

Title: Coexpression enables multi-study cellular trajectories of development and disease

Authors: Brian Hie¹, Hyunghoon Cho², Bryan Bryson³, and Bonnie Berger^{1,4*}

Author information: ¹Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139; ²Broad Institute of MIT and Harvard, Cambridge, MA 02142; ³Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139; ⁴Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA 02139; *Correspondence: bab@mit.edu.

Abstract: Single-cell transcriptomic studies of diverse and complex systems are becoming ubiquitous. Algorithms now attempt to integrate patterns across these studies by removing *all* study-specific information, without distinguishing unwanted technical bias from relevant biological variation. Integration remains difficult when capturing biological variation that is distributed *across* studies, as when combining disparate temporal snapshots into a panoramic, multi-study trajectory of cellular development. Here, we show that a fundamental analytic shift to gene *coexpression* within clusters of cells, rather than gene expression within individual cells, balances robustness to bias with preservation of meaningful inter-study differences. We leverage this insight in Trajectorama, an algorithm which we use to unify trajectories of neuronal development and hematopoiesis across studies that each profile separate developmental stages, a highly challenging task for existing methods. Trajectorama also reveals systems-level processes relevant to disease pathogenesis within the microglial response to myelin injury. Trajectorama benefits from efficiency and scalability, processing nearly one million cells in around an hour.

1 **Introduction**

2 Single-cell RNA-sequencing (scRNA-seq) studies now profile millions of transcriptomes
3 across diverse tissues, conditions, species, and ages¹⁻⁹. To enable integration of biological
4 patterns into multi-study insight, several algorithms have been developed to align common cell
5 types across studies and then transform the underlying data to remove any study-specific
6 differences¹⁰⁻¹⁷; cells deemed to be of the same cell type will thus have similar transcriptomic
7 signatures in downstream analysis.

8 Unfortunately, because current integrative algorithms do not distinguish technical bias
9 from real biological variation, they remove any meaningful change in a cell type across
10 experimental conditions. A major task within single-cell analysis, however, is to infer trajectories
11 and “pseudo-temporal” relationships among cells, thereby algorithmically reconstructing
12 important continuous processes like differentiation or disease progression¹⁸⁻²¹. Reconstructing
13 such trajectories across disparate studies, separated by both experimental bias and real cellular
14 change, remains difficult even with state-of-the-art integration. Single-cell trajectories, therefore,
15 remain practically limited to patterns observed within a single study.

16 Here, we unite both integration and trajectory inference, two major single-cell analytic
17 efforts that have largely remained separate because current algorithms fail to achieve a delicate
18 balance between robustness to unwanted bias and preservation of relevant multi-study variation
19 (**Figure 1a**). To reveal dynamic biological processes at an unprecedented scope, we aim to
20 construct *multi-study trajectories* of cellular change.

21 Our novel, key insight is that differences in coexpression could preserve enough
22 biological variation while still enabling integration. Coexpression is a conceptually favorable
23 paradigm for integration since it favors redundant signal consistent across many genes^{10,25-27} and

24 since common coexpression measures (e.g., Spearman correlation) are robust to many
25 transformations of the data resulting from technical bias. In previous studies, coexpression has
26 been used extensively to assess global gene expression changes in different biological conditions
27 using both single-cell and bulk transcriptomics^{22–24}; here, we show that analysis that respects
28 variation in coexpression, combined with coexpression’s integrative properties, achieves a
29 balance crucial to enabling multi-study trajectory inference.

30 We therefore introduce Trajectoryrama, a coexpression-based algorithm for integration that
31 preserves and highlights cellular change across studies. Using Trajectoryrama, we efficiently
32 integrate trajectories of neuronal development (across embryonic, neonatal, adolescent, and adult
33 neurons) and hematopoiesis (across bone marrow, cord blood, fetal thymus, and peripheral
34 blood) that no other integrative method is able to recover. Trajectoryrama’s coexpression feature
35 space is highly interpretable, allowing us to probe the poorly understood microglial response to
36 myelin injury, revealing a disease-associated gene network across demyelination models in mice
37 and multiple sclerosis in human patients that implicates contributors to neurodegeneration.

38 Our conceptual advances beyond multi-study coexpression include panresolution
39 clustering, in which we consider all clusters across a cellular hierarchy for downstream analysis,
40 and interpretation through dictionary learning and functional analysis of condition-specific
41 coexpression networks. Our algorithmic innovations and versatile applications—from
42 understanding development across an entire lifespan to probing cell state change in response to
43 disease—underscore the utility of coexpression-based trajectory integration.

44 Results

45 *Multi-study coexpression analysis: Key concepts*

46 In conventional single-cell transcriptomic analysis, the fundamental analytic unit is an
47 individual cell described by features that encode levels of gene expression. A crucial difference
48 in Trajectoryrama's coexpression-based analysis is that the fundamental analytic unit is a *cluster* of
49 cells; this cluster is in turn described by features that encode the *correlation* in expression
50 between *pairs* of genes.

51 First, therefore, we require cells to be assigned to clusters. Clusters can be determined
52 based on experimentally-determined properties or conditions, or such clusters can be determined
53 by algorithms that group cells based on relative similarity in an unsupervised fashion²⁸. While
54 many clustering algorithms partition the data such that each cell is assigned to a single cluster,
55 this need not be the case. Indeed, cells often belong to a hierarchy of biologically-meaningful
56 groups²⁰; for example, in brain tissue, it may be useful to separate neurons and glia, but within
57 each category are many neuronal or glial subtypes. Rather than cluster cells based on a single
58 level of a cellular hierarchy, i.e., a single clustering "resolution," it is also possible to consider *all*
59 clusters at *multiple* resolutions. This approach is particularly useful when determining clusters
60 for coexpression-based analysis, since coexpression may change with clustering resolution^{24,27}.
61 We refer to this strategy as panresolution clustering, or *panclustering*.

62 After we determine clusters, each cluster is considered as a single datapoint in subsequent
63 analysis. The features that describe a cluster are the correlations in expression (within that
64 cluster) between all pairs of genes (**Figure 1b**). If there are M genes, then there will be $\binom{M}{2} + M$
65 unique gene pairs, where we compute a correlation for each pair. In Trajectoryrama, we use the
66 Spearman rank correlation due to its invariance under monotonic transformations of the

67 underlying data and robustness to small numbers of large-magnitude outliers. Equivalently, we
68 can think of each cluster as being described by a single gene-by-gene correlation matrix.
69 Equivalently, we can also think of each cluster as being described by a different gene association
70 network, where the weights of edges connecting genes correspond to correlation strength. We
71 can impose additional quality control cutoffs by setting low correlations to zero, or “sparsifying”
72 the features, which helps reduce noise and improve computational efficiency, a property we take
73 advantage of in our analysis.

74 Once we have featurized our clusters by coexpression, we can perform downstream
75 analyses, many of which are analogous to standard expression-based analyses. For example, we
76 can form trajectories by constructing a k-nearest-neighbors (KNN) graph where each node is a
77 cluster and edges between nodes are added based on proximity in coexpression feature space.
78 We can also find similarities and differences in coexpression among clusters, which correspond
79 to stable or changing gene-gene associations. Correlations unique to a condition can in turn be
80 interpreted as edges in a condition-specific gene network.

81 Trajectorama leverages and implements all of these concepts, encompassing cell
82 clustering through coexpression featurization through downstream interpretation, within a single
83 analytic framework, illustrated in **Figure 1b**. In particular, we design Trajectorama to integrate
84 vast amounts of data while preserving relevant study-specific biological variation.

85 *Unified trajectory of neuronal development containing 932,301 cells*

86 We first assessed whether coexpression could achieve the difficult balance of preserving
87 continuously changing cellular phenotypes while overcoming study-specific bias. Given a wealth
88 of scRNA-seq datasets that profile the mouse brain at different developmental timepoints, we
89 reasoned that coexpression could construct a picture of neuronal development at an

90 unprecedented scale. Known developmental age would help us validate the structure found by
91 our analysis.

92 We therefore used Trajectoryrama to analyze five large-scale studies of mouse neurons
93 from embryonic to adult. The first study¹ used sci-RNA-seq3 to profile 562,272 cells
94 representing the neural tube and notochord collected at day-length intervals from embryonic day
95 (E)9.5 through E13.5. The second³ used Drop-seq and 10x Chromium v2 to profile 50,363
96 cortical neurons from late embryonic (E13.5 - E14.5) and postnatal day (P)10. The third² used
97 Microwell Seq to profile 10,796 cells across three developmental timepoints for E14.5, P1, and
98 P56. The fourth⁴ used 10x Chromium v1 to profile 101,213 neurons from multiple adolescent
99 timepoints from P12 through P27 and from a P60 adult. The fifth⁵ used Drop-seq to profile
100 207,657 neurons from P60 through P70 adults. This data was generated by laboratories spanning
101 both United States coasts and three continents using single-cell or single-nucleus transcriptomic
102 platforms and in total profiled more than 150 individual mice.

103 We obtained a panclustering of cells based on the Louvain community detection
104 algorithm²⁹, a common clustering method for scRNA-seq data. Louvain clustering iteratively
105 merges cells into cluster “communities” until convergence, which is controlled by a resolution
106 parameter³⁰ (higher resolutions tend to increase the number of communities). We also obtain
107 many possible realizations of a Louvain clustering by repeating the algorithm with multiple
108 resolution parameters and use cluster assignments across all agglomerative iterations (**Methods**).
109 To see if coexpression could directly overcome study-specific bias, we panclustered each study
110 separately before combining clusters across all studies during downstream analysis in
111 coexpression space.

