

TissueFormer: a neural network for labeling tissue from grouped single-cell RNA profiles

Ari S. Benjamin¹ and Anthony Zador¹

*Correspondence:
benjami@cshl.edu
(AB)

¹Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724

Abstract Single-cell RNA sequencing technologies have enabled unprecedented insights into gene expression and are poised to transform clinical diagnostics. At present, most computational approaches for interpreting single-cell data operate at the level of individual cells, predicting labels or properties based on isolated transcriptomic profiles. This approach overlooks a key class of signals: the composition of cells within a sample or defined population. Such signals are often critical for inferring tissue identity, disease state, or other sample-level phenotypes. To address this limitation, we introduce TissueFormer, a Transformer-based neural network that analyzes groups of single-cell RNA profiles to infer population-level labels while retaining single-cell resolution. Applied to predict the cortical area of groups of cells sampled from spatial transcriptomic data from mouse brains, TissueFormer outperformed both single-cell foundation models and machine learning methods applied to pseudobulk and cell type composition. This higher performance enables the automated construction of high-resolution brain region maps in individual animals directly from spatial transcriptomic data. More broadly, TissueFormer provides a framework for predicting any population-level phenotypes which are influenced by cellular diversity and tissue-level organization.

Introduction

Many clinically and biologically important features such as disease state, tissue identity, or treatment response are defined at the level of cell populations but are driven by patterns of single-cell phenotypes. Such features are often predictable from the composition and organization of cells in a population. For example, functionally distinct areas in the cerebral cortex can be distinguished by the relative abundance of cell types (*Yao et al., 2021b*). In clinical medicine, methods for detecting imbalances in cell type proportions such as complete blood count (CBC) tests form the basis of key diagnostics. For example, a high neutrophil-to-lymphocyte count is a widely used marker of inflammation (*Templeton et al., 2014*). Despite their clinical utility, readouts depending on cell composition remain limited to a handful of canonical patterns identified through decades of observation rather than discovered systematically through large-scale, data-driven analyses.

With the increasing availability of single-cell RNA sequencing, it is now possible to measure patterns in single-cell properties at high resolution across large numbers of cells per sample (*Svensson et al., 2018*). This rich granularity presents a tradeoff for the analysis of phenotypes across tissues, samples, or individuals, each of which contains many sequenced cells. At one extreme, many computational pipelines focus on single cells in isolation, taking the a cell profile as input and producing a label specific to that cell. This is the approach taken by most recent single-cell 'foundation models' of mRNA transcription data (*Lopez et al., 2018; Connell et al., 2022; Cui et al., 2023; Theodoris et al., 2023; Rosen et al., 2023;*

39 **Yuan et al., 2024; Schaar et al., 2024; Ito et al., 2025; Hsieh et al., 2024**), reviewed in (**Szalata et al.,**
40 **2024**). These models learn rich single-cell representations by first ‘pretraining’ with an unsupervised task
41 on large datasets before fine-tuning on specific tasks, often resulting in improved performance when
42 adapted for tasks such as cell type classification. However, these models are fundamentally single-cell
43 in design, and as result they are unable to detect signals that emerge only from the *composition* of cells
44 present.

45 An opposing strategy is to first compute summary statistics about cell populations across tissues or
46 samples, then use these for downstream analysis and classification. One such statistic is the average
47 expression, or psuedobulk expression profile, as is frequently used in differential state analysis (**Crowell**
48 **et al., 2020**). Alternatively, one may calculate compositional features about a population, such as cell
49 type frequencies, ratios of cell types, putative cell-cell communication networks, and cell type specific
50 pathway scores (**Cao et al., 2022**). While informative, both psuedobulk and cell composition measures
51 obscure single-cell granularity and rely on the researcher to specify which compressed representation of
52 a population of cells is likely to be informative about sample labels.

53 In order to maximize the potential of single-cell data, it is important that methods combine the
54 transfer-learning capabilities and end-to-end trainability of single-cell foundation models with the ability
55 to compare cells across a population. Certain models exist which are end-to-end trainable and can
56 compare across cells, such as CellCnn (**Arvaniti and Claassen, 2017**), scAGG (**Verlaan et al., 2025**),
57 and ScRAT (**Mao et al., 2024**), but these models are not compatible with existing foundation models.
58 The ability to incorporate knowledge from pretrained models is especially important given the high
59 dimensionality of the inputs when predicting sample-level phenotypes. In each independent sample, the
60 input dimensionality is on the order of the number of cells per sample times the number of genes per cell.
61 By first pretraining on self-supervised tasks, foundation models effectively reduce this dimensionality by
62 establishing a learned prior over the representation of each cell.

63 To bridge these approaches we present TissueFormer, a neural network architecture that analyzes
64 groups of single-cell RNA profiles collectively while utilizing knowledge embedded in pre-trained single-cell
65 foundation models. The model incorporates a module based on the Geneformer architecture for single-
66 cell analysis (**Theodoris et al., 2023**), which can be initialized with pre-trained weights. During processing,
67 each cell is given a learnable representation with this pretrained module. These cell representations
68 are then processed by several further layers of self-attention, allowing the model to learn to attend
69 to the relevant cells and their interactions. Importantly, the output of Tissueformer is invariant to the
70 order of cells presented as input. These design choices allow TissueFormer to learn arbitrary functions
71 from observations of many cells to sample-level phenotypes while simultaneously leveraging pretrained
72 single-cell foundation models.

73 To validate our approach, we applied TissueFormer to apply Allen Brain Atlas labels to regions of
74 the mouse cortex profiled with spatial transcriptomics. This task is a competitive benchmark for the
75 supervised classification of sampled tissues, though it is important to note that clustering or alignment-
76 style methods have also been developed for tissue annotation (e.g. **Biancalani et al. (2021); Hu et al.**
77 **(2021); Lee et al. (2025)**). We found that TissueFormer outperforms both standard supervised methods as
78 well as fine-tuned single-cell foundation models at predicting the labels due to previous brain annotations.
79 Furthermore, the performance of TissueFormer scaled roughly logarithmically with the number of
80 cells input as a group, demonstrating the utility of observing multiple cells. These results highlight the
81 importance of sample- or tissue-level features in transcriptomic analysis and establish a framework for
82 integrating single-cell foundation models into higher-order biological investigations.

83 Results

84 Our primary goal is to design and evaluate a model which is generally applicable to the supervised
85 problem of sample or tissue annotation from single-cell data. In what follows, we first describe the
86 particular characteristics of this problem that make it challenging for standard methods.

87 Multi-cell supervised problem setting

88 A typical problem in single cell analysis is relating mRNA expression levels in each cell to a cellular
89 property such as cell type identity (**Figure 1a**). If we denote the cell-by-transcript matrix of counts as C ,
90 where each row of C_i contains the mRNA transcript counts for a single cell, this problem seeks to find a
91 function f that maps each row C_i (a cell) to a particular value y_i corresponding to a label. For example,
92 y_i could represent the cortical area from which a given cell was obtained.

$$f(C_i) = y_i \quad (1)$$

93 In some problem settings, it is advantageous to use information from many cells. In this multi-cell
94 setting, the characteristic equation relates sets of cells to labels. For spatial transcriptomic data, for
95 example, a set might refer to a group of cells within some spatial distance of one another. Denoting
96 G_i as the i th such group, and $\{c \in G_i\}$ as the set of cells in this group, the classification or regression
97 problem seeks to find a function f that maps each set of cells to its label y_i .

$$f(\{c \in G_i\}) = y_i \quad (2)$$

98 The fundamental difficulty of the multi-cell problem is its high dimensionality. Relative to the single-cell
99 problem, the dimensionality of input examples increases from the number of genes D to $D \times |G|$, where
100 $|G|$ is the number of cells in a group. Meanwhile, the number of independent labels decreases from the
101 total number of cells to the number of disjoint groups, i.e. from N to $\frac{N}{|G|}$. The ratio of input dimensions to
102 independent examples, a proxy for the difficulty of the problem, thus increases as the group size squared,
103 $|G|^2$.

