatures for the
460 U5-hNSCs, and LGG cells. First, each cell cycle state's average expression of marker genes is
461 computed for every single cell or nuclei. Then, these average cell cycle expression patterns are
462 regressed out of the dataset during normalization. The ccSeurat regression method uses only
463 the S and G2/M cell cycle states, so we first tested regression with the S and G2/M cell cycle
464 states from ccAFv2. We found that the ccAFv2 marker gene derived average cell cycle
465 expression patterns could mitigate cell cycle transcriptional signatures as effectively as ccSeurat
466 (empirical p-value > 0.05, **Supplemental Table S9; Supplemental Figure S10**). Additionally,
467 we found that incorporating Late G1, S, S/G2, G2M, and M/Early G1 was also quite effective
468 and led to a more robust homogenization of the cell cycle states based on PCA plots
469 (**Supplemental Table S9; Supplemental Figure S10**). This approach enables researchers to
470 dissect complex gene expression patterns and uncover novel insights into cellular processes
471 beyond the cell cycle.
472

473 **Classifying neuroepithelial-derived cells in humans and mice**

474 It was crucial for ccAFv2 to be highly user-friendly, ensuring researchers can easily apply it
475 across a wide range of datasets. The model was designed to accept inputs for tissue source,
476 data type, and gene identifier, eliminating the need for manual data conversion (**Figure 6A**). In a
477 previous study (O'Connor et al. 2021) we applied the ccAF classifier to cells from the developing
478 human telencephalon (Nowakowski et al. 2017). We applied ccAFv2 to these same cells and
479 compared the ccAF and ccAFv2 predicted cell cycle proportions. We observed that the Neural
480 G0 state was less frequent in all cell types for ccAFv2 relative to ccAF (**Figure 6B**;
481 **Supplemental Table S10**). The Neural G0 state was distinctly less frequent in the neuronal cell
482 types. For ccAF Neural G0 made up most of the cell cycle states for EN-PFC and EN-V1, but in
483 ccAFv2 these two cell types classified primarily as G1 (**Figure 6B**). The glial cell types had the
484 largest Neural G0 subpopulations (**Figure 6B**).

485 We also applied ccAFv2 to cells from the ventricular-subventricular zone (V-SVZ) of the adult
486 mouse brain (Cebrian-Silla et al. 2021), a location known to contain neural stem and precursor
487 cells in the adult brain (Lim and Alvarez-Buylla 2016). The adult mouse V-SVZ validates
488 observations from the developing human telencephalon (**Figure 6C**). In the V-SVZ the glial cells
489 tended to have larger Neural G0 subpopulations, neuronal cell types tended to have less Neural
490 G0, and microglial had the smallest amount Neural G0 (**Figure 6C**). The results are similar
491 given the differences between species, developmental state, and anatomical origins. These
492 findings illustrate that the ccAFv2 classifier can be applied to cells originating from both humans
493 and mice.

494 **Classifying quiescent-like neural stem cells**

495 Previously we validated the Neural G0 state using two independent *in vivo* scRNA-seq profiling
496 studies of NSCs from adult neurogenesis in the subventricular zone (SVZ) that used
497 fluorescence activated cell sorting (FACS) to sort out quiescent and activated NSCs (Llorens-
498 Bobadilla et al. 2015; Dulken et al. 2017). We applied ccAFv2 to these same cells and
499 compared the ccAF and ccAFv2 predicted Neural G0 subpopulations. Overall, the qNSCs are
500 enriched with quiescent-like Neural G0 cells, and the aNSCs are at some stage of the cell cycle
501 (**Figure 6D**). The proportion of cells classified as Neural G0 decreased for ccAFv2 in the
502 quiescent NSCs (qNSCs) and was replaced by more G1, S/G2, and a small amount of G2/M
503 (**Figure 6D; Supplemental Figure S11**). For Llorens-Bobadilla et al., 2015 the active NSCs 1
504 (aNSC1) were more highly enriched with S phase cells, and aNSC2 were enriched with S/G2

505 and G2/M. A similar trend was observed for the Dulken et al., 2017 dataset. Additionally, we
506 used the G0 arrest signature from Wiecek et al., 2023 to validate the Neural G0 state in ccAFv2
507 (Wiecek et al. 2023). We found significant enrichment of the G0 arrest signature (i.e., QuieScore
508 G0) within the U5-hNSC Neural G0 and G1 states (**Supplemental Figure S12**). These results
509 continue to validate our assertion that Neural G0 represents a quiescent-like cell state, and that
510 ccAFv2 can accurately classify this quiescent-like G0 state.

511 **Accurate classification of cells and nuclei**

512 Tissues in single-cell studies can be processed into cells for scRNA-seq or nuclei for snRNA-
513 seq. Both methods are commonly used and have advantages and limitations (Slyper et al.
514 2020). Thus, it is important to demonstrate whether ccAFv2, which is trained on cells, can
515 accurately classify cell cycle states for single nuclei. We employed the Zhang et al., 2021
516 dataset which characterized developing human spinal cord tissue from five developmental time
517 points using both scRNA-seq and snRNA-seq from the same experimental conditions (Zhang et
518 al. 2021). The proportions of cells in each cell cycle state are similar between scRNA-seq and
519 snRNA-seq from the same condition (**Figure 6E**), illustrating the versatility of the ccAFv2
520 classifier in effectively analyzing both scRNA-seq and snRNA-seq profiles.

521 **Mapping proliferation onto tissue through spatial transcriptomics**

522 scRNA-seq and snRNA-seq provide valuable information about the transcriptional states of cells
523 and nuclei, but without contextual information, relating these states with previously described
524 biology can be challenging. Spatial transcriptomics captures the transcriptional activity of a
525 single-cell or a region containing a small number of cells at a position within an intact tissue,
526 offering structural information that can be referenced to anatomical atlases and established
527 histology. We applied ccAFv2 to the highest resolution sequencing-based spatial transcriptomic
528 dataset currently available, derived from a slice of a mouse E15.5 embryo binned to 8 μ m,
529 achieving near cellular resolution (**Figure 7A-B**). Despite rapid development at late prenatal
530 stages, Neural G0 classified spots were highly prominent, particularly across the midbrain
531 region. The proliferation marker *Mki67* was minimally expressed in these G0 enriched areas
532 and, in the brain, highlighted the described stem cell niches of the lower, medial, and caudal
533 ganglionic eminences (LGE, MGE, CGE), along with the stem cell migratory paths from these
534 regions along the sub-ventricular zone (SVZ; **Figure 7C**; Kriegstein and Noctor 2004).

535 From E12.5 to E17.5 the mouse cortex develops in a well-defined layers (**Figure 7D**). Applying
536 ccAFv2 to the cortex captured the layered patterning of cell cycle states that fit with our current
537 model of cortical development (**Figure 7E-Q**). Histology showed an outer layer of dermis (skin)
538 marked by *Krt5* expression (**Supplemental Figure S13A-B**), covering the developing skull
539 identified by *Col1a1* expression (**Supplemental Figure S13C**), followed by densely packed
540 cells of the brain (**Figure 7D**). Precursors of excitatory neurons migrate along and divide in the
541 SVZ with asymmetric division specifically occurring within the ventricular zone (VZ). A much
542 smaller population of inhibitory neuronal precursors migrate and divide along the intermediate
543 zone (IZ) and medial zone (MZ) before invading the cortical plate (CP) and differentiating into
544 their neuronal sub-type. By E15.5, the CP is already well populated with post-mitotic neurons
545 that previously migrated from the SVZ from E12.5-E14.5 and will form layers IV-VI of the adult
546 cortex. Upon reaching the border of the CP and MZ, neural stem cells receive maturation
547 factors from glial cell types including *Re1n* (**Figure 7F**), with canonically post-mitotic neurons
548 marked by *Tbr1* (**Figure 7G**; Englund et al. 2005). This post-mitotic region classifies as Neural
549 G0 by ccAFv2 (**Figure 7E & K**). Similarly, *Mki67* is sparse within the CP, but highly active in the
550 SVZ and VZ (**Figure 7J**). Intermediate progenitor cells (IPCs) in the SVZ, marked by *Eomes*,
551 undergo symmetric division before radiating outward (**Figure 7G**). Within the VZ radial glia
552 migrate further inward before asymmetric division, with newly divided IPCs radiating back up to
553 the SVZ and outward to populate the cortex. These events are captured by ccAFv2 as a single

554 enriched S and S/G2 band marking the SVZ (**Figure 7N & O**). We also observed two bands
555 containing G2/M classified spots, corresponding to regions of symmetric division in the SVZ and
556 asymmetric division in the VZ (**Figure 7P**). Cells committed to differentiation, immediately
557 migrate outward along and the few spots that classify as M/Early G1 were almost entirely in the
558 IZ and above (**Figure 7Q**). These results demonstrate that ccAFv2 can be effectively applied to
559 spatial transcriptomics of the developing cortex, accurately recapitulating known biological
560 insights into the spatial organization of proliferative activity.

561 Discussion

562 We designed ccAFv2 to use transcriptomic data to accurately classify cell cycle states and a
563 quiescent-like G0 state for single cells or nuclei. The performance of the updated classifier was
564 superior to its predecessor and demonstrated comparable or better performance than other
565 state-of-the-art cell cycle classifiers. The ccAFv2 classifies cells into a broader range of cell
566 cycle states than the contemporary state-of-the-art cell cycle classifiers (Hao et al. 2021; Zheng
567 et al. 2022; Schwabe et al. 2020; Liu et al. 2017; Hsiao et al. 2020; Scialdone et al. 2015) and
568 includes a quiescent-like G0 state. Moreover, ccAFv2 features a tunable parameter to filter out
569 less certain classifications. We showcased the versatility of ccAFv2 by successfully applying it
570 to classify cells, nuclei, and spatial transcriptomics data in humans and mice, using various
571 normalization methods and gene identifiers. The classifier can be used either as an R package
572 integrated with Seurat (https://github.com/plaisier-lab/ccafv2_R) or a PyPI package integrated
573 with SCANPY (<https://pypi.org/project/ccAF/>). We proved that ccAFv2 has enhanced accuracy,
574 flexibility, and adaptability across various experimental conditions, establishing ccAFv2 as a
575 powerful tool for exploring cell cycle dynamics in diverse biological contexts.