112 When we visualize the coexpression landscape with a force-directed embedding³¹ of the
113 KNN graph in which each node is a panresolution cluster, the graphical topology naturally
114 arranges according to biological age (**Figure 2a**) rather than study-specific structure (**Figure 2b**).
115 Analogous to assigning pseudotimes to cells in gene expression space, we can likewise run a
116 diffusion-based pseudotime (DPT) algorithm¹⁹ within the coexpression landscape using the
117 cluster with the lowest average age as the root of the diffusion process. Pseudotimes assigned to
118 panresolution clusters in coexpression space were significantly correlated with biological age
119 (Spearman $r = 0.87$, $P < 10^{-308}$, $n = 2,442$ panresolution clusters) (**Figure 2c**).

120 If instead we use gene expression to learn two-dimensional visualizations of these
121 datasets by plotting panresolution clusters using average gene expression, the datapoints arrange
122 according to study-of-origin, without conveying any continuous developmental structure (**Figure**
123 **2d,e**). Uniform Manifold Approximation and Projection (UMAP) visualization of cells, the key
124 algorithm underlying the Monocle 3 trajectory inference algorithm¹, also does not convey the
125 developmental relationships among the studies (**Supplementary Fig. 1**). Study-specific structure
126 is still present after applying existing integrative algorithms based on mutual nearest neighbors
127 matching³² (Scanorama) or on a latent space parameterized by a variational autoencoder (scVI)¹⁶
128 (**Fig. 2d,e**); these methods are representative of many others also based on nearest neighbors
129 matching¹¹⁻¹³ or on learning a joint latent space^{10,15,17}. Another integrative method, Harmony¹⁴,
130 removes nearly all study-specific signal, as designed (**Figure 2d,e**), which includes the valuable
131 development-related information that only the coexpression landscape captures.

132 *Interpretation of coexpression landscape yields insight into neuronal development*

133 Given this panoramic view into neuronal development, we facilitate further interpretation
134 by highlighting similar coexpression patterns across many panresolution clusters with *dictionary*

135 *learning*. In dictionary learning, we represent the coexpression matrix of each panresolution
136 cluster as a sparse weighted sum of a few basis coexpression matrices, or “dictionary entries.”
137 Each basis matrix can also be interpreted as a network, with edges between genes weighted by
138 coexpression. Dictionary learning for correlation matrices has been successfully applied to
139 diverse problems, including information retrieval³³ and functional brain profiling³⁴.

140 We looked for significant gene ontology (GO) process enrichments³⁵ within the set of
141 genes involved in “marker edges” unique to a particular dictionary entry, using a background set
142 of all genes considered in our coexpression analysis (around two thousand highly variable genes;
143 **Methods**). Within the embryonic portion of the coexpression landscape, we observe
144 differentiation and developmental processes (GO:0051094, false discovery rate [FDR] $q = 3.3 \times$
145 10^{-3}) and neuron fate commitment (GO:0048663, FDR $q = 8.4 \times 10^{-3}$). Late-fetal and early-
146 postnatal development includes neurogenesis (GO:0050767, FDR $q = 3.9 \times 10^{-4}$) and neuron
147 projection organization (GO:0030030, FDR $q = 0.018$). Adolescent and adult stages are enriched
148 for a more diverse set of processes from neurotransmission (GO:0001505, FDR $q = 1.5 \times 10^{-4}$) to
149 amyloid- β response (GO:1904646, FDR $q = 0.042$). The enriched processes for all of these
150 dictionary entries are consistent with their respective developmental stages, offering evidence
151 that Trajectory integration preserves inter-study patterns due to biological development.

152 We can also look at individual genes that are strongly associated with diffusion
153 pseudotime in the coexpression landscape and validate them with the Allen Developing Mouse
154 Brain Atlas (ADMBA), which spatially locates the expression of around 2000 preselected genes
155 using in situ hybridization (ISH) experiments³⁶. Genes with the strongest associations with
156 developmental pseudotime also showed strong developmental changes in ISH intensity in the
157 expected direction, i.e., increasing or decreasing with development. The top such positively

158 correlated gene is *Fos* (Spearman $r = 0.67$; $n = 2,442$ panresolution clusters), which encodes a
159 well-known marker of neuronal activity³⁷; the top such negatively correlated gene is *Eomes*
160 (Spearman $r = -0.45$; $n = 2,442$ panresolution clusters), which encodes an important transcription
161 factor in early neurogenesis³⁸ (**Figure 3b,c**).

162 Our analysis also reveals genes strongly associated with development, such as *Gm9945*
163 and *Pon3* (Spearman $r = 0.78$ and $r = -0.54$, respectively; $n = 2,442$ panresolution clusters), that
164 the ADMBA did not include in their list of assayed genes but may be important to include in
165 future developmental studies. We make these correlations and GO enrichments available as
166 **Supplementary Data**, which may be of further interest to developmental biologists.

167 *Neuronal developmental landscape is robust to parameter choice*

168 Two important parameters control the amount of information considered in our analysis
169 and can be thought of as “smoothing” parameters. The first is the correlation cutoff parameter
170 that controls the sparsity of underlying correlation matrices; lower values include more
171 information but may increase noise and computational burden. The second is the number of
172 nearest neighbors in the KNN graph representing the coexpression landscape, which impacts
173 both visualization and diffusion pseudotime; considering more nearest neighbors results in a
174 smoother trajectory. While we do introduce some smoothing into our analysis, the studies are
175 consistently arranged according to their developmental order even as these parameters vary.
176 With less smoothing, we also observe age-related branching of the developmental trajectory,
177 suggestive of neuronal subtype-related structure (**Supplementary Fig. 2**).

178 *Coexpression integrates neuronal subtypes across studies*

179 While the most pronounced signal captured within the neuronal trajectory is
180 developmental age, there is still substantial heterogeneity among neurons. We therefore sought to

181 determine if Trajectoryrama could provide multi-study insight into neuronal *subtypes* as well. To
182 do so, we relied on extensive expert labelling of neuronal subtypes from Zeisel *et al.*⁴ (adolescent
183 mice) and Saunders *et al.*⁵ (adult mice) to define neuronal clusters of interest. When comparing
184 these subtypes in gene expression space, subtypes group primarily according to study (**Figure**
185 **3d**). When we instead featurize by coexpression, the clusters group primarily according to
186 common subtypes, and only secondarily (since we do expect some differences due to real
187 biological change) according to study (**Figure 3d**).

188 Neuronal subtypes group according to three major coexpression-based patterns. Genes
189 most unique to the first group are enriched in glutamergic structures (GO:0098978, FDR $q = 1.5$
190 $\times 10^{-11}$) and glutamate signaling (GO:0035235, FDR $q = 6.4 \times 10^{-3}$). In contrast, the second
191 group has significant enrichments for both adrenergic (GO:0004935, FDR = 0.017) and
192 cholinergic (GO:0032224, FDR $q = 0.037$) processes. The third group, which also contains the
193 highest number of adolescent subtypes, is most significantly enriched for neurons with synaptic
194 plasticity (GO:0048167, FDR $q = 9.3 \times 10^{-8}$) and involved in cognition (GO:0007611, FDR $q =$
195 1.9×10^{-4}), learning, and memory (GO:0007611, FDR $q = 7.4 \times 10^{-5}$). The hierarchy of subtypes
196 has additional structure as well, though we focused on only the three largest, highest-level
197 groupings that each contain subtypes from both studies (**Figure 3d**).

198 *Trajectoryrama constructs a multi-tissue hematopoietic trajectory*

199 Based on the ability of Trajectoryrama to integrate neuronal studies while respecting
200 biological change, we next set out to establish if it could demonstrate similar capabilities within a
201 completely separate developmental system. To this end, we analyzed the coexpression landscape
202 of four hematopoietic datasets from the fetal thymus³⁹, bone marrow, cord blood⁷, and peripheral
203 blood⁶. Throughout these tissues, we expect to observe cells in many stages of hematopoiesis,

204 including stem cells and erythroid progenitors, mostly in the bone marrow and cord blood, to
205 more mature lymphocytes and myeloid cells, mostly as peripheral blood mononuclear cells
206 (PBMCs)⁴⁰.

207 Visualizing the coexpression landscape of panresolution clusters obtained across all
208 studies reveals an organization consistent with the three main branches of hematopoiesis
209 corresponding to erythropoiesis, myelopoiesis, and lymphopoiesis (**Figure 4a**). Such
210 organization (with similar developmental granularity) has been observed in the gene expression
211 space²⁰ and in the chromatin accessibility space⁴¹ of single studies in single tissues, but,
212 importantly, here we instead show a unified hematopoietic landscape across multiple tissues
213 generated by disparate laboratories.