104 Most strategies for the multi-cell problem explicitly reducing its dimensionality through aggregation
105 (as in pseudobulk analyses) or summary statistics (e.g. **Cao et al. (2022)**). However, these strategies
106 destroy information which could be critical in multi-cell problems. Pseudobulk methods obscure the
107 variance in the signal, whereas cell type composition or other such signals are fixed and non-adjustable
108 representations of single-cell profiles and are incapable of observing transcriptional dysregulation within
109 cell types. To take full advantage of single-cell data, a model should ideally be able to attend to any
110 aspect of the full RNA transcriptome in each cell while comparing across the population.

111 It is nevertheless crucial take steps to constrain the distribution of learned functions so as to cir-
112 cumvent the curse of dimensionality in multi-cell prediction problems. One way to address this is by
113 pre-training a foundation model on a self-supervised task before downstream supervised learning. If this
114 self-supervised task is sufficiently similar to the supervised task at hand, the pretrained model weights
115 will be close in weight space to a good solution for the supervised task. This reduces the number of
116 training steps and data needed in the supervised task.

117 Recent related work has introduced several foundation models tailored for spatial transcriptomic
118 data. Because these models observe groups of cells, they can be seen as solutions to a subcase of the
119 multi-cell problem in Equation 2. For example, scGPT-spatial is pretrained to predict gene expression in
120 held-out cells in the local neighborhood (**Wang et al., 2025**). Unlike TissueFormer, scGPT-spatial averages
121 the embedding of several cells in the local neighborhood for the inter-cell problem, rather than employing
122 self-attention across cells, and thus cannot compare compositional signals. HEIST (**Madhu et al., 2025**)
123 and CI-FM (**You et al., 2025**) define graph neural networks over a hierarchical graph spanning both nearby

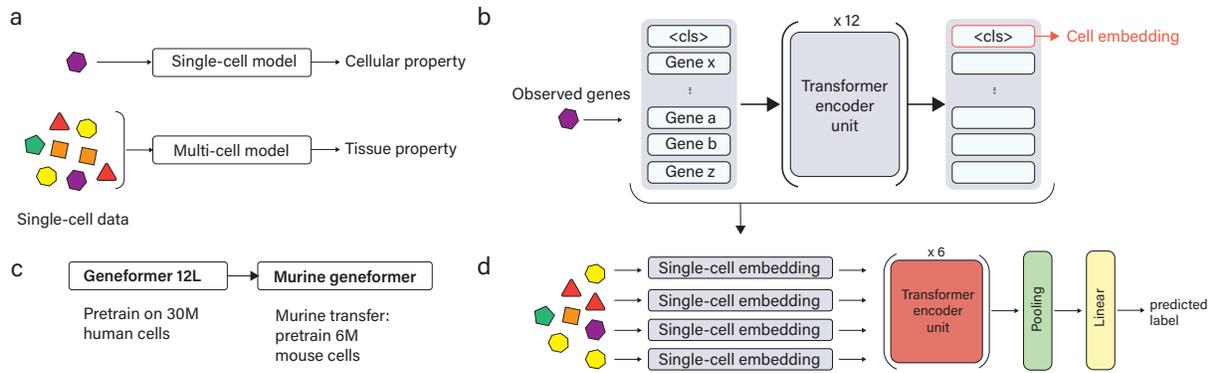


Figure 1. TissueFormer is a Transformer-based neural network to analyze groups of single-cell RNA profiles. a) Single-cell foundation models focus on cellular targets, yet many properties of tissues relate to cellular diversity. b) TissueFormer learns to extract critical information from each cell into a cell embedding vector. The module that extracts these is in architecture identical to Geneformer, a single-cell foundation model. First, cells are processed by a normalized rank-ordering of their genes into a 'sentence' of gene tokens, prefixed by a '<cls>' token to pool gene information, then processed by several transformer layers. c) Experiments in this manuscript use a pretrained single-cell model we call Murine Geneformer, created by adapting the 12-layer Geneformer model (pretrained on human cells) to mouse cells via pairing orthologous genes and further pretraining. d) TissueFormer architecture. Cell embeddings for all cells in a group are fed to 6 Transformer encoding units before average-pooling and a linear readout.

124 cells and their genes, which in principle enables comparative processing across cells. As explained
125 below, TissueFormer uses Transformer layers and not graph neural networks, making it compatible with
126 single-cell foundation models trained on traditional scRNA-seq data. Unlike these works, the design
127 of TissueFormer is not restricted to spatial transcriptomic data and is applicable to any sample-level
128 phenotype prediction problem.

129 TissueFormer architecture

130 TissueFormer is a neural network architecture tailored for groups of single-cell RNA transcriptomic
131 profiles (**Figure 1**). It is able to attend to information within each cell's relative mRNA transcription profile,
132 yet also look across cells to extract information from the diversity and frequency of mRNA profiles in the
133 group. Architecturally, this is accomplished by first applying a 'single-cell module' to extract information
134 from each cell, then attending across cells in further layers. This combination of intra- and inter-cell
135 expressivity creates a regression problem of high dimensionality, yet this is mitigated by the bottleneck
136 architecture and the easy loading of pre-trained single-cell foundation models into the cell embedding
137 module.

138 The building block of TissueFormer which performs these capabilities is the Transformer block
139 implementing self-attention (**Vaswani et al., 2017**). Compared to older neural network architectures like
140 Multi-Layer Perceptrons (MLP), Transformers are distinguished by the shape of their inputs. Whereas an
141 MLP observes vector-shaped inputs (i.e. a vector of mRNA transcript counts), Transformers operate on
142 lists of vectors, or equivalently, on matrix-shaped inputs. Self-attention acts to select which vectors in
143 the list are most relevant in the current context.

144 In order for single-cell data to be compatible with Transformers, each cell's vector of mRNA transcript
145 counts c_i must be expanded into a list of vectors. Multiple strategies for this conversion are possible
146 for single-cell data, as reviewed in **Szalata et al. (2024)**. Here we use a rank-ordering approach as in
147 Geneformer (**Theodoris et al., 2023**). Each cell's expressed genes are ordered by their relative expression
148 normalized by their median expression in a pretraining corpus. Each gene name in this ordered list is

149 then mapped to a high-dimensional embedding vector via a learnable dictionary. While alternatives to
150 rank encoding have been tested in other foundation models, rank encoding is common in bioinformatic
151 pipelines and has generally been found to be robust to common sources of experimental variability
152 (**Ballouz et al., 2015**).

153 The single-cell module within TissueFormer is designed to compress each cell's unique and relevant
154 information into a *cell embedding* for later processing (**Figure 1b**). This is achieved by prefixing to each
155 cell a special token named the <cls> token, following standard design choices (**Devlin et al., 2019**). The
156 <cls> token instantiates a working space for the accumulation of information from all genes. During
157 learning, any error-corrective information passed backwards is bottlenecked through the cell embedding,
158 ensuring it collects the necessary information from each cell. In sum, the cell embedding vector which is
159 passed to later layers can be written as a function of the mRNA transcript count vector c_i as:

$$\text{CellEmb}[c_i] = \text{SingleCellModule}[\text{RankEncoding}[c_i]][\text{<cls>}] \quad (3)$$

160 The single-cell module shares the Geneformer architecture, meaning that any trained single-cell
161 model can be loaded into this module. This gives TissueFormer the ability to incorporate model weights
162 and biological knowledge from pretrained single-cell foundation models. In this work, all experiments
163 unless otherwise indicated use weights from the Geneformer 12L model, a 12-layer BERT network trained
164 on a masked prediction objective on 30 million human cells (**Theodoris et al., 2023**). We then further
165 pre-trained these weights using an identical objective function but on a collection of datasets of mouse
166 tissue totaling 6 million cells (**Figure 1c**). To facilitate the transfer of knowledge across species, we
167 initialized the gene tokens for mouse genes with their orthologous human genes where available. We
168 call this resulting pretrained single-cell model 'Murine Geneformer'.