576 A major limitation of developing cell cycle classifiers is a lack of scRNA-seq datasets with
577 ground truth labels for each cell cycle state, including G0. Previous studies have used the DNA
578 stain Hoechst (Buettnner et al. 2015) or FUCCI (Leng et al. 2015) to sort cells into G1, S, and
579 G2M subpopulations. However, the limitations of these studies render them unsuitable for
580 constructing a classifier. These limitations include the use of embryonic stem cells which lack
581 distinct G1 or G0 phases as the model system (Ballabeni et al. 2011; White and Dalton 2005),
582 not of fixing the cells, and reliance on non-transcriptional markers. We demonstrated the ability
583 to isolate cells from specific cell cycle states, characterize them transcriptionally, and map their
584 signatures onto scRNA-seq data to identify the spatial distribution of those states. In future
585 studies, we aim to leverage this approach to create training datasets for developing more robust
586 and generalizable cell cycle classifiers.

587 Batch effects pose a significant challenge for single cell, nuclei, and spatial RNA-seq studies;
588 and we have addressed their impact on ccAFv2 in three ways. First, we strongly recommend
589 users apply ccAFv2 to each dataset separately prior to any integration or combining of datasets.
590 The ccAFv2 predictions in the metadata integrate very easily and preempt any issues caused by
591 batch effects between datasets. Second, we advise using SCTtransform normalization to
592 standardize datasets to the Pearson residual scale, which mitigates differences in magnitude
593 and variance (Hafemeister and Satija 2019). The ccAFv2 R package is specifically
594 parameterized to reapply SCTtransform, normalizing the expression of all genes, not just the
595 most variable ones, to maximize overlap with ccAFv2 marker genes, ensuring optimal
596 classification accuracy for each cell. Finally, ccAFv2 model incorporates expression of multiple
597 marker genes to predict a cell cycle state. This approach minimizes the influence of batch
598 effects, as a significant misclassification would require the simultaneous, directional alteration of
599 multiple marker genes, a highly improbable scenario for random batch effects. By adhering to
600 these recommendations and leveraging the design of ccAFv2, we provide effective strategies to
601 mitigate the impact of batch effects when using ccAFv2.

602 The ccAFv2 classifier will be most helpful in biological contexts where the cell cycle is active.
603 We utilized atlases of developing human and mouse embryos and fetuses because proliferation
604 is essential in developing organisms (Soufi and Dalton 2016; Pauklin and Vallier 2013).
605 Evidence is building to show that cell fate decisions are tightly coupled to cell cycle events and
606 machinery (Pauklin and Vallier 2013). In healthy adult organisms, proliferation plays critical roles
607 in several processes: maintenance of stem cell populations (Harada et al. 2021), clonal
608 expansion of both innate and adaptive immune cells (Adams et al. 2020), and germ cell meiosis,

609 encompassing oogenesis (Bukovsky et al. 2005) and spermatogenesis (Guo et al. 2018). Our
610 recommendation is to aggregate Neural G0, G1, and Late G1 into G0/G1 when applying ccAFv2
611 to cell types that are not neuroepithelial. Defects in cell cycle machinery or regulation can lead
612 to runaway proliferation characteristic of cancer (Hanahan and Weinberg 2000), or the lack of
613 proliferation of crucial cell types can lead to neurodegenerative disorders (Joseph et al. 2020).
614 Cell cycle classification would benefit any *in vitro*, *in vivo*, or *ex vivo* studies of proliferating cells.
615 On the other hand, we provide methods to regress the cell cycle expression patterns out of
616 single cell or nuclei data to uncover underlying biological signals. Overall, incorporating cell
617 cycle states into single-cell and nuclei studies enhances our ability to dissect complex biological
618 systems, unravel cellular heterogeneity, and decipher the molecular mechanisms by which
619 proliferation affects cellular processes.

620 The studies reported here demonstrate that the quiescent-like G0 state (Neural G0) in ccAFv2 is
621 detectable across all three germ layers in developing human fetal cells and provided a list of
622 putative marker genes that are common across cell types. This corroborates our previous
623 findings that Neural G0 was an active transcriptional signature executed by a subpopulation of
624 U5-hNSCs (O'Connor et al. 2021). However, support for novel G0 states was observed in the
625 growth factor deprived LGG cells, where an increased proportion of "Unknown" cells was
626 detected, hinting at novel quiescent-like state(s) missing from ccAFv2. Other studies have
627 identified multiple G0 states in a single cell type that are invoked in response to different stimuli
628 (e.g., spontaneous loss of mitogenic factors, serum starvation, drug treatment, etc.) (Stallaert et
629 al. 2022). Thus, we find it very likely that additional G0 states with distinct transcriptional
630 signatures will be identified. The ccAFv2 ANN and its associated training software are fully
631 equipped to integrate these additional G0 states. Future studies that extend the cell cycle
632 classifier to include novel G0 states holds immense potential for advancing our understanding of
633 quiescence in biological systems. By leveraging advanced computational methods, high-
634 throughput technologies, and interdisciplinary approaches, researchers can unravel the
635 complexities of cellular dormancy and pave the way for innovative strategies to manipulate
636 quiescent cell behavior to improve health and combat disease.

637

638 **Methods**

639 **Culture of LGG glioma neurospheres**

640 For the “no growth factors” condition, cells from LGG glioma neurospheres (LGG275, BT237)
641 were dissociated, seeded into two poly-D-lysine and laminin coated T25 cm² flasks at a density
642 of 40,000 cells/cm², and cultured for 4 days using medium without growth factors. For the “with
643 growth factors” condition, cells were cultured for 4 days as neurospheres with EGF and FGF2 at
644 10µg/L and heparin at 2mg/L in PolyHEMA coated flasks, and medium was replaced 1 day
645 before single cell sequencing.

646

647 **scRNA-seq library preparation and sequencing of LGG glioma neurospheres**

648 Cells were dissociated and single cell suspensions loaded onto the Chromium controller (10x
649 Genomics, Pleasanton, CA) to generate single-cell Gel Beads-in-Emulsion (GEMs). The single-
650 cell RNA-seq libraries were prepared using the Chromium Next GEM Single Cell 3' Reagent
651 Kits V3.1 (Dual Index, P/N 1000268, 10x Genomics). Briefly, reverse transcription was
652 performed at 53°C for 45 min followed by incubation at 85°C for 5 min. GEMs were then broken
653 and the single-stranded cDNAs were cleaned up with DynaBeads MyOne Silane Beads
654 (Thermo Fisher Scientific; P/N 37002D). The cDNAs were PCR amplified, cleaned up with
655 SPRIselect beads (SPRI P/N B23318), fragmented, end-repaired, A-tailed, and size-selected
656 with SPRIselect beads. Indexed adapters were ligated and cleaned up with SPRIselect beads.
657 The resulting DNA fragments were PCR amplified and size selected with SPRIselect beads.
658 The size distribution of the resulting libraries was monitored using a Fragment Analyzer (Agilent
659 Technologies, Santa Clara, CA, USA) and the libraries were quantified using the KAPA Library
660 quantification kit (Roche, Basel, Switzerland). The libraries were denatured with NaOH,
661 neutralized with Tris-HCl, and diluted to 150 pM. Clustering and sequencing were performed on
662 a NovaSeq 6000 (Illumina, San Diego, CA, USA) using the paired-end 28-90 nt protocol on one
663 lane of an SP flow cell and on one lane of an S4 flow cell. Sequencing data can be accessed
664 from NCBI SRA. Both library preparation and sequencing were performed at the Montpellier
665 GenomiX facility (MGX) in Montpellier, France.

666

667 **Data analysis**

668 Image analyses and base calling were performed using the NovaSeq Control Software and the
669 Real-Time Analysis component (Illumina). Demultiplexing was performed using the 10x
670 Genomics software Cellranger mkfastq (v7.1.0), a wrapper of Illumina's bcl2fastq (v2.20). The
671 quality of the raw data was assessed using FastQC (v0.11.9) from the Babraham Institute and
672 the Illumina software SAV (Sequencing Analysis Viewer). FastqScreen (v0.15.1) was used to
673 identify potential contamination. Alignment, gene expression quantification and statistical
674 analysis were performed using Cell Ranger count with the human's transcriptome (GRCh38). To
675 discard ambient RNA falsely identified as cells, Cell Ranger count was run a second time with
676 the option --force-cells to force the number of cells to detect. Cell Ranger aggr was used to
677 combine each sample results into one single analysis. Cell Ranger output files can be accessed
678 from NCBI GEO at GSE263796.

679

680 **scRNA-seq, snRNA-seq, and ST-seq neuroepithelial datasets**

681 In total 42 scRNA-seq, 11 snRNA-seq, and 8 ST-seq datasets were processed and employed in
682 the studies used to characterize the ccAFv2 classifier. Detailed descriptions of the source,
683 quality control, processing, normalization of each dataset can be found in the **Supplemental**
684 **Material**. Analyses were conducted in R (R Core Team 2025) using the specified packages.
685

686 **Implementation of the ccAFv2 ANN model**

687 The core algorithm of the ccAFv2 is a fully connected artificial neural network (ANN)
688 implemented using the Keras API (v2.12.0) that employs TensorFlow (v2.12.0) to construct
689 ANNs. A fully connected ANN model was developed to classify the cell cycle state of single cells
690 (**Figure 1A**). The input for a single cell is expression for the 861 most highly variable genes
691 ($\log_2(\text{FC}) > 0.25$; p value adj < 0.05) from O'Connor et al., 2021. A dense input layer takes in the
692 expression of the 861 and is fully connected to the first hidden layer comprised of 600 neurons.
693 The first hidden layer is fully connected to the second hidden layer comprised of 200 neurons
694 which then connects to the output layer of seven neurons (one for each cell cycle class: Neural
695 G0, G1, Late G1, S, S/G2, G2/M, and M/Early G1). A SoftMax regression function in the output
696 layer is used to compute the likelihood for each class. Overfitting in the ANN is prevented
697 through the incorporation of two dropout layers using a dropout rate of 50%. The first dropout
698 layer is positioned between the first and second hidden layers and the second dropout layer is
699 between the second hidden layer and the output layer (**Figure 1A**, Xie et al. 2019). Neuron
700 activation functions were modeled using the Rectified Linear Unit (ReLU) function. The loss
701 function for the ccAFv2 ANN was categorial cross-entropy and Stochastic Gradient Descent
702 (SGD) was used to optimize the learning. The predicted class for a single cell is identified as the
703 highest likelihood exceeding the specified threshold (**Figure 1B**). By default, the threshold is set
704 at 0.5 and can be adjusted within the range of 0 to 1. If a cell's likelihood falls below the
705 threshold it is classified as "Unknown."
706