214 We interpret different branches in the coexpression trajectory partially based on
215 experimentally-determined PBMC labels. Prior to scRNA-seq, a large number of the PBMCs
216 underwent fluorescence activated cell sorting (FACS) for progenitor-associated (CD34⁺),
217 myeloid-associated (CD14⁺), and lymphoid-associated (CD4⁺, CD8⁺, CD19⁺, CD56⁺) cell-
218 surface marker expression (**Supplementary Fig. 3**). Dictionary learning yielded four main
219 dictionary entries corresponding to the major regions within the landscape (**Figure 4b**). The first
220 dictionary entry, which we call progenitor-associated, corresponds to all of the CD34⁺-labeled
221 clusters. The second dictionary entry, which we call erythropoietic, includes GO enrichments
222 related to heme biosynthesis (GO:0006783, FDR $q = 0.04$) and strong metabolic signatures
223 (GO:0044237, FDR $q = 1.0 \times 10^{-12}$). The third dictionary entry, which we call lymphopoietic,
224 includes all lymphoid-specific (CD4⁺, CD8⁺, CD19⁺, CD56⁺) clusters. The fourth dictionary
225 entry, which we call myelopoietic, includes some CD14⁺ clusters and GO enrichments involving
226 myeloid differentiation (GO:0045637, FDR $q = 3.4 \times 10^{-4}$).

227 We also note that PBMCs largely exist at the periphery of the landscape, consistent with
228 such cells being the most mature within the hematopoietic lineage. In contrast, Harmony-based
229 integration removes all tissue-specific differences and obscures the lineage relationships among
230 the tissues (**Figure 4c**) while mean expression of clusters without correction, and even following
231 Scanorama and scVI correction, primarily exhibits study-specific structure (**Figure 4c**). Overall,
232 our hematopoietic analysis adds additional support for coexpression as an integrative strategy
233 that can preserve key biological differences among disparate studies.

234 *Trajectoryrama reveals a disease-specific microglial gene network*

235 While Trajectoryrama can yield panoramic views across long developmental scales, we
236 next wanted to assess if it could also reveal more fine-grained insight into biological systems that
237 are less well understood. In particular, recent work has begun to illuminate the key role of
238 microglia in neurodegenerative disease^{42,43}, for which coexpression provides a unique
239 opportunity to integrate information across multiple microglial studies while still preserving
240 disease-specific signal.

241 We therefore integrated microglia from mouse and human samples across three
242 studies^{5,8,9}, which together contained single-cell microglial transcriptomes from multiple points
243 along a mouse lifespan, from models of mouse brain injury (facial nerve axotomy and
244 demyelination), and from human donors with and without multiple sclerosis (MS). The
245 Trajectoryrama coexpression landscape includes a main age-related trajectory, from embryonic
246 (E14.5) through aged (P540) microglia, and off-trajectory outlier clusters from injured tissue
247 samples (**Figure 5a**); similarly, hierarchically grouping the known microglial conditions based
248 on similarity in coexpression space (**Methods**) obtains a clear outlier group consisting of
249 microglia from mice that had undergone artificial demyelination and from human MS patients

250 **(Figure 5b)**. We note that, in coexpression space, this injury-associated group naturally separates
251 from other microglial conditions without supervision.

252 We then constructed an injury-associated coexpression network by considering the gene
253 pairs with the highest increase in coexpression, combined across both mouse and human injury
254 conditions, relative to baseline microglial coexpression (**Methods**). GO enrichment analysis of
255 genes ranked by increased coexpression in disease state reveals three main functional categories:
256 lipid and protein clearance, leukocyte-mediated cytotoxicity, and cellular activation involved in
257 inflammation (**Figure 5c; Supplementary Data**). These processes are consistent with the
258 hypothesized role of microglia in MS as involved in clearance of damaged myelin via
259 phagocytosis⁴² and as drivers of neurodegenerative pathogenesis by inducing neuronal cell
260 death⁴⁴ and promoting local inflammation⁴².

261 The most valuable insight into microglial processes relevant to disease, and to myelin
262 injury in particular, comes from visualizing the injury-associated coexpression network itself
263 (**Figure 5d**). Two major connected components appear in the network: the first related to lipid
264 clearance and leukocyte-mediated cytotoxicity and the second related to inflammatory activation.

265 The first connected component recovers key gene modules that have been implicated in
266 neurodegeneration. Of special note, the network recovers an *APOE/TREM2/GM2A* gene module
267 that has been extensively linked to a microglial “sensor” of neurodegeneration^{9,43,45}. Another
268 high-degree module includes *SIRPA*, which regulates demyelination repair⁴⁶, and *MSRI*, which
269 has been implicated in myelin uptake in MS lesions⁴⁷. The network suggests a correlative link
270 from the *APOE/TREM2* neurodegenerative sensing module to the *SIRPA/MSRI* uptake and
271 clearance module through genes like *AXL*, which has also been suggested as essential to recovery
272 from myelin injury⁴⁸. While many of these genes have been *individually* implicated in

273 neurodegeneration, we note that our coexpression-based analysis suggests *links* among these key
274 genes that are useful for follow up study. Experimentally establishing the causal role these genes
275 play in disease pathogenesis is important future work.

276 The second major connected component centers on the *NEATI* long noncoding
277 (lnc)RNA, which recently has been linked to inflammatory activation of macrophages⁴⁹. These
278 results suggest that the observation of *NEATI* in MS serum, for which the mechanism was
279 previously unknown⁵⁰, is tied in part to microglial inflammation. Further experimentation is
280 needed to see if *NEATI* leads to or is a consequence of inflammation-mediated pathogenesis, or
281 it could also serve as a biomarker of MS disease or related inflammation.

282 More broadly, the injury-associated microglia network illustrates how coexpression-
283 based analysis across multiple experiments can generate further hypotheses that lead to novel
284 biological discovery. Not only can coexpression analysis elucidate broad developmental changes,
285 its rich feature space and inherent interpretability can also provide deep insight into cell state
286 changes such as those in health versus disease.

287 *Trajectorama is practical for datasets with millions of cells*

288 To enable consortium-scale analysis, we made algorithmic choices that allow scalability
289 to large numbers of cells, while preserving the ability to model complex phenomena. For
290 example, we choose to sparsify our coexpression matrices using a nominal cutoff rather than the
291 memory intensive strategy of preserving dense correlation matrices or the runtime intensive
292 strategy of learning sparse covariance matrices via regularization⁵¹ (**Supplementary Table 1**).
293 Since scRNA-seq experiments typically measure little to no signal for many genes, we also
294 limited analysis to around two thousand genes with highest statistical variability, a common
295 dimensionality reduction strategy in conventional expression analysis^{28,52} (**Methods**).

296 We performed all of our analyses in a practical amount of computational time and
297 resources. Our entire coexpression-based procedure, which includes panresolution clustering
298 through downstream analysis of the coexpression landscape, analyzes almost a million cells in a
299 little over an hour on a standard cloud instance with 16 cores (**Supplementary Table 2**). Our
300 pipeline has a runtime and memory usage with a close-to-linear asymptotic scaling in the number
301 of cells and a worst-case quadratic asymptotic scaling in the number of features (i.e., genes).
302 While the coexpression space may seem cumbersome quadratic, scRNA-seq experiments
303 typically measure only around one or two thousand genes with nontrivial variability⁵²; moreover,
304 the number of strong correlations is usually within the same order of magnitude as the number of
305 highly variable genes.

306 Once the data has been summarized as panresolution clusters, further downstream
307 analysis including visualization, pseudotime assignment, and dictionary learning becomes
308 extremely efficient due to the greatly reduced number of datapoints; in the case of mouse
309 neuronal development, analysis is done on just 2,442 panresolution clusters instead of 932,301
310 single cells. The resource requirements for different stages of our analytic pipeline on the mouse
311 neuronal development analysis are provided in **Supplementary Table 2**.

312 **Discussion**

313 Our work shows that researchers can analyze an unprecedented amount of information
314 across scRNA-seq studies, while retaining key biological variation, by focusing on the
315 coexpression matrix of a group of cells as the fundamental unit of analysis. While not intended
316 as a complete replacement for current integrative methods, as we have shown, Trajectorama can
317 be valuable when researchers wish to integrate data while preserving inter-study biological
318 variation. As laboratories continue to conduct single-cell experiments that explore heterogeneous
319 biological models and conditions, we expect such scenarios to be ubiquitous.

320 By leveraging coexpression, Trajectorama benefits from a number of additional
321 properties. Current integrative methods map cells into an arbitrary feature space that only
322 preserves *relative* meaning (for example, cell A is more similar to cell B than to cell C). In
323 contrast, coexpression has *intrinsic* meaning: each feature in coexpression space is simply the
324 correlation between two genes (for example, Spearman correlation⁵³), a fundamental and
325 intuitive data science concept. Trajectorama is also highly efficient, since it combines
326 information across many cells similar to existing algorithms that accelerate workflows via data
327 sketching or summarization^{54,55}.

328 Our results suggest many directions for future work. Our coexpression matrices are not
329 positive semidefinite (PSD) for practical reasons, but efficiently learning large numbers of
330 nontrivially sparse PSD matrices is an interesting and challenging task. If all coexpression
331 matrices are PSD, it may be possible to leverage the distance along the manifold represented by
332 all PSD matrices to obtain more natural dictionary learning-based decompositions³³ and nearest-
333 neighbor queries (which would involve designing new techniques for efficient nearest-neighbor
334 search). Additional methods might also enforce further constraints within the dictionary learning

335 objective (for example, basis matrices that are valid correlation matrices) or take other
336 approaches to interpreting large numbers of coexpression matrices like common principal
337 components analysis⁵⁶ or other kinds of tensor decomposition⁵⁷.

338 Other considerations include exploring alternative methods for measuring coexpression⁵⁸,
339 inferring causal gene regulatory networks, or exploring different clustering strategies,
340 panresolution or otherwise. A larger question is whether other feature spaces exist that enable
341 multi-study trajectories; for example, metric learning approaches could directly construct such a
342 space via known developmental metadata⁵⁹. Reasoning about the relationship between
343 coexpression and other functional associations within single cells, like those involving chromatin
344 accessibility or methylation, remains an important consideration.