169 TissueFormer attends to information across cells by processing each cell embedding with several
170 additional Transformer layers (**Figure 1d**). At this stage, TissueFormer is invariant to the order of cells
171 presented to the model, seeing them only as a bag or set of cells. The output of any Transformer is
172 inherently equivariant with respect to the order of the input embedding vectors (i.e. rearranging the inputs
173 causes the outputs to be rearranged in the same order). In order to attain invariance to cell order, such
174 that rearrangements to the order have no effect on the output, the output of the Transformer layers is
175 pooled by averaging over cells. A final step is a linear classification layer for classification problems.
176 Collectively, these design choices allow TissueFormer to attend to any relevant information in each cell's
177 transcriptome in a context-dependent manner based on the diversity of cells present and the problem at
178 hand.

179 Taken together, the equation representing TissueFormer can be summarized as:

$$y_i = \text{LinearClassifier}[\text{AvgPool}[\text{Transformers}[\{\text{CellEmb}[c_j] \text{ for } c_j \in G_i\}]]] \quad (4)$$

180 Cortical annotation

181 The mammalian brain is often partitioned into regions with distinct functional roles, connectivity, and
182 cellular composition. In any particular brain, variations in neuroanatomy may arise due to genetic factors
183 as well as due to the animal's environment. These variations underlie the evolution of innate behavior
184 and reflect the aptitude for adaptation within a lifetime.

185 Current methods for the annotation of single brains have certain well-known limitations. One common
186 approach is to register a brain into the average-brain template of the Common Coordinate Framework
187 (CCF), a 3D reference coordinate system constructed from over 1,000 reference mouse brains (**Wang et al.,**
188 **2020**). Once transformed into CCF coordinates, one can visualize the area boundaries of the Allen Brain
189 Atlas, which again reflect a canonical, average brain. The drawback of this approach is that it obscures
190 individual variation, especially the relative size of brain areas to one another. An alternative approach

191 which enables individual variability is to employ experimental methods that interrogate only a small
192 brain region in each animal—for example, focal viral projection tracing (**Xu et al., 2020**) or two-photon
193 functional imaging for a particular functional modality like visual responses (**Kalatsky and Stryker, 2003**).
194 While accurate, these methods' limited spatial scope prevents a holistic assessment of anatomy across
195 the entire cortex.

196 An alternative strategy that balances whole-brain coverage with single-animal individuality is to
197 perform *in situ* single-cell RNA sequencing across a single brain. This method modernizes a classic
198 approach of defining neuroanatomy through differences in cell type composition, an idea that goes back
199 at least to Brodmann's brain atlas (**Brodmann, 1909; Zilles and Amunts, 2010**). Currently, the technical
200 efficiency of spatial transcriptomics is now high enough to allow the profiling of cells across an entire
201 mouse brain, and furthermore cheaply enough to cover multiple brains in single studies (**Chen et al.,**
202 **2024**). Given such datasets, the remaining challenge in a robust pipeline for labeling brain areas is
203 computational in nature.

204 To investigate whether TissueFormer could serve this purpose, we applied TissueFormer to predict
205 the brain region of cells from their transcriptomes as characterized in a recent dataset of whole-brain
206 spatial transcriptomics of the brains of several animals (**Chen et al., 2024**). Over 1 million cells across
207 the brain, largely in the left hemisphere, were profiled in each of eight animals, four of which were raised
208 in normal rearing conditions. Each cell profile contains counts from a panel of 104 genes selected
209 for their ability to distinguish cell type clusters in excitatory cells. With multiple separate brains and
210 whole-cortex coverage in one hemisphere, this dataset enables a test of whether individual differences
211 in neuroanatomy can be resolved with spatial transcriptomics.

212 In addition to mRNA counts, this dataset also contains labels of the brain area of each cell as
213 established through a registration into the Common Coordinate Framework (CCF). Each brain in **Chen**
214 **et al. (2024)** was previously registered to the CCF using standard software that defines a non-rigid
215 smooth deformation from raw slice coordinates into the CCF. This step normalizes differences in gross
216 anatomy and brain size. Importantly, once transformed into CCF coordinates, one may associate the
217 topography of a new brain with the Allen Brain Atlas to obtain brain area annotations (**Wang et al., 2020**).
218 These CCF annotations provide the labels we use for training.

219 It is important to note that CCF boundaries may not reflect ground truth neuroanatomy. While CCF
220 registration accounts for overall brain shape, it does not account for more fine-grained differences
221 such as the relative size of brain areas or a shift in a single boundary in one animal versus another.
222 Nevertheless, these labels are sufficiently accurate to allow a comparison of computational methods. It
223 is possible as well that the 'errors' of a model trained on such data will in fact represent a more accurate
224 annotation than the CCF, a possibility we will return in **Figure 3** in which we examine the predictions of
225 TissueFormer on held-out brains.

226 Single-cell profiles are not informative of cortical area

227 The utility of attending across cells can be made clear by comparing with the performance of machine
228 learning methods which are trained to predict the area given only a single cell. To investigate this, we
229 first split the data into training data (90% of cells in 3 brains), validation data (the remaining 10% of the
230 same 3 brains), and a held-out test brain (**Figure 2b**). We then trained a range of methods on the task of
231 predicting the cortical area of a single cell from its transcriptome.

232 The methods tested ranged in complexity from simple heuristics to modern single-cell foundation
233 models (**Figure 2c**). We included three standard machine learning benchmarks, namely, a logistic regres-
234 sion model, a random forest model, and k-neighbors classifier. These models map the vector of mRNA
235 transcript counts to the label after a $\log(1 + x)$ transformation and z-scoring, and their hyperparameters
236 were tuned on the validation set. We also trained two models which make use of the cell types of the

237 cell in question. These cell types were previously constructed at three levels of granularity (*Chen et al.*
238 (2024)); we selected the finest level. The ‘cell type model’ predicts the area of a cell based on the most
239 common area of cells of the same type in the training set. The ‘cell type k-neighbors’ model labels the
240 target cell based on the area of the k-closest neighbors with the same cell type in the training set. Finally,
241 we fine-tuned Murine Geneformer to predict the area of each cell.

242 Despite the diversity of modeling approaches, all six models performed poorly, with classification
243 accuracies ranging from 20–31%. This highlights the challenge of predicting cortical area from single-cell
244 transcriptomes, consistent with previous findings that many cell types are broadly distributed across the
245 cortex (*Tasic et al., 2018; Yao et al., 2021b; Chen et al., 2024*). While some types are more localized and
246 thus were predicted with higher accuracy (for example, RSP/ACA or Retrosplenial/Anterior Cingulate
247 excitatory cells), overall performance was low. Among the models tested, the Transformer-based Murine
248 Geneformer achieved the highest accuracy, modestly outperforming the others. This is consistent with
249 recent results showing strong performance of large pretrained models on downstream tasks (*Theodoris*
250 *et al., 2023*). These results motivated our turn toward models that integrate information across cells.