707 **Training the ccAFv2 ANN classifier**

708 The ccAFv2 ANN model was trained on the 2,692 cells and 861 genes from the U5-hNSCs
709 dataset (O'Connor et al. 2021) using the labels from O'Connor et al., 2021. The training process
710 encompassed ten epochs repeated five times consecutively. In each epoch, the training data
711 was randomly partitioned into 80% for training and 20% for testing, with the testing subset held
712 out to assess training accuracy.
713

714 **Comparing ccAFv2 to other classification methods**

715 Classifiers were trained using the scRNA-seq gene expression of 2,962 cells with 861 genes
716 and cell cycle labels from the U5-hNSCs. The ccAFv2 classifier was tested against: (i) support
717 vector machine with reject option (SVMrej; classification cutoff ≥ 0.7), a general-purpose
718 classifier from the Scikit-learn library; (ii) random forest (RF), another general-purpose classifier
719 from the Scikit-learn library; (iii) k -nearest neighbor (KNN) from the SCANPY ingest method
720 (Wolf et al., 2018); and (iv) neural network (NN) ACTINN (Ma & Pellegrini, 2020). Classifier
721 performance was determined using F1 scores computed for each cell cycle state. Ten-fold
722 cross-validation with an 80% training and 20% testing split was used to determine the variance
723 of F1 scores for each cell cycle state from each classifier. A Student's *t*-test was used determine
724 if the mean of the F1 scores were significantly lower than ccAFv2.

725

726 **Optimizing the number of neurons in hidden layers**

727 The configuration of neurons in the two hidden layers is designed to reduce the number of
728 neurons at each layer from the 861 input genes down to the 7 cell cycle states. In total, 18
729 ccAFv2 models were trained using the U5-hNSCs dataset to determine the optimal number of
730 neurons for these hidden layers. This involved testing at increments of 100 the number of
731 neurons in the first hidden layer within the range of 200 to 700 neurons and in the second
732 hidden layer within the range of 100 to 400 neurons. For comparisons the F1 scores were
733 computed for each cell cycle state. Each model was also tested on pre-processed scRNA-seq
734 data of glioma stem cells (BT322, BT324, BT326, BT333, BT363, BT368) and tumor cells
735 (BT363, BT368) from Couturier et al. 2022 (**Supplemental Table S11**), along with Grade 2
736 Astrocytoma (LGG275; AUGUSTUS et al. 2021) (**Supplemental Table S12**). For these
737 datasets, Adjusted Mutual Information (AMI) scores, with the reference labels derived from
738 ccSeurat calls, and the number of cells predicted were calculated using the AMI function from
739 the aricode package in R. Barcodes with an “Unknown” ccAFv2 label were removed before
740 metrics were calculated.

741

742 **Computing feature importance**

743 Feature importance for all 861 ccAFv2 features was determined by permuting each feature's
744 expression, running ccAFv2 with the permuted expression matrix, and comparing the likelihoods
745 of all cells for a specific ccAFv2 state to the unpermuted likelihoods of the same cells. The
746 average difference in likelihood was computed for each feature in each ccAFv2 state. A
747 negative average difference in likelihood indicates that a feature was important, and the most
748 negative features are the most important.

749

750 **Comparing ccAFv2 to existing cell cycle classifiers**

751 The performance of ccAFv2 was compared with existing cell cycle state classifiers: ccAF (v1)
752 (O'Connor et al. 2021), Seurat (Hao et al. 2021), Tricycle (Zheng et al. 2022), SchwabeCC
753 (Schwabe et al. 2020; Zheng et al. 2022), reCAT (Liu et al. 2017), Peco (Hsiao et al. 2020) and
754 Cyclone (Scialdone et al. 2015). Each classifier was applied to the PCW 9 R1 (Zeng et al, 2023)
755 and BT322 (Couturier et al. 2020b) scRNA-seq datasets. Data was prepared as required to run
756 each classifier method. The quality of predicted cell cycle states for each classification method
757 was determined by computing the AMI score relative to reference cell cycle states. Ten-fold
758 cross-validation with a 20% hold-out testing set was used to determine the variance of AMI
759 scores for each cell cycle state from each classification method. For both datasets, the ccSeurat
760 predicted cell cycle states were used as the reference for computing AMI scores. Cells with
761 “Unknown” labels were excluded when computing AMI scores. The median AMI scores were
762 tabulated and plotted against the number of predicted states for each classifier. Representative
763 cell cycle state predictions for each classification method were also visualized as UMAPs.

764

765 Because each classifier predicts different numbers of cell cycle states (3 – 8 cell cycle states) it
766 was necessary to use simulated datasets to determine the range of AMI scores that correspond
767 to specific amounts of similarity to the reference. Predicted cell cycle states with 3 to 8 states
768 were simulated that contained specific 0 to 100% similarity to a simulated reference, at 10%

769 increments. The average AMI from 100 simulated cell cycle classifications was computed for
770 each specific amount of similarity to a simulated reference and plotted as a guide to assess the
771 quality between classification methods with different numbers of cell cycle states.
772

773 **Finding the optimal likelihood threshold**

774 A neuroepithelial dataset of *in vivo* hNSCs from fetal tissue at 3 to 12 weeks post-conception
775 from Zeng et al., 2023, that was independent of the ccAFv2 training data, was used to
776 determine the optimal likelihood threshold. Random sub-sampling of 90% of cells for each
777 timepoint was used to determine the variance of the classifications and ccAFv2 was applied with
778 likelihood thresholds ranging from 0.0 to 0.9 by increments of 0.1. For each iteration metrics
779 were collected including the number of cells predicted, and an AMI score computed using
780 ccSeurat cell cycle states as the reference. Cells with “Unknown” labels were excluded when
781 computing AMI scores. Metrics were not computed when 20 or fewer cells were predicted.
782 Student’s *t*-tests were used to compare AMIs computed at each examined likelihood threshold
783 with those derived from a likelihood threshold of 0.0, which is equivalent to not using a likelihood
784 threshold, and a significant difference was considered a p-value ≤ 0.05 . A baseline for
785 comparison was provided by random removal of an equivalent percentage of cells that were
786 classified as “Unknown” for each likelihood threshold, and an AMI was computed with the
787 remaining cells. Student’s *t*-tests were used to compare AMIs of the likelihood thresholded and
788 random removal at each likelihood threshold, and a significant difference was considered a p-
789 value ≤ 0.05 .
790

791 **Cell cycle state validation using hNSCs (PCW 9 R1)**

792 We used the hNSCs collected from whole fetal brain at nine weeks post-conception replicate 1
793 (PCW 9 R1) to validate the cell cycle states assigned by ccAFv2. After quality control
794 (**Supplemental Table S13**) and normalization with sctransform, 5575 cells were classified into
795 distinct cell cycle states using ccAFv2. We selected five key markers of cell cycle states:
796 *CCND1* (Late G1), *CCNE2* (S), *CCNA2* (G2/M), *CCNB1* (G2/M), and *CDK1* (G2/M) to assess
797 the expression patterns associated with these phases. The average expression levels of the
798 genes were calculated and visualized using violin plots, which were grouped according to the
799 cell cycle states predicted by ccAFv2. In addition, we monitored the dynamic changes in the
800 average expression of each key marker as cells transitioned between different cell cycle states.
801 Student’s *t*-tests were used to determine if the marker expression was significantly different at
802 each cell cycle state compared to the G1 state. Finally, relative expression levels of top marker
803 genes for each cell cycle state were identified using *FindAllMarkers()* and visualized using a
804 heatmap, with cells grouped by cell cycle state.
805

806 **Comparison of Neural G0 state with G0 arrest signature using QuieScore**

807 We applied the QuieScore algorithm (<https://github.com/dkornai/QuieScore>) to the U5-hNSCs
808 using the cancer type parameter of “LGG”. The G0 cells were identified by a *q_score_raw* of
809 greater than 3. We evaluated the similarity between the QuieScore-identified G0 cells with the
810 ccAFv2-identified Neural G0 cells using hypergeometric enrichment analysis.
811

812 **Determining the sensitivity of ccAFv2 to missing genes**

813 Sensitivity analysis was conducted on the U5-hNSC dataset by randomly setting a defined
814 percentage of classifier genes (1-90%) to zero and applying the ccAFv2 classifier. Each

815 percentage of classifier genes was subsampled ten times and for each iteration the metrics
816 error rate and percentage of cells predicted were recorded.

817

818 **Demonstrating the generalizability of ccAFv2**

819 The 245,906 human fetal cells 3 to 12 weeks post conception (Zeng et al., 2023) encompassing
820 fifteen cell types that represent all three germ layers (**Supplemental Table S14**) were classified
821 by ccAFv2. Positive marker genes for the Neural G0 cells were identified for each cell type
822 using the FindAllMarkers function (\log_2 fold change ≥ 0.25 ; adjusted p-value ≤ 0.05). The Neural
823 G0 markers were tabulated among each dataset and across all datasets to identify common
824 Neural G0 marker genes.

825

826 **Regressing out cell cycle transcriptional signatures using ccAFv2 marker genes**

827 The average expression from the marker genes for each cell cycle state (**Supplemental Table**
828 **S9**) was computed using the AddModuleScore function in Seurat. The S and G2/M or Late G1,
829 S, S/G2, G2/M, M/Early G1 module scores were regressed out in the SCTransform function in
830 Seurat. The variance explained by the first principal component of the marker genes was used
831 as a metric for co-expression of the cell cycle transcriptional signatures. Empirical p-values were
832 calculated by comparing the observed variance explained to the variance explained of 1,000
833 randomly sampled gene sets of the same size. Significantly regressing out the cell cycle
834 transcriptional signature was determined by a reduction in the variance explained that made the
835 empirical p-value non-significant (>0.05).