345 Trajectorama can be used to probe biological systems beyond those interrogated in this
346 study, providing an informative analysis that is complementary to existing integrative methods
347 for studying biological processes at single-cell resolution and at multi-institution scale. We make
348 our analysis pipelines and data available at <http://trajectorama.csail.mit.edu>.

349 **Methods**

350 *Mouse neuronal development dataset preprocessing*

351 We obtained publicly available datasets from five large-scale, published single-cell
352 transcriptomic studies of the mouse brain at different developmental timepoints¹⁻⁵. We used only
353 the cells that passed the filtering steps of each respective study and additionally removed low-
354 complexity or quiescent cells with less than 500 unique genes. For the embryonic dataset from
355 Cao *et al.*¹, we only considered cells that the study authors had assigned to the “neural tube and
356 notochord” trajectory. For the datasets from Zeisel *et al.*⁴ and Saunders *et al.*⁵ we only
357 considered cells that the study authors had labeled as neuronal. We then intersected the genes
358 with the highest variance-to-mean ratio (i.e., dispersion) within each study to obtain a total of
359 around 2,000 genes that were highly variable across all studies. All studies provided data as
360 digital gene expression (DGE) counts, which we further log transform after adding a pseudo-
361 count of 1.

362 *Human hematopoiesis dataset preprocessing*

363 We obtained publicly available datasets of cord blood and bone marrow cells from the
364 Human Cell Atlas⁷ (<https://preview.data.humancellatlas.org/>) and PBMCs from Zheng *et al.*⁶
365 (<https://support.10xgenomics.com/single-cell-gene-expression/datasets>). We removed cells with
366 less than 500 unique genes; we also noticed a large number of cells with high percentages of
367 ribosomal transcripts, which may indicate nontrivial amounts of ambient ribosomal RNA
368 contamination during the scRNA-seq experiment, so we only included cells with less than 50%
369 ribosomal transcripts in further analysis. As described previously, we intersected the genes with
370 the highest dispersions within each study to obtain a total of around 2,000 genes that were highly

371 variable across all studies. All studies provided data as digital gene expression (DGE) counts,
372 which we further log transform after adding a pseudo-count of 1.

373 *Microglia dataset preprocessing*

374 We obtained publicly available datasets from three single-cell transcriptomic studies of
375 microglia across a diverse set of conditions^{5,8,9}. We kept only the cells labeled by the original
376 studies as microglia and we additionally removed low-complexity or quiescent cells with less
377 than 500 unique genes. Mouse genes were mapped to human orthologs. As described previously,
378 we intersected the genes with highest dispersions within each study to obtain around 2,000 genes
379 that were highly variable across studies, followed by a log transformation after adding a pseudo-
380 count of 1.

381 *Panresolution clustering*

382 We modify the Louvain clustering algorithm^{29,30} ([https://github.com/vtraag/louvain-](https://github.com/vtraag/louvain-igraph)
383 [igraph](https://github.com/vtraag/louvain-igraph)) to store community information at each iteration. We choose Louvain clustering due to
384 its asymptotic efficiency, since its runtime and space usage scales with the size of the k -nearest
385 neighbor (KNN) graph of cells (i.e., each cell is a node in the graph), rather than quadratically in
386 the number of cells as in other hierarchical clustering algorithms. To capture a range of potential
387 clustering results, we rerun the Louvain clustering algorithm at a diverse range of clustering
388 resolutions (0.1, 1, and 10) on a 15-nearest neighbor graph, constructed using Euclidean
389 distances in gene expression space, storing the hierarchical cluster information for each run. The
390 three runs of Louvain clustering are done in parallel and we cluster each study individually. To
391 reduce the effect of noisy correlations, we consider clusters with a minimum of 500 cells, which,
392 combined with highly variable gene filtering (described below), reduces the chance that a strong
393 correlation is due to a few outlier cells.

394 *Computing coexpression matrices*

395 We compute the Spearman correlation matrix $\mathbf{R}^{(i)} \in [-1,1]^{M \times M}$ for each of the
396 panresolution clusters obtained as described above, where $i \in \{1,2, \dots, N\}$ with N denoting the
397 number of panresolution clusters and M denoting the number of highly variable genes. The entry
398 $\mathbf{R}_{ab}^{(i)}$ at row a and column b of $\mathbf{R}^{(i)}$, corresponding to the a^{th} and b^{th} genes, takes the value

$$399 \quad \mathbf{R}_{ab}^{(i)} = \begin{cases} r_{ab}^{(i)} & \text{if } |r_{ab}^{(i)}| > \eta \text{ and } \sigma_a > 0 \text{ and } \sigma_b > 0 \\ 0 & \text{otherwise,} \end{cases}$$

400 where $r_{ab}^{(i)}$ is the Spearman correlation coefficient⁵³ and σ_a and σ_b are the respective standard
401 deviations of the rank values of the gene expressions (which appear in the denominator of the
402 Spearman correlation expression). $\eta \in [0, 1]$ is a sparsification parameter that sets low
403 correlations to zero and can be interpreted as a smoothing parameter that preserves only the most
404 important associations. Low values of this parameter can introduce additional structure into the
405 analysis, but may also introduce larger amounts of noise (see **Supplementary Fig. 2**).

406 *Visualization and diffusion pseudotime analysis of panresolution clusters*

407 To visualize the coexpression landscape defined by the panresolution clusters, the
408 symmetric correlation matrices $\mathbf{R}^{(i)} \in [-1, 1]^{M \times M}$ are treated as vectors $\mathbf{r}^{(i)} \in [-1, 1]^{\binom{M}{2}+M}$ on
409 which we construct the k -nearest neighbors graph using the Euclidean distance in coexpression
410 space as the distance metric. This graph was visualized with a force-directed embedding using
411 the ForceAtlas2 algorithm³¹ (<https://github.com/bhargavchippada/forceatlas2>). For the mouse
412 neuronal development analysis, a diffusion pseudotime (DPT) algorithm¹⁹ was applied to this
413 graph using the panresolution cluster with the earliest average age as the root. Larger values of k
414 can also increase the amount of smoothing in the structure captured by the k -nearest-neighbors
415 graph and subsequent visualization and DPT analysis (see **Supplementary Fig. 2**). We used the

416 implementation in Scanpy⁶⁰ (<https://scanpy.readthedocs.io/en/stable/>) for the k -nearest neighbors
417 graph construction and DPT analysis.

418 We also visualized panresolution clusters in gene expression space, Harmony-integrated
419 expression space¹⁴, Scanorama-corrected expression space³², and scVI-integrated latent space¹⁶.
420 To summarize features across multiple cells into a single feature vector for each panresolution
421 cluster, we used the mean expression. We similarly constructed the k -nearest-neighbors graph
422 with panresolution clusters as nodes and Euclidean distance between the summarized gene
423 expression values as the distance metric.

424 *Coexpression matrix dictionary learning*

425 We formulated the dictionary learning problem for coexpression matrices by optimizing

$$426 \operatorname{argmin}_{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(N)}, \mathbf{V}} \left\{ \sum_{i=1}^N \|\mathbf{r}^{(i)} - \mathbf{V}\mathbf{u}^{(i)}\|_2^2 + \alpha \|\mathbf{u}^{(i)}\|_1 \right\}$$

427 subject to $\|\mathbf{v}_j\|_2 = 1$ for all $j \in [\kappa]$

428 where $\mathbf{u}^{(i)} \in \mathbb{R}_{\geq 0}^{\kappa}$ is a sparse code of weights for panresolution cluster i , α is a sparsity-
429 controlling parameter, $\mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_j \cdots \mathbf{v}_{\kappa}] \in \mathbb{R}_{\geq 0}^{((M) + M) \times \kappa}$ is a dictionary of κ (vectorized)
430 coexpression matrices, and κ is a user-defined parameter indicating the number of dictionary
431 entries to learn. We used an iterative optimization algorithm that alternatively estimated
432 dictionary weights and dictionary entries using least angle regression-based optimization⁶¹ until
433 convergence. We tune κ by plotting the objective function error versus values of κ and manually
434 selecting a value after which there are relatively smaller drops in objective function values, a
435 parameter selection procedure often referred to as the “elbow method.”

436 *Interpretation of dictionary entries*

437 We can interpret each dictionary entry \mathbf{v}_j as a coexpression network in which genes are
438 nodes and elements of \mathbf{v}_j define edge weights between those genes. We use the networkx Python
439 package⁶² to represent graphs and compute various graph statistics. Using genes that are involved
440 in edges that are unique to a given coexpression network, we look for gene ontology (GO)
441 process enrichments using a background set of all highly variable genes considered in the
442 analysis, for which P -values can be computed using a hypergeometric null model followed by
443 subsequent FDR q -value computation⁶³. We use the GOrilla webtool ([http://cbl-](http://cbl-gorilla.cs.technion.ac.il/)
444 [gorilla.cs.technion.ac.il/](http://cbl-gorilla.cs.technion.ac.il/))³⁵ with default parameters, which reports all enrichments more
445 significant than a nominal P -value of $1e-3$. We use the REVIGO webtool (<http://revigo.irb.hr/>)
446 with default parameters, which consolidates similar GO terms and visualizes terms in a two-
447 dimensional “semantic space” that places similar terms closer together⁶⁴. We only consider
448 dictionary entries that have nonzero weights in at least ten panresolution clusters.