251 **Annotating cell groups with TissueFormer**

252 To apply TissueFormer, we first created a pipeline to group cell sharing a label. Specifically, we grouped
253 cells into cylindrical columns which fall inwards from the surface of the cortex (*Figure 2d-e*). Note that
254 this strategy can easily be generalized to other shapes for spatial data from other tissues, or to other
255 sampling strategies for non-spatial data. The cerebral cortex can be viewed as a layered cake draped
256 over the exterior of the mouse brain, with functional areas as slices (*Figure 2a*). Thus, cortical columns
257 contain cells in the same functional area. Our pipeline programmatically selects a single cell, then grows
258 a cylinder centered on that cell up to the minimum radius needed to contain N total cells. Each columnar
259 group of N cells is labeled with the most common CCF annotation of cells in that group.

260 We next trained TissueFormer to predict the group’s label from the set of transcriptomic profiles
261 within the group. We again cross-validated by training on 90% of cells from three brains, validating on
262 the remaining 10% of cells, and testing on one held-out brain. As we varied the number of cells in each
263 group from $N = 1$ to $N = 256$, we observed that the validation accuracy of TissueFormer increased from
264 around 35% to nearly 80% accuracy (*Figure 2f*). On test brains, the accuracy Tissueformer also increased
265 markedly with group size from 30% to over 60% accuracy (*Figure 2g*). Providing the cortical depth of
266 each cell as additional input to the model did not improve results, suggesting that depth represents
267 redundant information (*Figure 2—figure Supplement 1f*). Note that random guessing on this 1-of-42
268 classification task yields 8% accuracy due to imbalanced classes. Similar relative accuracies were seen
269 when training with a class-balancing weighted objective (*Figure 2—figure Supplement 1c*). These results
270 demonstrate that TissueFormer is able to dramatically outperform single-cell models at predicting tissue
271 identity.

272 As the group size increased from $N = 1$ to $N = 256$, the accuracy increased smoothly and roughly
273 logarithmically with group size, indicating a predictable scaling law with more cells. Performance
274 saturated at $N = 128$ cells, possibly because at this scale and at this density of cells per brain the groups
275 sometimes straddle boundaries, at which point group identity begins to lose meaning. Indeed,
276 at $N = 256$ over 70% of groups contained a boundary (*Figure 2—figure Supplement 1g*). Higher sampling
277 densities would enable larger group sizes in smaller locales, and thus possibly an even higher saturating
278 accuracy. Due to this problem of spatial density, it is not possible to verify the maximum group size at
279 which cells contain no further information about area. Overall, the considerable increase of performance
280 with group size highlights the importance of the ability to integrate and compare information across
281 single cells.

282 The cell embedding module in the above experiment was previously pretrained on a masked mRNA

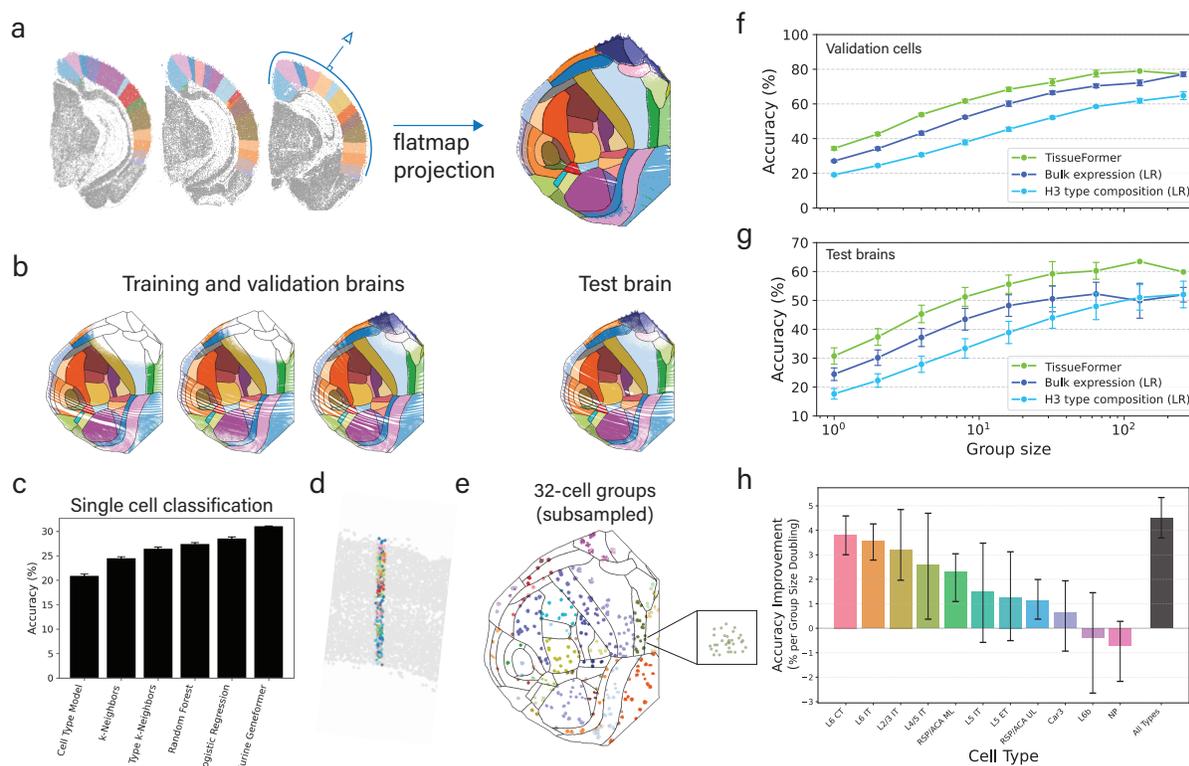


Figure 2. Comparison of accuracy at predicting brain regions shows that TissueFormer outperforms other methods. **a**) Cortical cells from coronal slices of one hemisphere (left, colored by area) are visualized as a 'flatmap', a top-down projection. Color legend is in **Figure 3.b**) Models were evaluated using 4-fold crossvalidation with 3 or 4 brains used for training and validation labels, 1 of 4 brains held out for testing. Shown here are true labels for each cell. **c**) Comparison of accuracy when predicting the area of a single cell in a test brain from its transcriptome shows that all methods perform poorly, though with Murine Geneformer performing best. **d**) An example cell group ($N = 256$) plotted in slice coordinates, colored by cell type. **e**) A sample of groups ($N = 32$) plotted in flatmap coordinates, here colored by the modal area of the group. **f**) Validation accuracy as a function of group size shows that TissueFormer outperforms logistic regression (LR) given either pseudobulk transcription or a histogram of cell types. **g**) Test brain accuracy. Error bars in f-g represent standard deviation across 4 folds. **h**) The average rate of test accuracy improvement with group size for the above TissueFormer (black bar) compared to the accuracy curve slopes of several TissueFormer models trained on homogeneous groups containing only a single cell type. Accuracy curves are visible in supplement 1.

Figure 2—figure supplement 1. Effect of pretraining on TissueFormer, performance of random forest benchmarks, and controls for unequal area size and sampling density.

283 transcript prediction task on several other datasets containing tens of millions of human and mouse cells.
284 In principle, this pretraining step could help, hinder, or have little effect on our current supervised task of
285 predicting brain area. To investigate its impact, we also trained a TissueFormer model from a random
286 initialization without any pretraining. Surprisingly, the performance of this randomly initialized model was
287 indistinguishable from the pretrained model for all group sizes (**Figure 2—figure Supplement 1d**). This is
288 likely due to the large number of cells in our labeled dataset. To confirm this, we varied the number of
289 training cells used to train both a pretrained and a randomly initialized TissueFormer from one thousand
290 cells to the full dataset of over two million. We found that pretraining indeed offered an advantage, but
291 only in the intermediate range of less than one million cells (**Figure 2—figure Supplement 1e**). Thus, this
292 spatial transcriptomic dataset is large enough that transfer learning from Murine Geneformer offers no
293 advantage.