836

837 **Application of ccAFv2 to neuroepithelial scRNA-seq and snRNA-seq profiling studies**

838 To maximize overlap with the ccAFv2 input genes, we enabled the option to apply SCTransform
839 (do_sctransform) for SCTransformed datasets. The species ('human' or 'mouse') and gene ID
840 ('Ensembl' or 'symbol') options were configured based on the specifications of each dataset.
841 Predicted cell cycle states were collected from each dataset and integrated with meta
842 information.

843

844 **Application of ccAFv2 to ST-seq data**

845 We downloaded the transcriptome profiles for a 5 μm section of a male C57BL/6 mouse embryo
846 taken from an FFPE tissue block obtained from Charles River Laboratories that was made
847 public by 10x Genomics (<https://www.10xgenomics.com/datasets/visium-hd-cytassist-gene-expression-libraries-of-mouse-embryo>). The 10x Visium HD Gene Expression Library
848 preparation kit afforded a resolution of 2 μm^2 spots and details about sample preparation and
849 library performance and QC can be found on the 10x website linked above. In Seurat the 2 μm^2
850 spots were binned into 8 μm^2 bins, the data log normalized, and ccAFv2 was applied to predict
851 cell cycle states for each spot. Expression of key genes was plotted using the normalized and
852 scaled values.

853

854 **R and Python package for ccAFv2**

855 The ccAFv2 classifier has been implemented as an R package (https://github.com/plaisier-lab/ccafv2_R) that can be installed and used as part of a Seurat workflow, and works for both
856 Seurat version 4 and 5 (**Supplemental Figure S14**). Due to differences in the Seurat v5

859 SCTransform function it was necessary to set the vst.flavor equal to “v1” to make it equivalent to
860 Seurat v4.3.0.1, and leaving the vst.flavor as the default in v5 leads to only small differences
861 (**Supplemental Figure S14**). For the Seurat v5.0.2 the matrixStats package was required to be
862 v1.1.0. Additionally, the ccAFv2 classifier has been implemented as a Python PyPI installed
863 package (<https://pypi.org/project/ccAF/>) that can be installed and used as part of a SCANPY
864 workflow. It should be noted that SCTransform normalization is the suggested method for
865 preparing data that will be classified by ccAFv2, and as of now there is no SCTransform option
866 in SCANPY.

867

868 **Culture of human skeletal muscle satellite cells (hSkMSCs)**

869 hSkMSCs were purchased from ScienCell Research Laboratories (P/N 3510, ScienCell) and
870 were grown in Skeletal Muscle Cell Medium (P/N 3501; ScienCell) on Nunclon Delta-treated cell
871 culture flasks and passaged according to vendor protocols. Cells were detached from their
872 plates using Trypsin/EDTA Solution (P/N 183; ScienCell) and collected with Trypsin Neutralizing
873 Solution (P/N 113; ScienCell).

874

875 **scRNA-seq characterization of hSkMSCs**

876 hSkMSCs were grown up to 80% confluence, washed with Molecular Biology Grade PBS (P/N
877 45001-130, VWR), dissociated with Trypsin/EDTA Solution, and collected in Trypsin
878 Neutralizing Solution. After centrifugation at 300 x g for 5 minutes, cells were resuspended in
879 Molecular Biology Grade PBS containing 0.04% BSA and counted using an automated cell
880 counter. Cells were then diluted to 1,000 cells/µl. scRNA-seq library preparation was performed
881 by the ASU Genomics Core facility. Samples were processed using the 10x Chromium Single
882 Cell 3' Gene Expression v3.1 kit into a single library (10x Genomics). The quality of the library
883 was determined using Agilent TapeStation automated electrophoresis. Samples were
884 sequenced at an average read depth of 100,000 reads per cell (Illumina, Novogene). The 10x
885 Genomics CellRanger v7.0.1 was used to align to the Human reference genome GRCh38-2020-
886 A (GRCh38), quantify, and provide basic quality control metrics for the scRNA-seq data. The
887 10x CellRanger outputs for 7,795 hSkMSCs was loaded into Seurat. Filtering and downstream
888 analyses was done using quality control and downstream processing code templates provided
889 in <https://github.com/plaisier-lab/ccAFv2>. Standard Seurat filters were applied requiring that the
890 cells had to have a least 200 features per cell, and transcripts need to be expressed in at least 3
891 cells. Then the cells were further filtered to 7,207 hSkMSCs by requiring the number of UMIs
892 per cell to fall within the range of 4,000 to 100,000, and the percentage of mitochondrial genes
893 expressed relative to total expression per cell was required to fall within the range of 0.9 to 10%.
894 The filtered cells were then normalized using SCTransform (Hafemeister and Satija 2019),
895 principal components were calculated, and a UMAP was generated.

896

897 **Staining for quiescent-like G0 cells**

898 This staining protocol is based on a protocol developed by Gookin et al., 2017 (Gookin et al.
899 2017) to identify a G0/quiescent subpopulation and has been adapted for fluorescence-
900 activated cell sorting (FACS) and to preserve RNA integrity. Cells were expanded in Nunclon
901 Delta-treated cell culture flasks to achieve the desired cell count, accounting for a 50% loss
902 during staining preparation before downstream FACS.

903
904 Replication stain: Replicating cells were labeled with a synthetic nucleotide
905 Tetramethylrhodamine-dUTP (P/N 17023, AAT Bioquest) transported into the cells using a
906 synthetic nucleotide triphosphate transporter (SNTT) the BioTracker NTP-Transporter Molecule
907 (P/N SCT064, Millipore Sigma) (Zawada et al. 2018; Gookin et al. 2017). When cells reached
908 70% confluence, they were washed with tricine buffer, the SNTT and synthetic nucleotide were
909 diluted so each component was 20 μ M in tricine buffer, added to cells, and incubated at 37°C
910 and 5% CO₂ for 5 minutes to transport fluorescently labeled synthetic nucleotide into the cells.
911 The stain was then aspirated and replaced with complete culture medium, and the cells were
912 incubated at 37°C and 5% CO₂ for 1 hour to allow time for replicating cells to incorporate the
913 fluorescently labeled synthetic nucleotides into their genome's.
914
915 Viability stain: Cells were then washed with Molecular Biology Grade PBS, dissociated with
916 Trypsin/EDTA Solution, and collected in Trypsin Neutralizing Solution. After centrifugation at
917 300 x g for 5 minutes, cells were resuspended in Molecular Biology Grade PBS and counted
918 using an automated cell counter. Cells were centrifuged again, PBS was removed, resuspended
919 in 1:1000 LIVE/DEAD™ Fixable Near-IR Dead Cell Stain Kit (P/N L34975, ThermoFisher
920 Scientific) using manufacturer instructions, and incubated for 30 minutes at room temperature in
921 the dark.
922
923 Fixation and rehydration: Cells were washed with 0.5% Ultra-Pure BSA (P/N AM2616,
924 ThermoFisher Scientific) in Molecular Biology Grade PBS and centrifuged at 300 x g two times.
925 The cells were fixed by first resuspending them in ice-cold Molecular Biology Grade PBS at a
926 volume of 200 μ l per 1 million cells. Then, ice-cold 100% methanol was added dropwise at a
927 volume of 800 μ l per 1 million cells, with gentle shaking. Cells were then incubated for at least
928 30 minutes at -20°C. After fixation, cells were kept on ice. Cells were rehydrated with cold 3X
929 SSC Rehydration Cocktail (Chen et al. 2018), followed by centrifugation at 500 x g for 5
930 minutes. Cells were washed one more time with the SSC Rehydration Cocktail, and one time
931 with 0.5% Ultra-Pure BSA in Molecular Biology Grade PBS.
932
933 Phosphorylated RB (pRB) staining: Hypo-phosphorylation of pRB is an established indicator of
934 a cell being in a quiescent G0 state (Gookin et al. 2017). Primary antibody for pRB (Ser807/811)
935 (P/N 8516T, Cell Signaling Technologies) was added at a dilution of 1:200 in 0.5% Ultra-Pure
936 BSA in Molecular Biology Grade PBS and incubated overnight at 4°C. Cells were washed three
937 times with 0.5% Ultra-Pure BSA in Molecular Biology Grade PBS and then fluorescently labeled
938 secondary antibody (P/N 4412, Cell Signaling) was added at a dilution of 1:1000 dilution for 30
939 minutes. Samples were then washed two times with 0.5% Ultra-Pure BSA in Molecular Biology
940 Grade PBS.
941
942 DNA staining: The ploidy of cells was determined using Hoechst DNA stain (Gookin et al.
943 2017). Prior to FACS cells were stained with 2 ug/ml of Hoechst DNA stain (P/N 561908, BD)
944 diluted in Molecular Biology Grade PBS, without BSA.
945
946 **Fluorescence-activated cell sorting of G0 cells**

947 Cells were filtered using sterile CellTrics 30 μ m filters (P/N 04-004-2326, Sysmex) into sterile,
948 nuclease-free 5 ml polystyrene round-bottom tubes for sorting (P/N 352235, Corning) and kept
949 covered from light and on ice until sorting. Cells were stained using the following experimental
950 design to define gates and have adequate controls: 1) Hoechst only, 2) Live/Dead only, 3)
951 replication only, 4) pRB only, 5) replication fluorescence minus one (FMO), 6) pRB FMO, and 7)
952 all stains. Example of gating can be seen in **Figure 4A**. Cell sorting was performed using the
953 FACSymphony flow cytometer (BD). G0 cells were defined as viable cells, that were diploid
954 (2N), with low pRB. FACS data analysis was performed using FlowJo (BD).

955

956 **RNA-sequencing of G0 cells**

957 RNA was extracted from 400,000 sorted cells from two biological replicates using the Qiagen
958 RNeasy Micro Kit (P/N 74004, Qiagen). The concentration and quality of RNA was determined
959 by Nanodrop (Thermo Scientific) and High Sensitivity RNA TapeStation (Agilent). Both G0
960 samples had more than 300 ng of RNA and RIN scores of greater than 9. Samples were sent
961 for sequencing on the NovaSeq X Plus (Illumina, Novogene). A docker RNA-seq pipeline
962 (cplaisier/star_2_7_1a_grch38_p21; DOI = <https://doi.org/10.5281/zenodo.5519663>) was employed to align reads from FASTQ
963 files to the genome using STAR v2.7.1a (Dobin et al. 2013) and GENCODE genome build
964 GRCh38 Release 31 (Frankish et al. 2023). Counts were tabulated using htseq-count (Putri et
965 al. 2022). DESeq2 (Love, Huber, and Anders 2014) was used for subsequent differential gene
966 expression analysis.