449 *Neuronal subtype hierarchical grouping and interpretation*

450 Cellular subtypes were determined according to expert curated labels provided by the
451 original studies^{4,5}. Each subtype was featurized by coexpression and by mean expression for
452 benchmarking purposes. Agglomerative hierarchical clustering of the subtypes was then
453 performed using the scipy Python library⁶⁵. To interpret genes unique to a group of subtypes, we
454 computed the mean coexpression within the group and sorted each dimension according to the
455 highest increase in correlation from the mean coexpression of all subtypes. Genes were then
456 ranked according to the first appearance of the gene within the sorted list of coexpression
457 dimensions; this gene ranking was used as input into the GOrilla webtool for GO enrichment
458 analysis.

459 *Microglial subtype analysis and interpretation*

460 Microglial subtypes were determined based on unique combinations of age, species, and
461 tissue injury status. For the coexpression landscape analysis, each of these subtypes was
462 considered as a separate study. Fewer clusters enabled a lower sparsification threshold of $\eta =$
463 0.1. All other methods and parameters remained the same.

464 As in the neuronal subtype analysis, we also hierarchically clustered the microglial
465 subtypes and observed an injury-associated group of microglial subtypes. We took the
466 coexpression mean of this injury-associated microglial group, including both mouse and human
467 clusters, and compared it to the mean of all microglial subtypes. Coexpression dimensions were
468 sorted according to the highest increase in correlation within the injury-associated group. This
469 sorted list was used to rank genes as input into the GOrilla webtool for GO enrichment analysis
470 and the first 150 edges in this list (all with an increase in correlation greater than 0.26) was used
471 to visualize the disease-specific coexpression network. We used Gephi version 0.9.2
472 (<https://gephi.org/>) to visualize the network⁶⁶.

473 *Statistical analysis and implementation*

474 We use the scientific Python toolkit, including the scipy and numpy Python packages⁶⁵,
475 to compute the statistical tests described in the manuscript, including Spearman correlation and
476 associate P -values. P -values listed as less than 10^{-308} indicate values returned by the statistical
477 software below the minimum nonzero floating-point value representable by the machine.

478 *Runtime and memory profiling*

479 We used Python's time module to obtain runtime measurements and used the top
480 program in Linux (Ubuntu 17.04) to make periodic memory measurements. We made use of
481 default scientific Python parallelism. We benchmarked our pipelines on a Google Cloud

482 Enterprise instance with 16 logical cores and 104 gigabytes of memory and, for memory-
483 inefficient alternative algorithms (**Supplementary Table 1**), on a local 2.30 GHz Intel Xeon E5-
484 2650v3 with 48 logical cores and 384 GB of RAM. scVI was trained on a Nvidia Tesla V100-
485 SXM2 with 16 GB of RAM.

Data Availability

We used the following publicly available datasets:

- Notochord and neural plate cells from Cao *et al.*¹ (GSE119945)
- Neurons from Mayer *et al.*² (GSE104158)
- Neurons from Han *et al.*³ (https://figshare.com/articles/MCA_DGE_Data/5435866)
- Neurons from Zeisel *et al.*⁴ (<http://mousebrain.org/>)
- Neurons and microglia from Saunders *et al.*⁵ (GSE116470)
- In-situ hybridization images from the Allen Developing Mouse Brain Atlas³⁶
(<https://developingmouse.brain-map.org/>)
- Bone marrow and cord blood cells from the Human Cell Atlas
(<https://preview.data.humancellatlas.org/>)
- PBMCs from Zheng *et al.*⁶ (<https://support.10xgenomics.com/single-cell-gene-expression/datasets>)
- Fetal thymus hematopoietic cells from Zeng *et al.*³⁹ (GSE133341)
- Microglia from Hammond *et al.*⁸ (GSE121654)
- Microglia from Masuda *et al.*⁹ (GSE124335)

Acknowledgements

We thank R. Chun, B. DeMeo, C. Mak, S. Nyquist, C. Wong-Fillman, and the Berger and Bryson laboratory members for valuable discussions and feedback. B.H. is partially supported by NIH grant R01 GM081871 (to B. Berger) and by the Department of Defense (DoD) through the National Defense Science and Engineering Graduate Fellowship (NDSEG)

Author Contributions

All authors conceived the algorithm. B.H. implemented the algorithm and performed the computational experiments. All authors interpreted the results and wrote the manuscript.

References

1. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**, 496–502 (2019).
2. Mayer, C. *et al.* Developmental diversification of cortical inhibitory interneurons. *Nature* **555**, 457–462 (2018).
3. Han, X. *et al.* Mapping the Mouse Cell Atlas by Microwell-Seq. *Cell* **175**, P1091-1107.e17 (2018).
4. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell* **174**, 999-1014.e22 (2018).
5. Saunders, A. *et al.* Molecular Diversity and Specializations among the Cells of the Adult Mouse Brain. *Cell* **174**, 1015-1030.e16 (2018).
6. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, (2017).
7. Li, B. & Regev, A. HCA data portal - census of immune cells.
8. Hammond, T. R. *et al.* Single-Cell RNA Sequencing of Microglia throughout the Mouse Lifespan and in the Injured Brain Reveals Complex Cell-State Changes. *Immunity* **50**, 253–271 (2019).
9. Masuda, T. *et al.* Spatial and temporal heterogeneity of mouse and human microglia at single-cell resolution. *Nature* **566**, 388–392 (2019).
10. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* **36**, 411–420 (2018).
11. Haghverdi, L., Lun, A., Morgan, M. & Marioni, J. Batch effects in single-cell RNA-

- sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
12. Stuart, T. *et al.* Comprehensive integration of single cell data. *Cell* **177**, 1888-1902.E21 (2019).
 13. Barkas, N. *et al.* Joint analysis of heterogeneous single-cell RNA-seq dataset collections. *Nat. Methods* **16**, 695–698 (2019).
 14. Korsunsky, I. *et al.* Fast, sensitive, and accurate integration of single cell data with Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
 15. Welch, J. *et al.* Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell* **177**, 1873-1887.E17 (2019).
 16. Xu, C. *et al.* Harmonization and Annotation of Single-cell Transcriptomics data with Deep Generative Models. *bioRxiv* (2019). doi:10.1101/532895
 17. Lotfollahi, M., Wolf, F. A. & Theis, F. J. scGen predicts single-cell perturbation responses. *Nat. Methods* **16**, 715–721 (2019).
 18. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
 19. Haghverdi, L., Büttner, M., Wolf, F. A., Buettner, F. & Theis, F. J. Diffusion pseudotime robustly reconstructs lineage branching. *Nat. Methods* **13**, 845–848 (2016).
 20. Wolf, F. A. *et al.* PAGA: Graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome Biol.* **29**, Article number: 59 (2019).
 21. Saelens, W., Cannoodt, R., Helena, T. & Saeys, Y. A comparison of single-cell trajectory inference methods. *Nat. Biotechnol.* **37**, 547–554 (2019).

22. Choi, J. K., Yu, U., Yoo, O. J. & Kim, S. Differential coexpression analysis using microarray data and its application to human cancer. *Bioinformatics* **21**, 4348–4355 (2005).
23. de la Fuente, A. From ‘differential expression’ to ‘differential networking’ - identification of dysfunctional regulatory networks in diseases. *Trends Genet.* **26**, 326–333 (2010).
24. Feigelman, J., Theis, F. J. & Marr, C. MCA: Multiresolution Correlation Analysis, a graphical tool for subpopulation identification in single-cell gene expression data. *BMC Bioinformatics* **15**, Article number: 240 (2014).
25. Crow, M., Paul, A., Ballouz, S., Huang, Z. J. & Gillis, J. Exploiting single-cell expression to characterize co-expression replicability. *Genome Biol.* **17**, Article number: 101 (2016) (2016).
26. Crow, M., Paul, A., Ballouz, S., Huang, Z. J. & Gillis, J. Characterizing the replicability of cell types defined by single cell RNA-sequencing data using MetaNeighbor. *Nat. Commun.* **9**, Article number: 884 (2018).
27. Crow, M. & Gillis, J. Co-expression in Single-Cell Analysis: Saving Grace or Original Sin? *Trends Genet.* **34**, 823–831 (2018).
28. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol. Syst. Biol.* **15**, e8746 (2019).
29. Blondel, V. D., Guillaume, J. L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech. Theory Exp.* (2008). doi:10.1088/1742-5468/2008/10/P10008
30. Lambiotte, R., Delvenne, J. C. & Barahona, M. Random walks, Markov processes and the multiscale modular organization of complex networks. *IEEE Trans. Netw. Sci. Eng.* **1**, 76–

- 90 (2014).
31. Jacomy, M., Venturini, T., Heymann, S. & Bastian, M. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PLoS One* **9**, Article number: 6 (2014).
 32. Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nat. Biotechnol.* **37**, 685–691 (2019).
 33. Cherian, A. & Sra, S. Riemannian Dictionary Learning and Sparse Coding for Positive Definite Matrices. *IEEE Trans. Neural Networks Learn. Syst.* **28**, 2859–2871 (2017).
 34. Eavani, H., Satterthwaite, T., Gur, R., Gur, R. & Davatzikos, C. Unsupervised Learning of Functional Network Dynamics in Resting State fMRI. *Inf. Process. Med. Imaging* 426–437 (2013).
 35. Eden, E., Navon, R., Steinfeld, I., Lipson, D. & Yakhini, Z. GOrilla: A tool for discovery and visualization of enriched GO terms in ranked gene lists. *BMC Bioinformatics* **10**, Article number: 48 (2009).
 36. Thompson, C. L. *et al.* A high-resolution spatiotemporal atlas of gene expression of the developing mouse brain. *Neuron* **83**, 309–323 (2014).
 37. Chung, L. A Brief Introduction to the Transduction of Neural Activity into Fos Signal. *Dev. Reprod.* **19**, 61–67 (2015).
 38. Arnold, S. J. *et al.* The T-box transcription factor Eomes/Tbr2 regulates neurogenesis in the cortical subventricular zone. *Genes Dev.* **22**, 2479–2484 (2008).
 39. Zeng, Y. *et al.* Single-Cell RNA Sequencing Resolves Spatiotemporal Development of Pre-thymic Lymphoid Progenitors and Thymus Organogenesis in Human Embryos. *Immunity* **51**, 930–948 (2019).