294 We next compared the performance of TissueFormer with two standard methods for processing
295 transcriptomic data of groups of single cells. We first examined pseudobulk analyses, which average
296 the single-cell profiles in the group to approximate the measurements of a traditional bulk sequencing
297 experiment. Logistic regression trained on pseudobulk vectors underperformed TissueFormer, yet still
298 showed a nearly logarithmic increase of accuracy with group size (**Figure 2f**). A random forest model
299 trained on the same data underperformed logistic regression (**Figure 2—figure Supplement 1b**). Next,
300 we trained logistic regression and random forests to predict area from the cell type composition of a
301 group, which was represented as a histogram of cell types. We used the finest-grain categorization of
302 cell types from **Chen et al. (2024)** for this analysis. Overall, type composition was less informative of
303 area than pseudobulk expression (**Figure 2f**). Intriguingly, each method displayed a logarithmic increase
304 in accuracy with group size despite the differences in their representations of RNA transcription.

305 The logarithmic increase in accuracy with group size across all methodologies could in principle
306 be driven by two effects. One possibility is that the diversity and particular composition of cells pro-
307 vided a new signal not available to single cell models. An alternative possibility is that cells contained
308 independent measurement noise which can be averaged away. Some evidence to the first possibility
309 is that benchmarks based on cell type composition alone showed such an increase (**Figure 2f-g**). To
310 further distinguish these factors within TissueFormer, we reasoned that we could artificially restrict
311 the diversity of cells in each group and then re-test a method's performance. If the scaling were due to
312 measurement noise alone, then this would have little effect and performance would likewise increase
313 with a similar slope on a log-linear plot. We therefore constructed a new method of constructing groups
314 such that they contain only a single cell type from the second level of granularity in the cell type hierarchy
315 ("H2" types). Note that there is still significant diversity within each H2 type, albeit much less than the
316 general population. After training a TissueFormer model from scratch on these homogeneous groups, we
317 found that the test set accuracy given a group of homogeneous type scaled more slowly with increasing
318 group size N , on average (**Figure 2h**). Groups of cells of certain types such as NP cells showed little
319 improvement with larger size. Other types, such as L6 IT cells, nevertheless show comparable increases,
320 likely indicating higher and more informative diversity in that population. The areas of some cell types
321 are more easily predictable than other types due regional localization (**Figure 2—figure Supplement 1h**).
322 These analyses confirmed that the increase in accuracy with group size was at least partially driven by
323 comparative signals across a cell group rather than solely by averaging away technical noise.

324 The decrease in accuracy between validation and test brain accuracy could have arisen due to multiple
325 reasons. A first reason is the differences in spatial coverage of data from each brain, visible in **Figure 2b**,
326 which may cause test brains to contain areas not in the training data. However, correcting for this effect
327 by only testing data with high density in the training set yielded improvements in accuracy of at most
328 a few percentage points (**Figure 2—figure Supplement 1a**). Secondly, this could have reflected other
329 batch effects across brains affecting single-cell measurements, effectively making this test brain an

330 out-of-distribution test for the methods. Finally, it is also possible that ‘errors’ in classification were due
331 to actual changes in neuroanatomy which were correctly predicted by the model but were mislabeled by
332 the pipeline of registering each brain to the CCF. We investigate this possibility in the following section.

333 **Predicted cortical maps**

334 TissueFormer can be used as an automated pipeline to create maps of cortical anatomy from single-cell
335 transcription after being trained on reference brains. Note that annotating from reference brains is a
336 supervised task, not to be confused with the unsupervised task of spatially clustering single-cell data
337 (see Discussion). In **Figure 3a**, we demonstrate this capability and display the predicted cortical maps for
338 each of 4 animals treated as a held-out test set. Pixels in these maps were colored by a weighted average
339 of nearby cells’ predictions (see Methods). For visualization we used the predictions of TissueFormer
340 trained with a class-size-balancing objective ensuring that the prior probability was uniform across
341 all areas regardless of area size. We found that all four maps were coarsely similar to the reference
342 annotation (see **Figure 2a**), supporting the finding in **Figure 2** that TissueFormer is able to predict brain
343 area more accurately than other approaches.

344 A close inspection of these cortical maps revealed several notable differences in the predicted
345 anatomy of each animal. For example, the boundary between the somatosensory cortex and the motor
346 cortex is consistently shifted in the posterior-lateral direction in brain 1 relative to the Common Coordinate
347 Framework (CCF) annotations, yet is consistently shifted in the opposite (anterior-medial) direction in
348 brain 2 (**Figure 3b**). Additionally, the primary visual area is shifted medially in brain 2, yet laterally in brains
349 3 and 4. Likewise, the primary auditory area is shifted medially in brain 2 yet is expanded and shifted
350 laterally in brain 4. These shifts are a predominant cause of why test brain ‘accuracy’ in **Figure 2g** is lower
351 than validation accuracy.

352 While it is possible that the discrepancies between the CCF annotation and the predicted annotation
353 are ‘errors’ of the model, it is also possible that this reflects true differences in anatomy between animals
354 or errors in registration with the CCF. The CCF represents the average anatomy of over 1,000 mouse brains,
355 obscuring potential individual variability between brains. To investigate this possibility, we examined the
356 somatosensory-motor boundary in more detail. This boundary can be canonically identified by differences
357 in cell type composition. In particular, motor cortex is classically considered to lack a Layer 4 (**Brodmann,**
358 **1909**), and thus should show a large decrease in the density of Layer 4/5 IT excitatory cells relative to
359 Layer 5 IT cells. This easily verifiable boundary was qualitatively visible in coronal slices colored by cell
360 type (**Figure 3c**). To confirm its location, we examined the density of these two cell types along slices
361 taken from brain 1 and brain 2 which intersect the sensory/motor boundary (**Figure 3d**). We found that
362 the CCF boundary was inconsistent across brains relative to the ratio of cell types, but TissueFormer’s
363 predicted maps showed a consistent 1:1. Furthermore, in the CCF boundaries, we observed a large shift
364 between animals in the ratio for large distances on either side (**Figure 3e**). In contrast, the density ratios
365 were consistent in their alignment to the predicted boundaries. Thus, a closer inspection of cell type
366 distributions around this verifiable boundary revealed a closer alignment with our predicted atlas than
367 with the boundaries due to CCF alignment.