967

968 **Correlating RNA-seq and scRNA-seq data**

969 DESeq2-normalized RNA-seq data and sctransform-normalized scRNA-seq data were loaded
970 into R. Marker genes were selected by identifying highly variable scRNA-seq genes with more
971 than 10 counts in the bulk G0 subpopulations. Additionally, 861 ccAFv2 marker genes present
972 in the RNA-seq data were included. Both scRNA-seq and RNA-seq datasets were filtered to
973 include these 3,120 marker genes. The expression profiles of individual cells in the scRNA-seq
974 data were correlated with the G0 RNA-seq profiles using the corr package in R (Makowski et al.,
975 2020), using the Spearman method.

976

977

978

979

980 **Data Access**

981 All raw and processed sequencing data for the hSkMSC FACs sorted RNA-seq and scRNA-seq
982 generated in this study have been submitted to NCIB Gene Expression Omnibus (GEO;
983 <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE285220. All raw and
984 processed sequencing data for the four LGG scRNA-seq generated in this study have been
985 submitted to GEO under accession number GSE263796.

986
987 All other data used in our analyses are available on Zenodo
988 (<https://zenodo.org/doi/10.5281/zenodo.10963136>). We also provide all code on github.com
989 (<https://github.com/plaisier-lab/ccAFv2>) and Docker images on DockerHub that were used to run
990 all analyses (https://hub.docker.com/r/cplaisier/ccafv2_extra and
991 <https://hub.docker.com/r/cplaisier/ccnn>).

992
993 **R package for ccAFv2**

994 We have developed an R package that can be installed using devtools from github. The
995 instructions for installation and usage can be found on github: <https://github.com/plaisier->
996 [lab/ccafv2_R](https://github.com/plaisier-lab/ccafv2_R)

997
998 **Python package for ccAFv2**

999 We have also developed a Python package that can be installed using pip from PyPi. The
1000 instructions for installation and usage can be found on PyPi and github:
1001 <https://pypi.org/project/ccAFv2/> and https://github.com/plaisier-lab/ccAFv2_py

1002
1003 **Docker images for ccAFv2**

1004 We also provide Docker images that include all dependencies and ccAFv2 preinstalled to make
1005 the package more user friendly. Please see the github repositories for information about how to
1006 get, run, and use the Docker images.

- 1007 - R package:
 - 1008 o Seurat v4: https://hub.docker.com/r/cplaisier/ccafv2_seurat4
 - 1009 o Seurat v5: https://hub.docker.com/r/cplaisier/ccafv2_seurat5
- 1010 - Python package: https://hub.docker.com/r/cplaisier/ccafv2_py

1011 **Competing interest statement**

1012 The authors declare no competing interests.

1013 **Acknowledgments**

1014 This work was supported by the following grants: NINDS/NIH (R01NS119650: A.P., P.P.,
1015 C.L.P.) and (R01NS123038: CLP); NCI/NIH (R01CA190957; R21CA170722; P30CA15704:
1016 P.P.); DoD Translational New Investigator Award (CA100735: P.P.); and the Pew Biomedical
1017 Scholars Program (P.P.). The authors acknowledge Joy Blain, Anna Engelbrektson, and Adam
1018 Kindelin for their assistance in scRNA-seq library preparation, quality control, and FACS. The
1019 BT237 cell line was provided by Keith Ligon from the Dana-Farber Cancer Institute Center for
1020 Patient Derived Models. Some components of Figure 5 were made using Biorender.com.

1021 **Author contributions**

1022 Project conception and experimental design were carried out by CLP, SAO, PJP, and AP.
1023 Implementation of ccAFv2 and the R and Python packages was performed jointly by SAO, RH
1024 and CLP. scRNA-seq of LGG primary cell lines was performed by LG and J-PH. Single
1025 cell/nuclei dataset curation was performed by SAO under the supervision of CLP. Testing of
1026 ccAFv2 was performed by SAO under the supervision of CLP. ST-seq dataset curation and
1027 testing was performed by CLP. Interpretation of ST-seq results was performed by BBB and
1028 CLP. The manuscript was written by SAO, BBB, and CLP with input from all authors.
1029

1030 **References**

1031 Adams NM, Grassmann S, Sun JC. 2020. Clonal expansion of innate and adaptive
1032 lymphocytes. *Nat Rev Immunol* **20**: 694–707.

1033 Arnold SJ, Huang G-J, Cheung AFP, Era T, Nishikawa S-I, Bikoff EK, Molnár Z, Robertson EJ,
1034 Groszer M. 2008. The T-box transcription factor Eomes/Tbr2 regulates neurogenesis in
1035 the cortical subventricular zone. *Genes Dev* **22**: 2479–2484.

1036 Augustus M, Pineau D, Aimond F, Azar S, Lecca D, Scamps F, Muxel S, Darlix A, Ritchie W,
1037 Gozé C, et al. 2021. Identification of CRYAB+ KCNN3+ SOX9+ Astrocyte-Like and
1038 EGFR+ PDGFRA+ OLIG1+ Oligodendrocyte-Like Tumoral Cells in Diffuse IDH1-Mutant
1039 Gliomas and Implication of NOTCH1 Signalling in Their Genesis. *Cancers* **13**: 2107.

1040 Ballabeni A, Park I-H, Zhao R, Wang W, Lerou PH, Daley GQ, Kirschner MW. 2011. Cell cycle
1041 adaptations of embryonic stem cells. *Proc Natl Acad Sci U S A* **108**: 19252–19257.

1042 Belmonte-Mateos C, Pujades C. 2021. From Cell States to Cell Fates: How Cell Proliferation
1043 and Neuronal Differentiation Are Coordinated During Embryonic Development. *Front
1044 Neurosci* **15**: 781160.

1045 Bressan RB, Dewari PS, Kalantzaki M, Gangoso E, Matjusaitis M, Garcia-Diaz C, Blin C, Grant
1046 V, Bulstrode H, Gogolok S, et al. 2017. Efficient CRISPR/Cas9-assisted gene targeting
1047 enables rapid and precise genetic manipulation of mammalian neural stem cells. *Dev
1048 Camb Engl* **144**: 635–648.

1049 Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA,
1050 Marioni JC, Stegle O. 2015. Computational analysis of cell-to-cell heterogeneity in
1051 single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol*
1052 **33**: 155–160.

1053 Bukovsky A, Caudle MR, Svetlikova M, Wimalasena J, Ayala ME, Dominguez R. 2005.
1054 Oogenesis in adult mammals, including humans. *Endocrine* **26**: 301–316.

1055 Butler A, Hoffman P, Smibert P, Papalex E, Satija R. 2018. Integrating single-cell transcriptomic
1056 data across different conditions, technologies, and species. *Nat Biotechnol* **36**: 411–420.

1057 Cebrian-Silla A, Nascimento MA, Redmond SA, Mansky B, Wu D, Obernier K, Romero
1058 Rodriguez R, Gonzalez-Granero S, García-Verdugo JM, Lim DA, et al. 2021. Single-cell
1059 analysis of the ventricular-subventricular zone reveals signatures of dorsal and ventral
1060 adult neurogenesis eds. J.G. Gleeson, M.E. Bronner, and D. Van der Kooy. *eLife* **10**:
1061 e67436.

1062 Chen J, Cheung F, Shi R, Zhou H, Lu W, Candia J, Kotliarov Y, Stagliano KR, Tsang JS, CHI
1063 Consortium. 2018. PBMC fixation and processing for Chromium single-cell RNA
1064 sequencing. *J Transl Med* **16**: 198.

1065 Couturier CP, Ayyadhyury S, Le PU, Nadaf J, Monlong J, Riva G, Allache R, Baig S, Yan X,
1066 Bourgey M, et al. 2020. Single-cell RNA-seq reveals that glioblastoma recapitulates a
1067 normal neurodevelopmental hierarchy. *Nat Commun* **11**: 3406.

1068 Davis AA, Temple S. 1994. A self-renewing multipotential stem cell in embryonic rat cerebral
1069 cortex. *Nature* **372**: 263–266.

1070 Doetsch F. 2003. A niche for adult neural stem cells. *Curr Opin Genet Dev* **13**: 543–550.

1071 Dulken BW, Leeman DS, Boutet SC, Hebestreit K, Brunet A. 2017. Single-Cell Transcriptomic
1072 Analysis Defines Heterogeneity and Transcriptional Dynamics in the Adult Neural Stem
1073 Cell Lineage. *Cell Rep* **18**: 777–790.

1074 Englund C, Fink A, Lau C, Pham D, Daza RAM, Bulfone A, Kowalczyk T, Hevner RF. 2005.
1075 Pax6, Tbr2, and Tbr1 are expressed sequentially by radial glia, intermediate progenitor
1076 cells, and postmitotic neurons in developing neocortex. *J Neurosci Off J Soc Neurosci*
1077 **25**: 247–251.

1078 Gookin S, Min M, Phadke H, Chung M, Moser J, Miller I, Carter D, Spencer SL. 2017. A map of
1079 protein dynamics during cell-cycle progression and cell-cycle exit. *PLOS Biol* **15**:
1080 e2003268.

1081 Guo J, Grow EJ, Mlcochova H, Maher GJ, Lindskog C, Nie X, Guo Y, Takei Y, Yun J, Cai L, et
1082 al. 2018. The adult human testis transcriptional cell atlas. *Cell Res* **28**: 1141–1157.

1083 Hanahan D, Weinberg RA. 2000. The Hallmarks of Cancer. *Cell* **100**: 57–70.

1084 Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby C,
1085 Zager M, et al. 2021. Integrated analysis of multimodal single-cell data. *Cell* **184**: 3573–
1086 3587.e29.

1087 Harada Y, Yamada M, Imayoshi I, Kageyama R, Suzuki Y, Kuniya T, Furutachi S, Kawaguchi D,
1088 Gotoh Y. 2021. Cell cycle arrest determines adult neural stem cell ontogeny by an
1089 embryonic Notch-nonoscillatory Hey1 module. *Nat Commun* **12**: 6562.