40. Zhang, Y., Gao, S., Xia, J. & Liu, F. Hematopoietic Hierarchy – An Updated Roadmap. *Trends Cell Biol.* **28**, 976–986 (2018).
41. Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548.e16 (2018).
42. Luo, C. *et al.* The role of microglia in multiple sclerosis. *Neuropsychiatr. Dis. Treat.* **13**, 1661–1667 (2017).
43. Deczkowska, A. *et al.* Disease-Associated Microglia: A Universal Immune Sensor of Neurodegeneration. *Cell* **173**, 1073–1081 (2018).
44. Dhib-Jalbut, S. & Kalvakolanu, D. V. Microglia and necroptosis: The culprits of neuronal cell death in multiple sclerosis. *Cytokine* **76**, 583–594 (2015).
45. Krasemann, S. *et al.* The TREM2-APOE Pathway Drives the Transcriptional Phenotype of Dysfunctional Microglia in Neurodegenerative Diseases. *Immunity* **47**, 566–581 (2017).
46. Sato-Hashimoto, M. *et al.* Microglial SIRPα regulates the emergence of CD11c⁺ microglia and demyelination damage in white matter. *Elife* **8**, e42025 (2019).
47. Hendrickx, D. A. E. *et al.* Gene expression profiling of multiple sclerosis pathology identifies early patterns of demyelination surrounding chronic active lesions. *Front. Immunol.* **8**, 1810 (2017).
48. Weinger, J. G. *et al.* Loss of the receptor tyrosine kinase Axl leads to enhanced inflammation in the CNS and delayed removal of myelin debris during Experimental Autoimmune Encephalomyelitis. *J. Neuroinflammation* **8**, 49 (2011).
49. Zhang, P., Cao, L., Zhou, R., Yang, X. & Wu, M. The lncRNA Neat1 promotes activation of inflammasomes in macrophages. *Nat. Commun.* **10**, Article number: 1495 (2019).
50. Santoro, M. *et al.* Expression Profile of Long Non-Coding RNAs in Serum of Patients

- with Multiple Sclerosis. *J. Mol. Neurosci.* **59**, 18–23 (2016).
51. Friedman, J., Hastie, T. & Tibshirani, R. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**, 432–441 (2008).
 52. Yip, S. H., Sham, P. C. & Wang, J. Evaluation of tools for highly variable gene discovery from single-cell RNA-seq data. *Brief. Bioinform.* **20**, 1583–1589 (2019).
 53. Zwillinger, D. & Kokoska, S. *CRC Standard Probability and Statistics Tables and Formulae. CRC Standard Probability and Statistics Tables and Formulae* (1999).
doi:10.1201/9781420050264
 54. Hie, B., Cho, H., DeMeo, B., Bryson, B. & Berger, B. Geometric Sketching Compactly Summarizes the Single-Cell Transcriptomic Landscape. *Cell Syst.* **8**, 483-493.E7 (2019).
 55. Baran, Y. *et al.* MetaCell: analysis of single cell RNA-seq data using k-NN graph partitions. *Genome Biol.* **20**, Article number: 206 (2019).
 56. Trendafilov, N. T. Stepwise estimation of common principal components. *Comput. Stat. Data Anal.* **54**, 3446–3457 (2010).
 57. Bergqvist, G. & Larsson, E. The higher-order singular value decomposition: Theory and an application. *IEEE Signal Process. Mag.* **27**, 151–154 (2010).
 58. Skinnider, M. A., Squair, J. W. & Foster, L. J. Evaluating measures of association for single-cell transcriptomics. *Nat. Methods* **16**, 381–386 (2019).
 59. Singh, R., Narayan, A., Hie, B. & Berger, B. Schema: A general framework for integrating heterogeneous single-cell modalities. *bioRxiv* (2019). doi:10.1101/834549
 60. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene expression data analysis. *Genome Biol.* **19**, 15 (2018).
 61. Efron, B. *et al.* Least angle regression. *Ann. Stat.* **32**, 407–499 (2004).

62. Hagberg, A. A., Schult, D. A. & Swart, P. J. Exploring network structure, dynamics, and function using NetworkX. *Proc. 7th Python Sci. Conf.* 11–16 (2008).
63. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society B* **57**, 289–300 (1995).
64. Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. Revigo summarizes and visualizes long lists of gene ontology terms. *PLoS One* **6**, Article number: 7 (2011).
65. Oliphant, T. E. SciPy: Open source scientific tools for Python. *Comput. Sci. Eng.* **9**, 10–20 (2007).
66. Bastian, M., Heymann, S. & Jacomy, M. Gephi: An Open Source Software for Exploring and Manipulating Networks. *Third Int. AAAI Conf. Weblogs Soc. Media* (2009).
doi:10.1136/qshc.2004.010033

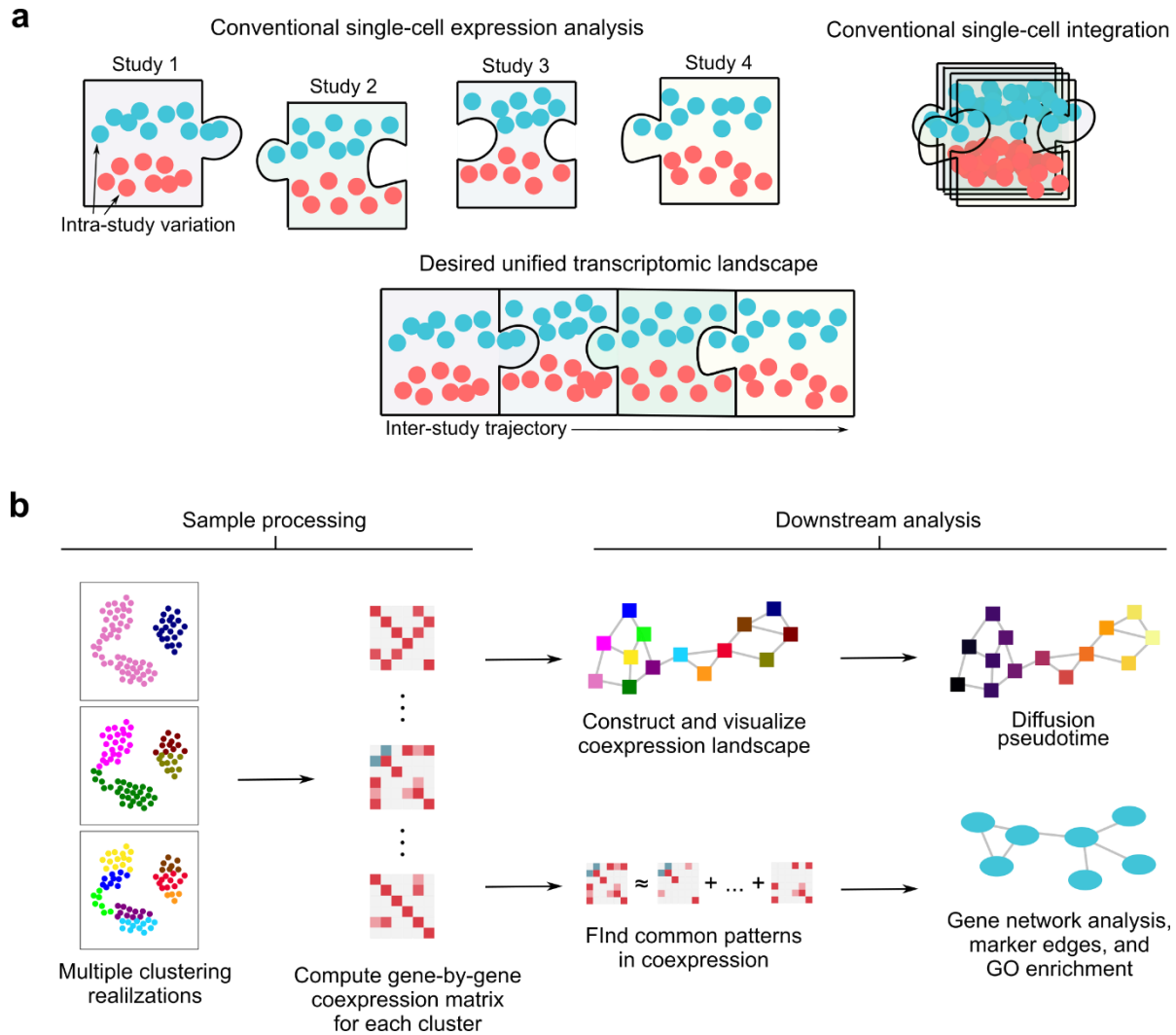


Figure 1. Overview of coexpression-based single-cell transcriptomic analysis.