368 To provide further verification of the predicted cortical maps, we examined the spatial distributions of
369 individual genes known to correlate with area identity (**Figure 3—figure Supplement 2**). For example, *Tshz2*
370 forms a striking medial-to-lateral gradient with high abundance in anterior- and posterior-cingulate/retrosplenial
371 areas with steep drops at the boundaries to motor, somatosensory and visual fields (**Yao et al., 2021b**).
372 Across all four brains, the inter-animal differences in *Tshz2* aligned better with the predicted boundary
373 than the CCF area boundaries. Other genes which are also selectively expressed in retrosplenial areas,
374 such as *Coro6* and *Zfpm2* (**Chen et al., 2022**), shared this pattern, as did genes with selectively
375 expressed in the dorsal but not the ventral retrosplenial area, such as *Nell1* and *Zmat4* (**Hashikawa et al.,**

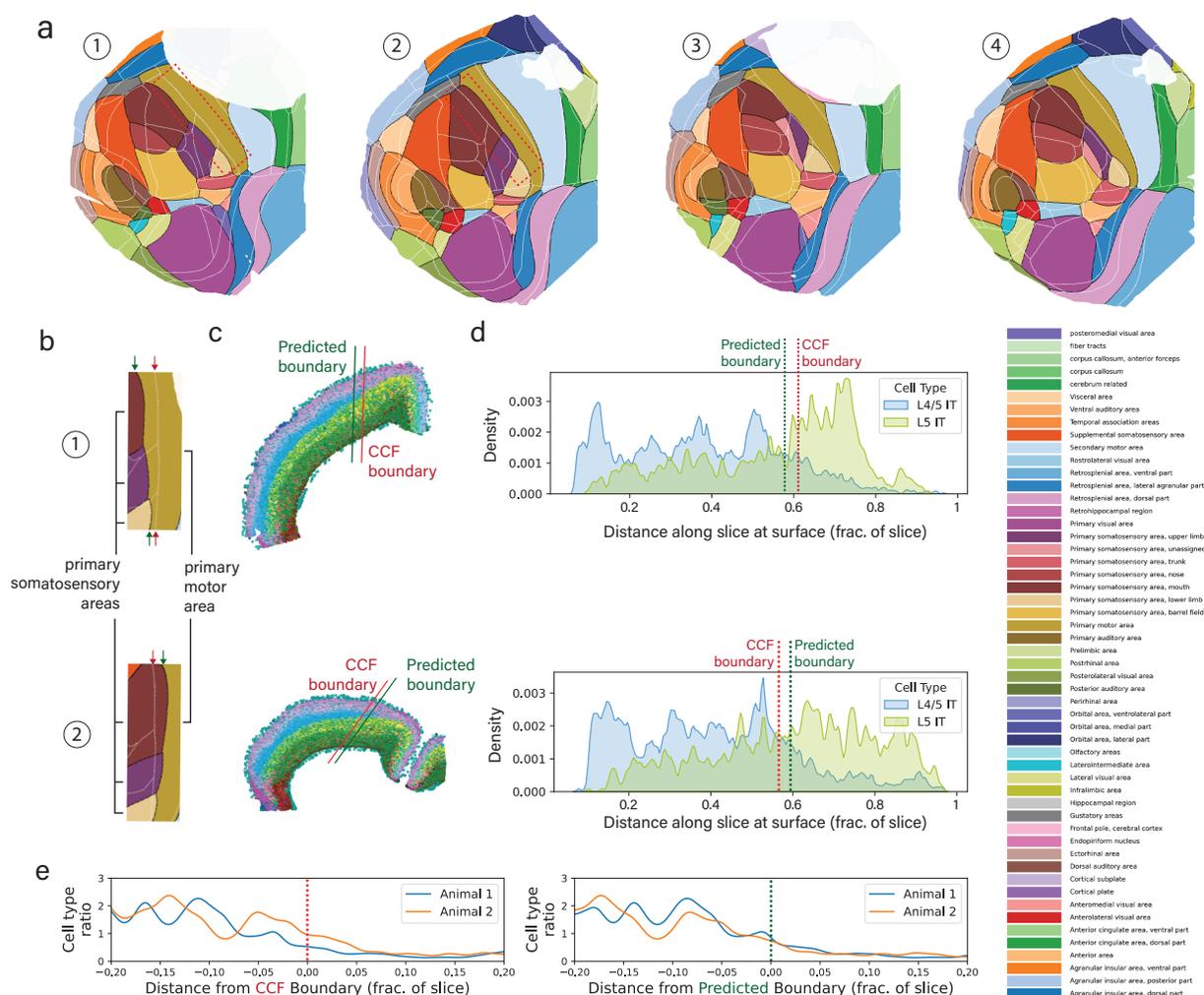


Figure 3. Creation of cortical maps from spatial transcriptomics using TissueFormer. **a)** Predicted cortical areas for each brain when held-out as a test brain. Predicted boundaries are in black, and the reference Common Coordinate Framework (CCF) boundaries are overlaid in white. Areas with low spatial density are masked. **b)** In brains 1 (top) and 2 (bottom), the somatosensory-motor boundary is displaced relative to the CCF, but in opposite directions. **c)** Example slices from brains 1 and 2 which contain the somatosensory-motor boundary boundary. Points are cells, colored by cell type with similar cell types assigned similar colors (see Methods). **d)** Spatial density along the slices of two cell types, L4/5 IT and L5 IT, in the same colors and slices as (c). **e)** The cell type ratio (L4/5 IT / L5 IT) aligned to the CCF boundary (left) or predicted boundary (right), showing an inter-animal shift in CCF vs. inter-animal consistency in the predicted maps.

Figure 3—figure supplement 1. Predicted areas of single cells, highlighting discrepancies (‘errors’) with CCF labels.

Figure 3—figure supplement 2. Comparison to maps of mRNA transcription of single genes.

376 **2020**). Meanwhile, the sensory/motor border analyzed in **Figure 3** was correlated to the expression
377 density of *Rorb1*, *Brinp3*, and *Rcan2* (**Cederquist et al., 2013**), all of which showed consistent shifts
378 across animals with the predicted maps but not with CCF boundaries. These examples of genes with
379 clear spatial boundaries in their patterns of expression lend support to the predicted maps' veracity.

380 Discussion

381 We developed TissueFormer, a neural network architecture that attends to groups of single-cell tran-
382 scriptomic profiles in order to predict a label or annotation common to that group. In tests of predicting
383 the brain area of a group of nearby cells in a multi-brain spatial transcriptomics dataset, TissueFormer
384 achieved much higher accuracy rates compared to methods which predict the area of a single cell from
385 its transcriptomic profile. For example, on held-out test brains TissueFormer achieved 60% accuracy with
386 groups of 64 cells compared to 30% for a matched single-cell model. Furthermore, accuracy steadily
387 increased with group size. TissueFormer also outperformed machine learning methods trained on the
388 pseudobulk expression or cell type composition of the same group of cells. These results highlight the
389 importance of attending across cells when predicting tissue-level properties.

390 Because TissueFormer accepts arbitrary cell groupings, it can flexibly model phenotypes defined by
391 anatomical contiguity (as in spatial transcriptomics), by patient sample (e.g., blood draws, biopsies), or by
392 dynamic windows in longitudinal sampling. These applications include a wide range of use cases, from
393 immunomonitoring in inflammatory disease to potentially the detection of cancer states (**Weinkauff et al.,**
394 **1999; Sun et al., 2025**). As increasingly comprehensive healthy and disease cell atlases become available,
395 population-aware models like TissueFormer offer a principled route to integrate those resources into
396 predictive diagnostics that bridge cellular resolution and clinical decision-making (**Dann et al., 2023**).

397 TissueFormer provides a promising tool for automated brain mapping from spatial transcriptomics.
398 Compared to the labels due to Common Coordinate Framework (CCF) registration, TissueFormer's
399 predicted labels displayed individual variability that better correlated with inter-animal differences in local
400 gene expression and cell type distributions, despite being trained to predict CCF labels. This tool could
401 thus expedite the study of the individual differences in the size and location of transcriptomically-defined
402 regions in the brain.

403 Our approach to brain mapping is a supervised approach in which reference brains supply ground
404 truth labels. Note that this is distinct from the approach of aligning tissue without regard to area labels, as
405 in Tangram (**Biancalani et al., 2021**) and CAST (**Tang et al., 2024**). An important caveat of the supervised
406 approach is the lack of perfect training data. No datasets of spatial transcriptomics are yet available
407 which have animal-by-animal annotations of neuroanatomy verified by additional modalities such as
408 projection tracing or functional imaging. Without access to such datasets, we instead trained to predict
409 the imperfect CCF labels. In general, there is no formal guarantee that the predictions on held-out brains
410 should reflect a better assessment than the CCF itself. While we are compelled by our empirical analysis
411 to trust the model, this must be done with caution until a model can be trained on a dataset with multiple
412 brains individually annotated with complementing experimental modalities.