1090 Hoffmann I, Balling R. 1995. Cloning and expression analysis of a novel mesodermally
1091 expressed cadherin. *Dev Biol* **169**: 337–346.

1092 Hsiao CJ, Tung P, Blischak JD, Burnett JE, Barr KA, Dey KK, Stephens M, Gilad Y. 2020.
1093 Characterizing and inferring quantitative cell cycle phase in single-cell RNA-seq data
1094 analysis. *Genome Res* **30**: 611–621.

1095 Johe KK, Hazel TG, Muller T, Dugich-Djordjevic MM, McKay RD. 1996. Single factors direct the
1096 differentiation of stem cells from the fetal and adult central nervous system. *Genes Dev*
1097 **10**: 3129–3140.

1098 Joseph C, Mangani AS, Gupta V, Chitranshi N, Shen T, Dheer Y, KB D, Mirzaei M, You Y,
1099 Graham SL, et al. 2020. Cell Cycle Deficits in Neurodegenerative Disorders: Uncovering
1100 Molecular Mechanisms to Drive Innovative Therapeutic Development. *Aging Dis* **11**:
1101 946–966.

1102 Kriegstein AR, Noctor SC. 2004. Patterns of neuronal migration in the embryonic cortex. *Trends
1103 Neurosci* **27**: 392–399.

1104 Leng N, Chu L-F, Barry C, Li Y, Choi J, Li X, Jiang P, Stewart RM, Thomson JA, Kendziora C.
1105 2015. Oscope identifies oscillatory genes in unsynchronized single cell RNA-seq
1106 experiments. *Nat Methods* **12**: 947–950.

1107 Liang G, Xu E, Yang C, Zhang C, Sheng X, Zhou X. 2015. Nucleosome-binding protein HMGN2
1108 exhibits antitumor activity in human SaO2 and U2-OS osteosarcoma cell lines. *Oncol*
1109 *Rep* **33**: 1300–1306.

1110 Lim DA, Alvarez-Buylla A. 2016. The Adult Ventricular–Subventricular Zone (V-SVZ) and
1111 Olfactory Bulb (OB) Neurogenesis. *Cold Spring Harb Perspect Biol* **8**: a018820.

1112 Liu Z, Lou H, Xie K, Wang H, Chen N, Aparicio OM, Zhang MQ, Jiang R, Chen T. 2017.
1113 Reconstructing cell cycle pseudo time-series via single-cell transcriptome data. *Nat*
1114 *Commun* **8**: 22.

1115 Llorens-Bobadilla E, Zhao S, Baser A, Saiz-Castro G, Zwadlo K, Martin-Villalba A. 2015. Single-
1116 Cell Transcriptomics Reveals a Population of Dormant Neural Stem Cells that Become
1117 Activated upon Brain Injury. *Cell Stem Cell* **17**: 329–340.

1118 Luecken MD, Theis FJ. 2019. Current best practices in single-cell RNA-seq analysis: a tutorial.
1119 *Mol Syst Biol* **15**: e8746.

1120 Ma F, Pellegrini M. 2020. ACTINN: automated identification of cell types in single cell RNA
1121 sequencing. *Bioinforma Oxf Engl* **36**: 533–538.

1122 Martínez-Cerdeño V, Noctor SC, Kriegstein AR. 2006. The role of intermediate progenitor cells
1123 in the evolutionary expansion of the cerebral cortex. *Cereb Cortex N Y N* **16 Suppl**
1124 1: i152-161.

1125 Maxwell CA, Keats JJ, Belch AR, Pilarski LM, Reiman T. 2005. Receptor for hyaluronan-
1126 mediated motility correlates with centrosome abnormalities in multiple myeloma and
1127 maintains mitotic integrity. *Cancer Res* **65**: 850–860.

1128 Nowakowski TJ, Bhaduri A, Pollen AA, Alvarado B, Mostajo-Radji MA, Di Lullo E, Haeussler M,
1129 Sandoval-Espinosa C, Liu SJ, Velmeshev D, et al. 2017. Spatiotemporal gene
1130 expression trajectories reveal developmental hierarchies of the human cortex. *Science*
1131 **358**: 1318–1323.

1132 Obernier K, Cebrian-Silla A, Thomson M, Parraguez JI, Anderson R, Guinto C, Rodriguez JR,
1133 Garcia-Verdugo J-M, Alvarez-Buylla A. 2018. Adult neurogenesis is sustained by
1134 symmetric self-renewal and differentiation. *Cell Stem Cell* **22**: 221-234.e8.

1135 O'Connor SA, Feldman HM, Arora S, Hoellerbauer P, Toledo CM, Corrin P, Carter L, Kufeld M,
1136 Bolouri H, Basom R, et al. 2021. Neural G0: a quiescent-like state found in
1137 neuroepithelial-derived cells and glioma. *Mol Syst Biol* **17**: e9522.

1138 Pauklin S, Vallier L. 2013. The Cell-Cycle State of Stem Cells Determines Cell Fate Propensity.
1139 *Cell* **155**: 135–147.

1140 Qiu P. 2020. Embracing the dropouts in single-cell RNA-seq analysis. *Nat Commun* **11**: 1169.

1141 1142 R Core Team (2025). R: A language and environment for statistical computing. R Foundation for
Statistical Computing, Vienna, Austria.

1143 1144 Scholzen T, Gerdes J. 2000. The Ki-67 protein: from the known and the unknown. *J Cell Physiol*
182: 311–322.

1145 1146 Schwabe D, Formichetti S, Junker JP, Falcke M, Rajewsky N. 2020. The transcriptome
dynamics of single cells during the cell cycle. *Mol Syst Biol* **16**: e9946.

1147 1148 1149 Scialdone A, Natarajan KN, Saraiva LR, Proserpio V, Teichmann SA, Stegle O, Marioni JC,
Buettner F. 2015. Computational assignment of cell-cycle stage from single-cell
transcriptome data. *Methods* **85**: 54–61.

1150 1151 1152 Slyper M, Porter CBM, Ashenberg O, Waldman J, Drokhlyansky E, Wakiro I, Smillie C, Smith-
Rosario G, Wu J, Dionne D, et al. 2020. A single-cell and single-nucleus RNA-Seq
toolbox for fresh and frozen human tumors. *Nat Med* **26**: 792–802.

1153 1154 1155 Soufi A, Dalton S. 2016. Cycling through developmental decisions: how cell cycle dynamics
control pluripotency, differentiation and reprogramming. *Dev Camb Engl* **143**: 4301–
4311.

1156 1157 1158 Stallaert W, Taylor SR, Kedziora KM, Taylor CD, Sobon HK, Young CL, Limas JC, Varblow
Holloway J, Johnson MS, Cook JG, et al. 2022. The molecular architecture of cell cycle
arrest. *Mol Syst Biol* **18**: e11087.

1159 1160 1161 1162 Tejero R, Huang Y, Katsyv I, Kluge M, Lin J-Y, Tome-Garcia J, Daviaud N, Wang Y, Zhang B,
Tsankova NM, et al. 2019. Gene signatures of quiescent glioblastoma cells reveal
mesenchymal shift and interactions with niche microenvironment. *EBioMedicine* **42**:
252–269.

1163 1164 1165 Tullai JW, Schaffer ME, Mullenbrock S, Sholder G, Kasif S, Cooper GM. 2007. Immediate-early
and delayed primary response genes are distinct in function and genomic architecture. *J
Biol Chem* **282**: 23981–23995.

1166 1167 Vong QP, Cao K, Li HY, Iglesias PA, Zheng Y. 2005. Chromosome alignment and segregation
regulated by ubiquitination of survivin. *Science* **310**: 1499–1504.

1168 White J, Dalton S. 2005. Cell cycle control of embryonic stem cells. *Stem Cell Rev* **1**: 131–138.

1169 1170 1171 Wiecek AJ, Cutty SJ, Kornai D, Parreno-Centeno M, Gourmet LE, Tagliazucchi GM, Jacobson
DH, Zhang P, Xiong L, Bond GL, et al. 2023. Genomic hallmarks and therapeutic
implications of G0 cell cycle arrest in cancer. *Genome Biol* **24**: 128.

1172 1173 1174 Xie P, Gao M, Wang C, Zhang J, Noel P, Yang C, Von Hoff D, Han H, Zhang MQ, Lin W. 2019.
SuperCT: a supervised-learning framework for enhanced characterization of single-cell
transcriptomic profiles. *Nucleic Acids Res* **47**: e48.

1175 1176 1177 Zawada Z, Tatar A, Mocilac P, Buděšínský M, Kraus T. 2018. Transport of Nucleoside
Triphosphates into Cells by Artificial Molecular Transporters. *Angew Chem Int Ed Engl*
57: 9891–9895.

1178 Zeng B, Liu Z, Lu Y, Zhong S, Qin S, Huang L, Zeng Y, Li Z, Dong H, Shi Y, et al. 2023. The
1179 single-cell and spatial transcriptional landscape of human gastrulation and early brain
1180 development. *Cell Stem Cell* **30**: 851-866.e7.

1181 Zhang Q, Wu X, Fan Y, Jiang P, Zhao Y, Yang Y, Han S, Xu B, Chen B, Han J, et al. 2021.
1182 Single-cell analysis reveals dynamic changes of neural cells in developing human spinal
1183 cord. *EMBO Rep* **22**: e52728.

1184 Zheng SC, Stein-O'Brien G, Augustin JJ, Slosberg J, Carosso GA, Winer B, Shin G, Bjornsson
1185 HT, Goff LA, Hansen KD. 2022. Universal prediction of cell-cycle position using transfer
1186 learning. *Genome Biol* **23**: 41.