(a) A conceptual illustration of the difference between attempting to extract biological information from single-studies, each profiling different parts of a larger biological system (“Conventional single-cell expression analysis”); integrative algorithms that attempt to minimize inter-study variation but may also remove overarching biological structure, including temporal dynamics (“Conventional single-cell integration”); and piecing together structure across multiple studies of complex and dynamic biological systems, which we accomplish with single-cell coexpression (“Desired unified transcriptomic landscape”). (b) Overview of coexpression-based analysis, in which the fundamental analytic unit is a group of cells featured by coexpression,

rather than a single cell featurized by expression. Many downstream analyses have analogs in standard single-cell expression analyses.

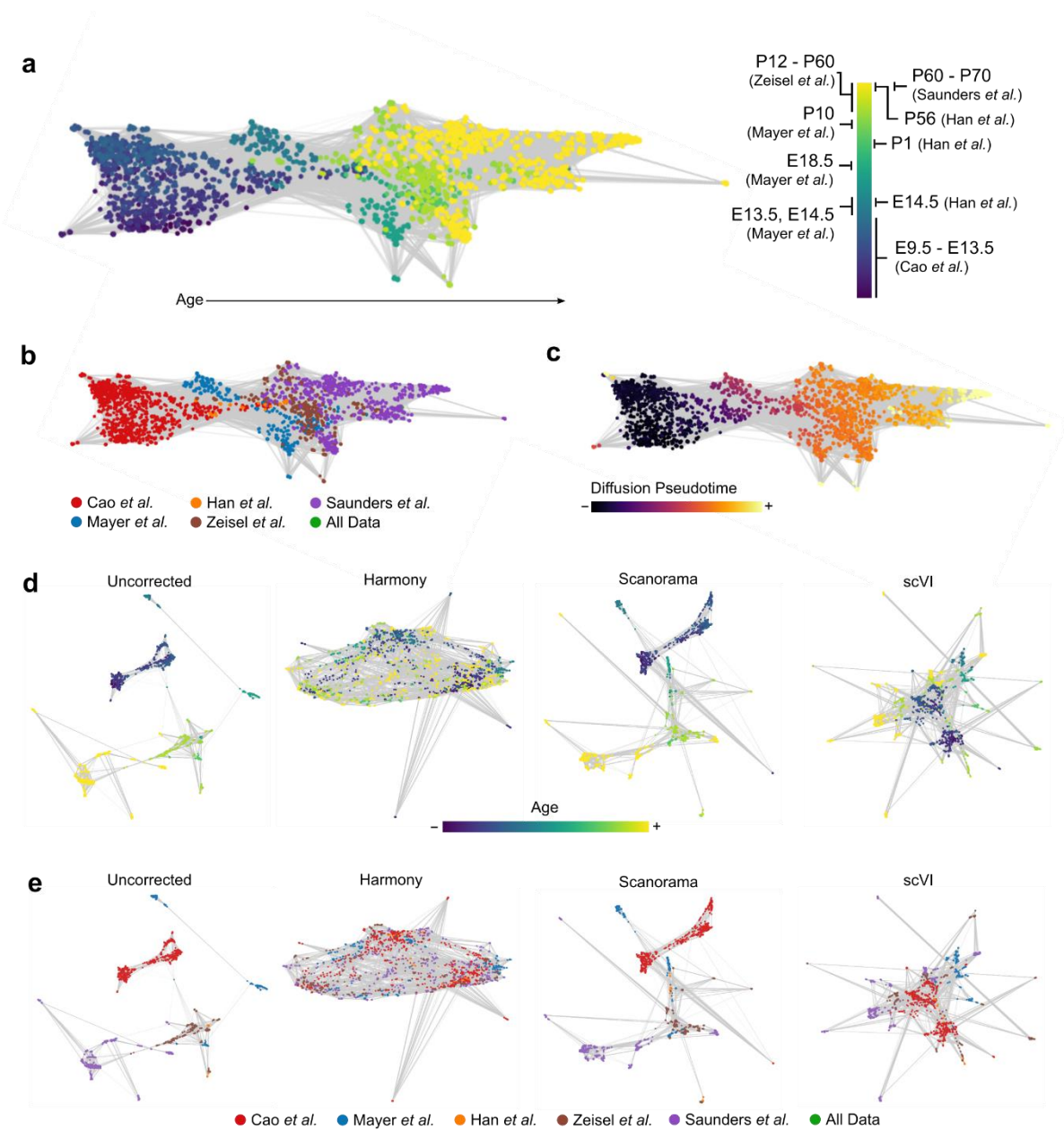


Figure 2. Coexpression landscape of mouse neuronal development.

(a) A force-directed layout of the k -nearest-neighbors graph of panresolution clusters in coexpression space, which we refer to as the “coexpression landscape,” reveals a trajectory consistent with developmental age. (b) Studies are arranged according to order in developmental time, without removing all study-specific signal. (c) Diffusion pseudotime starting from the lowest-age node is strongly associated (Spearman $r = 0.87$, $P < 10^{-308}$, $n = 2,442$ panresolution

clusters) with biological age. **(d,e)** Panresolution clusters in uncorrected expression space and after correction with Scanorama or scVI still show large study-specific patterns without clear age-related structure. Harmony integration removes all study-specific differences including those related to developmental age.

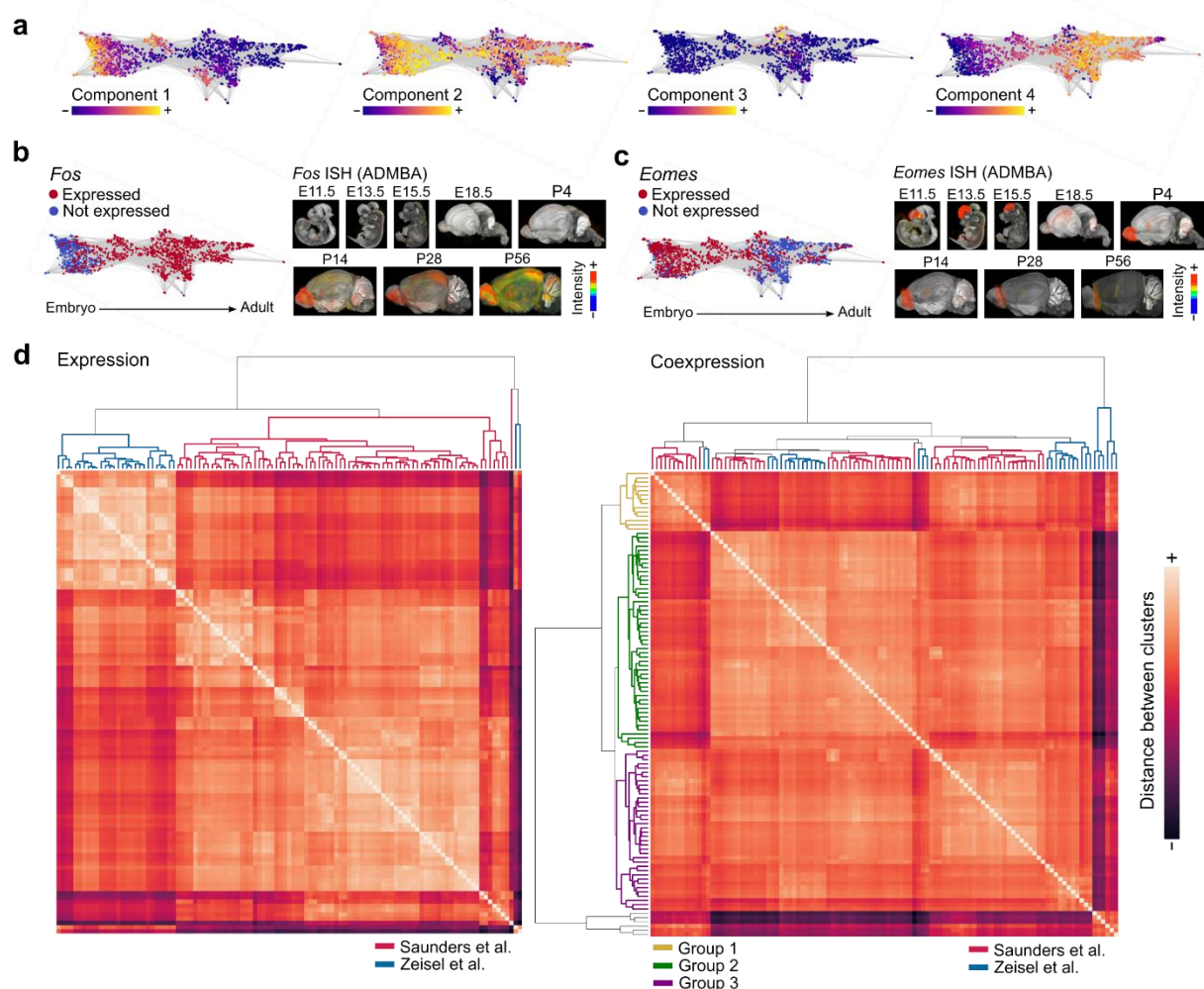


Figure 3. Neuronal trajectory interpretation and cross-study subtype integration.