413 In the current study, we did not take up the question of whether different boundaries than the Allen Brain
414 Atlas would be in any sense better, as in recent spatial clustering methods. Current boundaries are widely
415 used and important for shared nomenclature. However, the automated clustering and identification of
416 spatially homogeneous regions within spatial transcriptomic dataset is an interesting and open research
417 question. Several algorithms have already been developed for this purpose (**Dries et al., 2021; Zhao**
418 **et al., 2021; Chitra et al., 2025; Hu et al., 2021; Dong and Zhang, 2022; Singhal et al., 2024; Jackson**
419 **et al., 2024**). When algorithms from this family are applied to the mouse cortex, the resulting regions
420 correspond more to cortical layers than to functional regions (**Ortiz et al., 2020; Partel et al., 2020; Lee**
421 **et al., 2025**). This is consistent with there being larger differences in gene expression across layers

422 than across cortical areas. In contrast, the functional specialization of neurons in cortex varies more
423 across area than across layers within an area; neurons within a cortical column typically mediate the
424 same cognitive tasks and furthermore respond to similar features within that task (**Mountcastle, 1997**;
425 **Callaway, 1998**). Amid this dichotomy between transcriptomic and functional contiguity in space, we
426 chose in this manuscript to classify functional areas for practical utility and for the reason of it being a
427 benchmark task that leverages transcriptomic diversity within cellular ensembles. Nevertheless, it would
428 be interesting in future work to train TissueFormer with self-supervised or contrastive objectives so as to
429 discover brain areas *de novo*.

430 **Methods**

431 **TissueFormer architecture**

432 TissueFormer contains two main components: a single-cell module and an aggregation module. In
433 this sense it can be viewed as an end-to-end trainable hierarchical model above single-cell foundation
434 models.

435 The single-cell module is based on the BERT architecture, an encoder-only Transformer-based neural
436 network (**Devlin et al., 2019**). We use identical hyperparameter settings as used by Geneformer 12L
437 (**Theodoris et al., 2023**), allowing pretrained weights to be readily employed. Specifically, we use a
438 12-layer BERT network with ReLU activations, hidden layer sizes of 512, and a context window of up to
439 2048 genes. Each Transformer uses multi-head attention with 8 heads, and the inputs are appended
440 with sinusoidal positional embeddings.

441 The aggregation model pools across single-cell embeddings. As a first step, the first token from
442 the single-cell module for each cell, which is always the `<c1s>` token in the input token list, is extracted.
443 This group is then fed as input to a stack of Transformer layers (without positional embedding so as to
444 preserve order equivariance). Each layer employs LayerNorm normalization and GELU nonlinearities.
445 Finally, the output is averaged over cells, resulting in a single vector for the group of cells. This is given as
446 input to a single linear classifier layer, outputting unnormalized log probabilities of each of the potential
447 labels. All classification tasks use a cross-entropy objective.

448 In the present manuscript, we used the following hyperparameters. The aggregation module had 6
449 Transformer layers, each with a hidden size of 768 units. All training runs on the classification task used
450 the Adam optimizer with a learning rate of 0.001 and a linear learning rate decay following a learning rate
451 warmup with warmup ratio 0.1. Our implementation relies on the Hugging Face Transformers library
452 (**Wolf et al., 2020**), based on Pytorch (**Paszke et al., 2019**). We used Hydra for configuration management
453 (**Yadan, 2019**).

454 **Single-cell pretraining**

455 The single-cell module can be pretrained on single-cell data using a self-supervised objective. In our
456 experiments we use a version of Geneformer adapted for mouse brain tissue. We first downloaded the
457 Geneformer 12L model, which was trained to predict the identity of randomly masked genes (15% of
458 total genes per cell) on a corpus of 30 million human cells. We then adapted this to data mouse tissue
459 by repeating this methodology on a hand-curated dataset containing over 6 million mouse cells. Using
460 the CellXGene Census explorer (**Program et al., 2025**), we aggregated data from several publications
461 totaling over 6 million cells (**Yao et al., 2021b**; **Govek et al., 2022**; **Kozareva et al., 2021**; **Steuernagel**
462 **et al., 2022**; **Yao et al., 2021a**). After tokenizing (see below), we then trained the Geneformer architecture
463 on this dataset with an identical masked gene prediction objective. Models were trained for 3 epochs
464 using AdamW, a learning rate of 0.001, and a batch size of 32 cells on two Nvidia V100 cards. We
465 tested training on the murine dataset from scratch as well as finetuning the human model. In order to

466 finetune the human model, we gave mouse genes the same initial embedding vector as their human
467 orthologs, where available, and gave unmatched genes random initializations. Orthologs were obtained
468 from Ensembl (**Harrison et al., 2024**), with the first selected if multiple orthologous genes existed. This
469 transfer from human genes significantly improved the model compared to training on mouse genes
470 alone as measured with the cross-entropy loss on 100,000 held-out cells from the mouse brain datasets.
471 We call this finetuned Geneformer pretrained model ‘Murine Geneformer’.

472 Data and tokenization

473 Data

474 The data published in **Chen et al. (2024)** contain spatial transcriptomic data obtained via BARseq. BARseq
475 is a probe-based *in site* sequencing technology based on Illumina chemistry. BARseq establishes mRNA
476 transcript counts in cells based on the identification of probe sequences at spatially resolved locations,
477 followed by cell segmentation and assignment of genes to cells. In this dataset, probes were selected to
478 resolve 104 marker genes chosen to resolve excitatory cells in cortex. After quality control, the overall
479 dataset contains 10.3 million cells across nine brains, including a pilot brain, four brains of mice raised in
480 normal rearing conditions, and four brains of mice raised following monocular enucleation.

481 After downloading spatial transcriptomic data were downloaded from **Chen et al. (2024)**, we first
482 converted all files into AnnData format (**Virshup et al., 2023**) with a custom script. We selected the
483 four animals raised in control conditions, then selected only cortical cells, resulting in over 1.6 million
484 cortical cells across 4 animals. Several metadata annotations for each cell were computed in the original
485 publication. These include cell type annotations at 3 hierarchical levels, and cellular location within each
486 brain slice in ‘slice coordinates’ as well as in CCF coordinates after each brain was registered to the Allen
487 Brain Atlas. This dataset also contains the location within warped ‘flatmap’ coordinates, which use the
488 butterfly projection from the `ccf-streamlines` python package to place cells in a 3D coordinate space
489 in which cortical depth is the third axis (**Wang et al., 2020**).

490 For plots in which cells were colored by cell type (e.g. **Figure 3**), we assigned to each cell type a color
491 so that similar cell types would have similar colors. Specifically, the perceptual distance between any two
492 cell types’ colors was chosen to match the correlation between the cell types’ pseudobulk expression
493 vectors. Unlike standard color maps, this ensures that perceptual saliency aligns with biological meaning,
494 with no salient colors (e.g. red) arbitrarily assigned. Technically, this involves computing the cell type
495 similarity matrix, embedding this matrix into 3D via multidimensional scaling (MDS) to best preserve
496 distances, then interpreting this 3D space as the LUV axes of the perceptually uniform LUV color space.
497 We released a Python package which computes such colormaps called `Co1orMyCe11` to accompany this
498 manuscript (**Benjamin, 2025**).