1187

1188

1189 **Figure Legends**

1190

1191 **Figure 1.** Implementing and testing the ccAFv2 classifier. **A.** The design of the Artificial Neural
1192 Network (ANN) implemented for the ccAFv2. Expr. = expression, ReLU = Rectified Linear Units.
1193 **B.** Method designed to determine the predicted class from the likelihoods generated by running
1194 expression data from a single cell through the ccAFv2 ANN. **C.** Comparison of five different
1195 classification methods using F1 scores (a metric that integrates precision and recall, and has a
1196 maximum value of 1), from the 10-fold cross validation analysis of training on the U5-hNSCs.
1197 The F1 scores are computed for each cell cycle state from each of the 10 testing datasets. **D.**
1198 Determining the optimal number of neurons in each hidden layer using average U5-hNSC F1
1199 score across cell cycle states on the x-axis, and the average AMI score across the remaining
1200 datasets (U5-hNSCs; glioma stem cells: BT322, BT324, BT326, BT333, BT363, BT368; tumor
1201 cells: BT363, BT368; and Grade 2 Astrocytoma: LGG275). Each combination of hidden layer
1202 neurons is labeled using: number of hidden layer one neurons / number of hidden layer two
1203 neurons. The chosen optimal configuration of 600 hidden layer 1 neurons and 200 hidden layer
1204 2 neurons (600 / 200) is denoted in red. **E.** UMAP of U5-hNSCs with cells colored by the labels
1205 from O'Connor et al., 2021. **F-L.** The top 15 most important features for the ccAFv2 classifier
1206 were identified based on the mean change (Δ) in likelihood after permuting each feature's
1207 expression. A negative mean change in likelihood indicates that the feature increased the
1208 likelihood of predicting a ccAFv2 state.

1209

1210 **Figure 2.** Comparing the performance of ccAFv2 to existing cell cycle state classifiers. **A.**
1211 Median AMI score for each cell cycle classifier's predictions of the hNSCs from a whole fetal
1212 brain at 9 weeks post conception (PCW 9 R1; Zeng et al., 2023) relative to the ccSeurat cell
1213 cycle states is plotted against the number of cell cycle states predicted by the classifier. The
1214 average similarity to the reference was computed, based on the number of cell cycle states in
1215 the reference and predicted by the classifier, and were plotted at 10 percent intervals to facilitate
1216 comparison between classifiers with differing numbers of predicted cell cycle states. **B.** Overlay
1217 of representative cell cycle state predictions on the hNSCs from a whole fetal brain at PCW 9
1218 R1. **C.** Median AMI score for each cell cycle classifier's predictions of the glioma stem cell line
1219 BT322 relative to the ccSeurat cell cycle states is plotted against the number of cell cycle states
1220 predicted by the classifier. Again, average similarity to the reference was computed based on
1221 the number of cell cycle states in the reference and predicted by the classifier and were plotted
1222 at 10 percent intervals to facilitate comparison between classifiers with differing numbers of
1223 predicted cell cycle states. **D.** Overlay of representative cell cycle state predictions on the tumor
1224 cells of BT322.

1225

1226 **Figure 3.** Application of ccAFv2 to *in vivo* hNSCs from fetal tissue 3 to 12 weeks post
1227 conception. **A.** Proportions of cell cycle states in U5-hNSCs which were grown *in vitro* and were
1228 derived from a human fetus at 8 PCW for both ccAFv2 and ccSeurat. **B.** Proportions of cell
1229 cycle states of hNSCs extracted from 3 to 12 PCW fetal tissue for both ccAFv2 and ccSeurat
1230 (Zeng et al., 2023). **C-F.** Distribution of cyclin expression in the *in vivo* hNSCs from a whole
1231 human fetal brain at PCW 9 R1 grouped by cell cycle phase. **G.** Mean expression of cyclins
1232 across the ccAFv2 cell cycle phases in cells from a whole human fetal brain at PCW 9 R1. Red

1233 points denote the ccAFv2 cell cycle state with the highest average expression. Gene expression
1234 levels at each cell cycle state were compared to those in G1 cells using Student's *t*-test (****
1235 indicates $p \leq 0.0001$). **H.** Expression of ccAFv2 marker genes for each cell cycle state in hNSCs
1236 from a human whole fetal brain at PCW 9 R1. Important genes names are denoted in dark red.
1237 **I.** Testing different likelihood thresholds 0.0 to 0.9 using AMI score and percent of cells
1238 predicted as the metrics. Dashed red line indicates 90 percent of cells were predicted, and red
1239 dot indicates significantly improved AMI score due to applying threshold. **J.** Comparison of
1240 likelihood threshold application to random removal of the same number of cell predictions for *in*
1241 *vivo* hNSCs from a human whole fetal brain at PCW 9 R1. Metric used for assessment is the
1242 AMI score. Likelihood thresholds start at 0.3 on the x-axis because AMI values at likelihood
1243 thresholds 0 to 0.3 are the same. AMI scores at each likelihood threshold were compared using
1244 Student's *t*-test (**** indicates $p \leq 0.0001$). Rep. = biological replicate.
1245

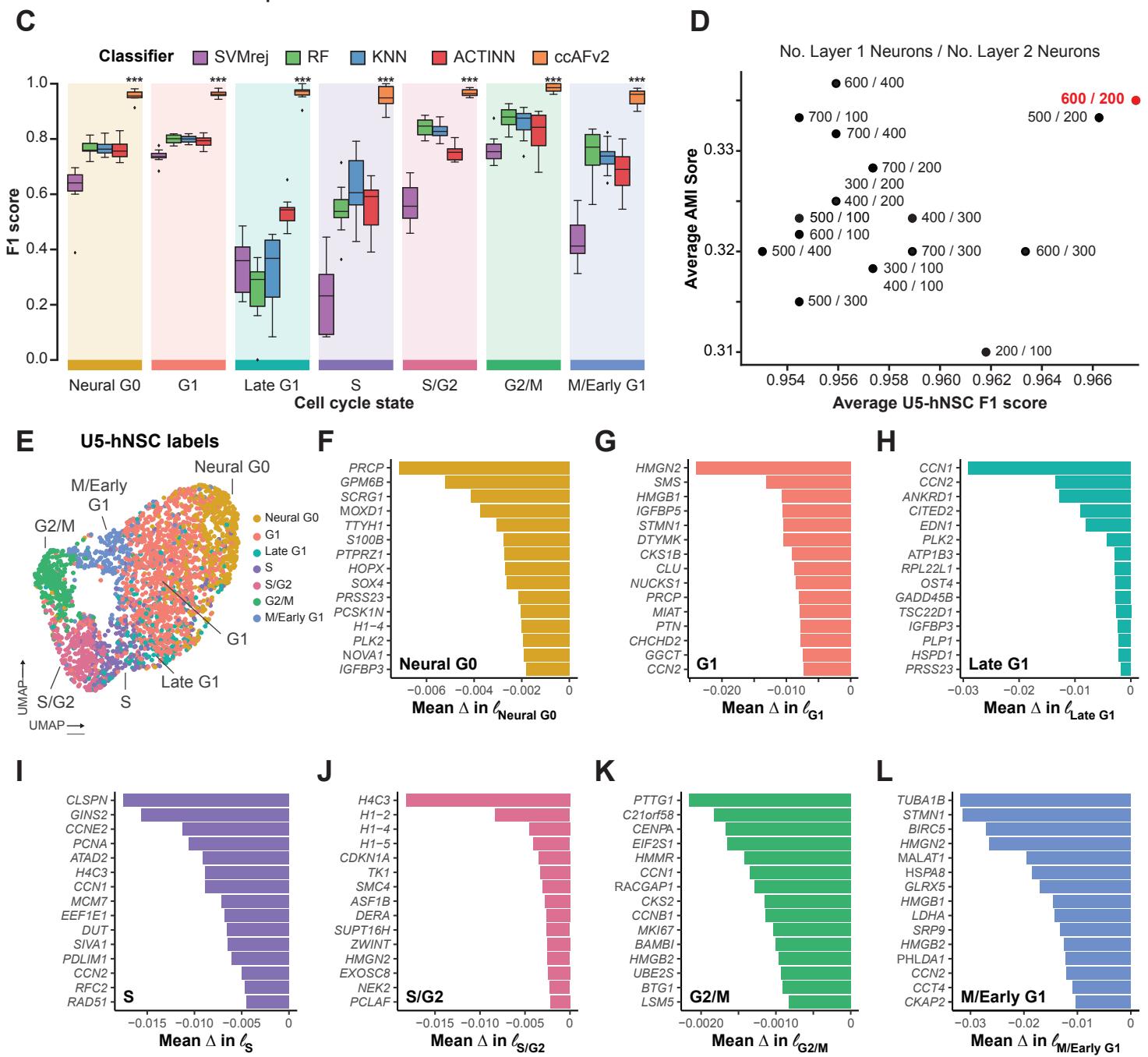
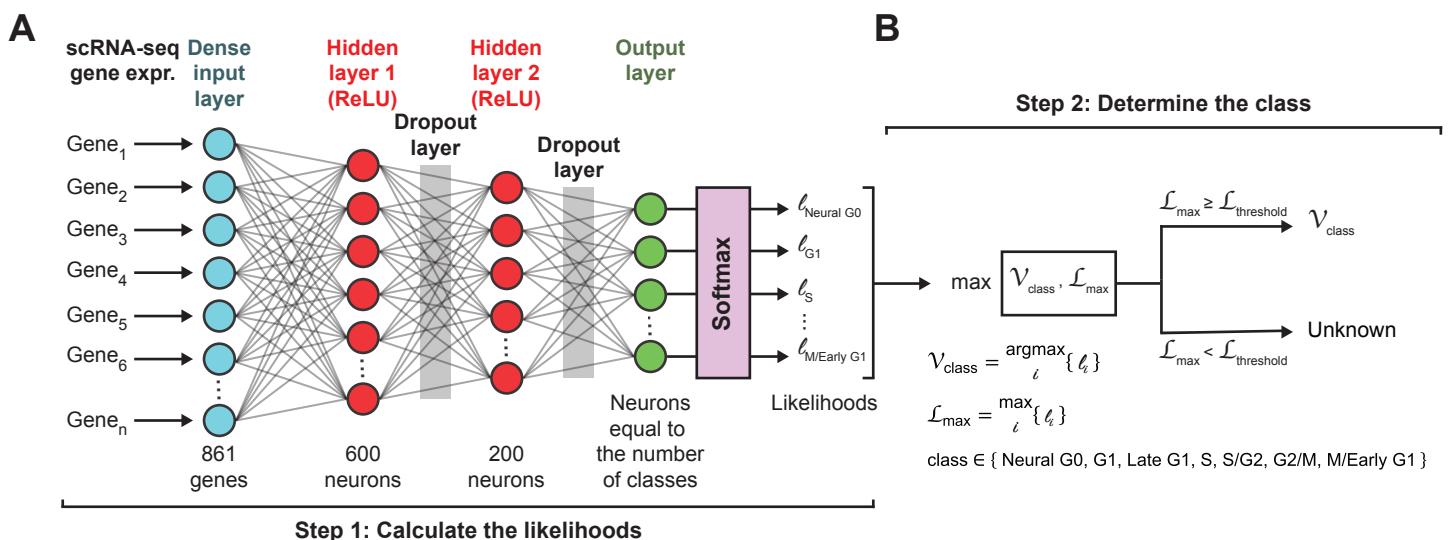
1246 **Figure 4.** Experimental enrichment of mesenchymal G0 cells from hSkMSCs using FACs. **A.**
1247 Gating strategy for isolating mesenchymal G0 cells from hSkMSCs. First gated on single cells,
1248 then live cells, next diploid cells that are not replicating are selected, and finally cells with hypo-
1249 phosphorylation of RB are selected. **B.** ccAFv2 applied to unsorted hSkMSCs. The red dashed
1250 area encompasses the UMAP area containing the vast majority of Neural G0, G1, and Late G1
1251 cells, these three states do not exhibit consistent clustering. **C.** Neural G0 cells are highlighted
1252 in color, while all other states are shown in gray. **D-E.** Cells are colored based on their
1253 Spearman correlation coefficient with the scRNA-seq expression profiles and the RNA-seq
1254 profiles of flow-sorted mesenchymal G0 cells (pRB-, diploid, non-replicating) from two biological
1255 replicates. **F.** Test of which ccAFv2 cell cycle states were significantly enriched with
1256 mesenchymal G0 cells, and not mesenchymal G0 cells. Mesenchymal G0 cells are defined by a
1257 Spearman correlation ≥ 0.1 in both replicates, while non-G0 cells are defined by a correlation \leq
1258 0.1 in one or both replicates. Values are represented as the negative logarithm of the p-value.
1259 **G.** Percentage of ccAFv2 states in mesenchymal G0 and not G0 cells. **H.** Differential expression
1260 of genes between mesenchymal G0 cells versus not G0 cells. Each dot represents one gene (n
1261 = 22,845). Dotted lines denote $\log_2(\text{fold change})$ and adjusted p-value cutoffs to identify
1262 significant marker genes ($\log_2\text{FC} \geq 0.5$; $p\text{-adj} \leq 0.05$). Red dots denote ccAFv2 Neural G0
1263 marker genes. Labeled genes are marker genes for hSkMSC G0 cells that overlap with ccAFv2
1264 Neural G0 marker genes.
1265