(a) Dictionary entries highlight different stages of neuronal development. (b,c) We observe positive correlations between diffusion pseudotime, corresponding to development, with the expression of genes such as *Fos* and negative correlation with the expression of *Eomes*. Changes in expression of these genes over development are validated and spatially located by the Allen Developing Mouse Brain Atlas (ADMBA)³⁶. Images show locations and levels of gene expression intensity measured by in situ hybridization (ISH); blue-green is low, yellow-orange is medium, and red is high. (d) Neuronal subtypes featured by mean expression group primarily

according to study while subtypes featured by coexpression group primarily according to three main groups, followed secondarily by study.

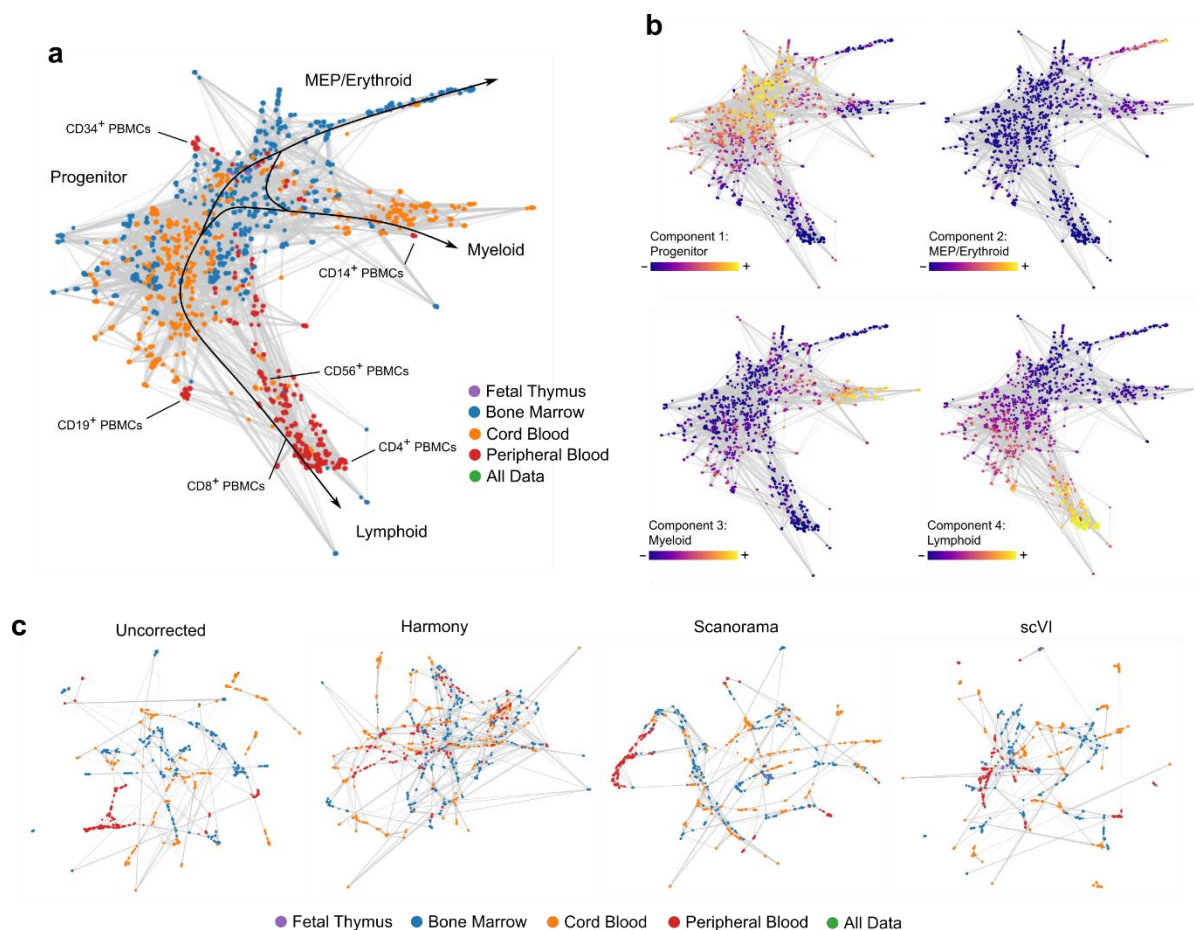
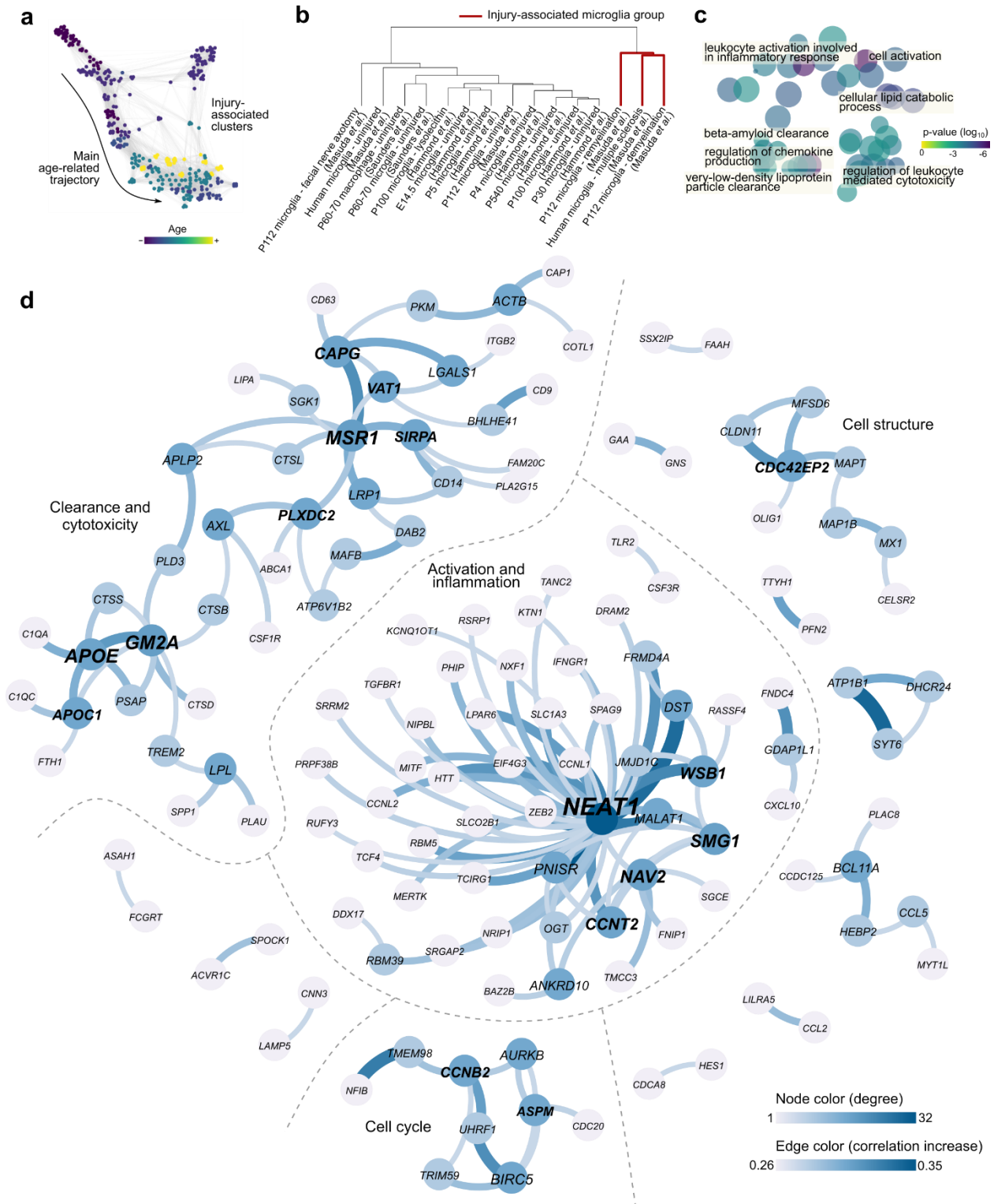


Figure 4. Coexpression landscape of human hematopoiesis.

(a) The coexpression landscape of immune cells from bone marrow, cord blood, and peripheral blood organizes largely according to erythropoietic, lymphopoietic, and myelopoietic lineages. Some of the PBMCs have FACS-derived labels, enabling us to place clusters with known surface markers in various regions of the coexpression landscape (also see **Supplementary Fig. 3**). (b) Dictionary learning of the coexpression matrices separates the coexpression landscape into four main regions; FACS labels and GO process enrichments suggests that these dictionary entries correspond to the different, main stages of hematopoiesis. (c) Existing integrative methods either do not overcome study specific bias (Scanorama and scVI) or obscure the lineage relationships among the four tissues (Harmony).



(a) The coexpression landscape of panresolution clusters reveals a main age-related trajectory as well as off-trajectory outlier clusters from injured tissue. (b) Grouping microglial subtypes reveals a cluster containing injury-associated conditions in both mouse and human microglia. (c) GO enrichment terms are visualized in two dimensional “semantic space” with key terms relevant to disease-associated microglia also displayed. (d) The disease-specific coexpression network reveals functional gene modules related to myelin injury. The top 150 edges in which coexpression increases from a baseline microglial state are arranged into a disease-associated coexpression network; almost all of these associations have not been described by previous studies. Major subgraphs are labeled according to GO terms associated with internal genes. Nodes are colored darker blue with higher degree; gene labels are larger and bolder with higher degree; and edges are thicker and darker blue with a higher increase in correlation from the baseline microglial state.