1266 **Figure 5.** Application of ccAFv2 to the transcriptomes of 245,906 single cells derived from
1267 human fetuses aged 3 to 12 PCW. **A.** The 15 different cell types included in the analysis
1268 encompass all three germ layers. For each cell type the number of cells is given. **B.** Percentage
1269 of each ccAFv2 predicted state for each cell type. **C.** Percentage of each ccAFv2 predicted state
1270 for each cell type when Neural G0, G1, and Late G1 are binned. **D-G.** Z-score normalized cyclin
1271 expression across 15 cell types. Thin lines represent individual cell types, while thick lines
1272 indicate the average Z-score normalized cyclin expression for each germ layer. Lines are color-
1273 coded according to their corresponding germ layer.
1274

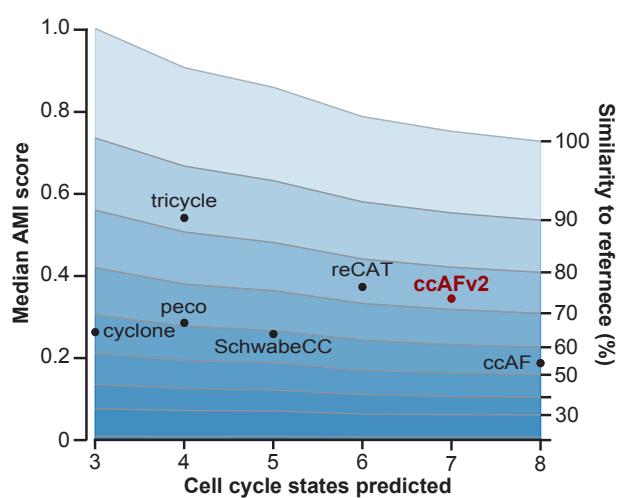
1275 **Figure 6.** Application of ccAFv2 to single cells and nuclei from human and mice. **A.** Summary
1276 schematic of data ccAFv2 can be applied to and suggested data preparation. **B.** Proportion of

1277 cells assigned to each cell cycle state for scRNA-seq data from the developing human
1278 telencephalon (Nowakowski et al, 2017). **C.** Proportions of cell cycle states from scRNA-seq
1279 from the ventricular (V)-SVZ of the adult mouse brain (Cebrian-Silla et al, 2021). Prog. =
1280 progenitors, Inh = inhibitor, Ex = excitatory, NSPCs = neural stem/progenitor cells, IntProg. =
1281 intermediate progenitor cells. **D.** Proportions of cell cycle states from scRNA-seq from GLAST
1282 and PROM1 flow-sorted cells from the subventricular zone (SVZ) of mice (Llorens-Bobadilla et
1283 al, 2015), and EGFR, GFAP, and PROM1 flow-sorted cells from the subventricular zone (SVZ)
1284 of adult mice (Dulken et al, 2017). qNSC1 = dormant quiescent neural stem cell, qNSC2 =
1285 primed-quiescent neural stem cell, aNSC1 = active neural stem cell, aNSC2 = actively dividing
1286 neural stem cell. qNSC = quiescent neural stem cell, aNSC = active neural stem cell. **E.**
1287 Proportions of cell cycle states from scRNA-seq (C) and snRNA-seq (N) from spinal, cervical,
1288 lumbar, and thoracic regions from the developing human spinal cord at 8, 10, 11, 20, and 23
1289 PCW (Zhang et al, 2021).
1290

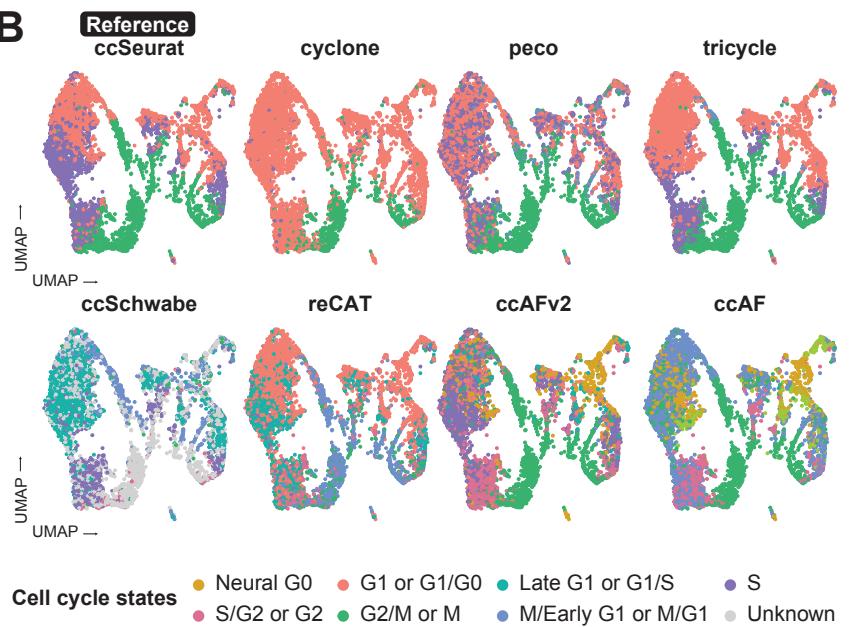
1291 **Figure 7.** Application of ccAFv2 to spatial transcriptomics data from a male C57BL/6 mouse
1292 embryo at E15.5. **A.** H&E staining for the whole embryo. **B.** Spatial overlay of the predicted
1293 ccAFv2 states onto the whole embryo. **C.** Spatial expression of the cell cycle marker gene
1294 *Mki67* for the whole embryo. The black boxes in panels **A** through **C** indicate the region of the
1295 developing cortex that was magnified in panels **D** through **Q**. The developmental regions of the
1296 developing cortex are denoted on the side: Dermis = developing skin, Skull = developing skull,
1297 CP/MZ = cortical plate and marginal zone, IZ = intermediate zone, SVZ = subventricular zone,
1298 VZ = ventricular zone. **D.** H&E staining for the developing embryo cortex. **E.** Spatial overlay of
1299 the predicted ccAFv2 states onto the developing embryo cortex. **F-I.** Expression of key marker
1300 genes describing the developmental regions in the developing embryo cortex. **J.** Spatial
1301 expression of the cell cycle marker gene *Mki67* in the developing embryo cortex. **K-Q.**
1302 Likelihoods for each of the cell cycle states spatially overlayed onto the developing embryo
1303 cortex. The magnitude of the likelihood indicates the probability that a cell with that cell cycle
1304 state underlies that spot of the spatial array.



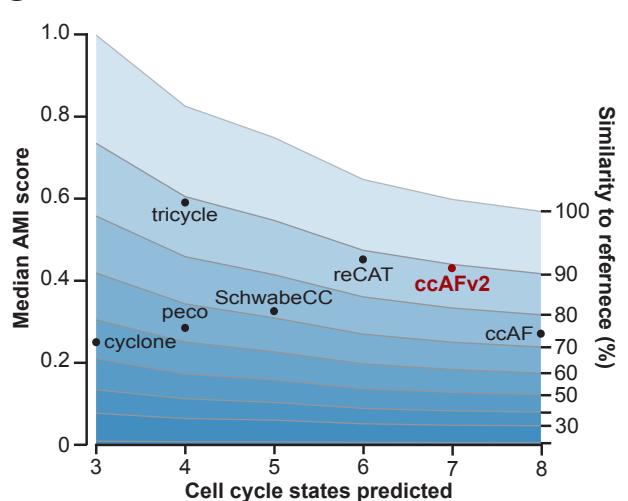
A



B



C



D