499 This dataset also contains labels of the brain area of each cell. These are inherited from the CCF
500 coordinate space. We used these annotations as training data when predicting area from cellular
501 transcription. As downloaded, these areas were at the finest level of the hierarchical tree of brain
502 annotations and distinguished between cortical layers. In order to obtain the cortical area without regard
503 to layers, we ascended one level up the hierarchical ‘ancestor’ tree of annotation layers in the Allen Brain
504 Atlas. The resulting areas, displayed in e.g. **Figure 3**, comprise 42 functionally distinct cortical areas.

505 Tokenization

506 All cells’ mRNA transcript vectors were ‘tokenized’, or converted into a format readable by Transformer
507 models, with the following procedure. We used an identical strategy as performed by Geneformer,
508 but wrote custom code to accept data in anndata format (**Virshup et al., 2023**). First, each cell was
509 normalized by the total counts in each cell. All mRNA counts were then normalized by a fixed vector
510 representing the median per-gene counts in a reference data corpus. When fine-tuning Geneformer, it

511 was empirically better to use the same median count vector used to train Geneformer (i.e. on human
512 data) rather than re-establish median counts in mouse data. These normalized counts were then rank
513 sorted, and the indices of the sort (i.e. the output of argsort) was used as the model input. For example,
514 a model input of [gene 5, gene 2, ...] describes that the gene 5 had highest relative transcription, followed
515 by gene 2, and so on. In a modification of the Geneformer pipeline, we ensured genes with zero counts
516 were given a special token `<pad>`. We also modified the Geneformer pipeline to prepend a `<c1s>` token.
517 The 'sentences' seen by the model were thus of the form [`<cls>`, gene 5, gene 2, ..., `<pad>`]. The tokens
518 are mapped to high-dimensional embedding vector via a dictionary learned within TissueFormer.

519 Group construction strategy

520 In order to train and predict with TissueFormer, it is necessary to deliver to it groups of cells which share a
521 label. This requires code infrastructure to appropriately group cells according to the proper criteria. Since
522 single-cell datasets are already quite large, this would ideally occur without requiring the duplication of
523 data for cells that occur in multiple potential groups. Thus, we took an online strategy which loaded the
524 single-cell dataset (after tokenization) and then in an online manner constructed groups to deliver to
525 TissueFormer in batches. To accomplish this we designed a custom sampler and data collator within
526 the Hugging Face code ecosystem for training TissueFormer.

527 Our strategy was to select cells in groups which are approximately cortical columns. In order to cover
528 the cortical surface with a roughly equal density, our pipeline first selected a seed cell at random from
529 the training data, which potentially covers multiple brains. We next extracted the N cells in the same
530 brain which were closest to that cell in the XY plane of the flatmap coordinate system. This resulted in a
531 group of N cells which lie within a column. The width of the column differed in each group depending
532 on the local density of cells. Note that because each brain contains a comparable number of cells, the
533 column size in the test brain was comparable to the column size during training despite there being 3
534 training brains and 3 times the overall number of cells.

535 Predicting area from single cells

536 To provide a baseline for the multi-cell methods, we trained several machine learning models to predict
537 area from single-cell data. Benchmark models either saw cell type (H3 types), raw transcriptomic data, or
538 both in the case of the k-nearest neighbor cell type model. To finetune Murine Geneformer, we removed
539 the last layer and replaced it with a linear classifier of area. We found best performance with a progressive
540 training strategy popular in supervised finetuning. For one epoch, we froze all but the last layer, then each
541 epoch progressively unfroze one Transformer layer in Murine Geneformer. We used a batch size of 32
542 and learning rate 0.0005. Results of training this model from scratch without pretraining are visible in
543 **Figure 2—figure Supplement 1e**.

544 TissueFormer training

545 We trained TissueFormer on an NVIDIA H100 card with a batch size of 4096 total cells divided into N
546 groups, each with $|G| = 4096/N$ cells. A single model trains to completion for this dataset set (1.4 million
547 cells) in 10 epochs in less than an hour. The most efficient training strategy involves saturating the GPU
548 memory, and this is largely driven by the overall number of cells rather than the number of distinct groups.
549 We trained all models with varying group size N to convergence to ensure a fair comparison. However,
550 as convergence can be difficult to verify, we attempted to further match the training resources for each
551 model by comparing models with an equal number of training steps, which is proportional to the number
552 of cell groups seen.

553 Due to the imbalance in the size of brain areas, some classes have many more cells than others. A
554 Bayes optimal model trained on this data would learn an unequal prior over classes, adding a slight bias

555 to classify cell groups as belonging to the largest areas. This may not be desirable in certain applications
556 of brain mapping. For example, cells on the boundary between a large area and a small area which have
557 zero evidence towards one or the other will be classified into the large area due to the prior. This would
558 cause an inflation of large areas. In **Figure 3** we display the results of a model trained with an objective
559 designed to eliminate this bias and enforce an equal prior. Specifically we down-weighted examples
560 from large areas with the factor as computed by the `compute_class_weight` function in scikit-learn
561 (**Pedregosa et al., 2011**). For M total cells, C classes, and M_i cells in class i , this factor is $\frac{M}{C * M_i}$. This
562 ensures that the magnitude of the loss function equally reflects examples from all areas.

563 **Pseudobulk and type composition benchmarks**

564 To ensure a fair comparison with TissueFormer, we constructed train and test sets using the same
565 data loading pipeline with the same random seed. For the pseudobulk results, we then averaged the
566 mRNA transcript counts within each group of cells, then trained our benchmark methods. For the
567 cell type composition models, we extracted the cell type of each cell in the group, then constructed
568 a normalized histogram reflecting the relative composition of cell types in that group. For both data
569 modalities, our benchmark methods were logistic regression and random forests. Specifically, we trained
570 logistic regression with the `lbfgs` solver. We used nested cross-validation to determine the regularization
571 parameter, and found that no regularization was optimal at this dataset size. We also used a random
572 forest classifier with 200 estimators, max depth of 15, 33% of features seen per tree, and a minimum
573 samples per split and 5 and per leaf at 2. We used the scikit-learn implementation of both classifiers
574 (**Pedregosa et al., 2011**).

575 **Brain map visualization**

576 To create **Figure 3a**, we first trained 4 models in a cross-validated scheme, with 1 brain held out each fold.
577 To achieve a high-resolution coverage, we then predicted on the test brain using every cell in the test
578 brain once as a center seed of group selection. The prediction on each cell was stored as its label. These
579 are visualized in **Figure 3—figure Supplement 1**. We then constructed a method to present smooth brain
580 maps with clear boundaries between areas. Namely, we trained a support vector classifier (SVC) with
581 a radial basis function kernel on the cells and their predicted labels, then visualized the predictions on
582 pixels. We used the cuML implementation, a GPU-accelerated method, and used $\gamma = 0.00001$ and $C = 1$
583 as hyperparameters of the SVC (**Raschka et al., 2020**). Intuitively, this method acts to color each pixel
584 in **Figure 3a** by a weighted average of the TissueFormer-predicted labels of the cells surrounding that
585 pixel. The weights are selected due to a combination of proximity to the pixel and a factor determined by
586 the algorithm to minimize the rate at which cells are misclassified, i.e. located on a pixel with a different
587 color (area) than that cell's TissueFormer-predicted area.

588 **Data Availability**

589 All code to train TissueFormer as well as to create the figures in the manuscript is available at <https://github.com/ZadorL>
590 annotation.

591 **Acknowledgements**

592 This project was funded with support from NIH grants 1RF1 MH126883-01A1, 5U19NS123716-04, and
593 7U01NS132363-02. Computational resources were provided with support from NIH grant S100D028632-
594 01.

References

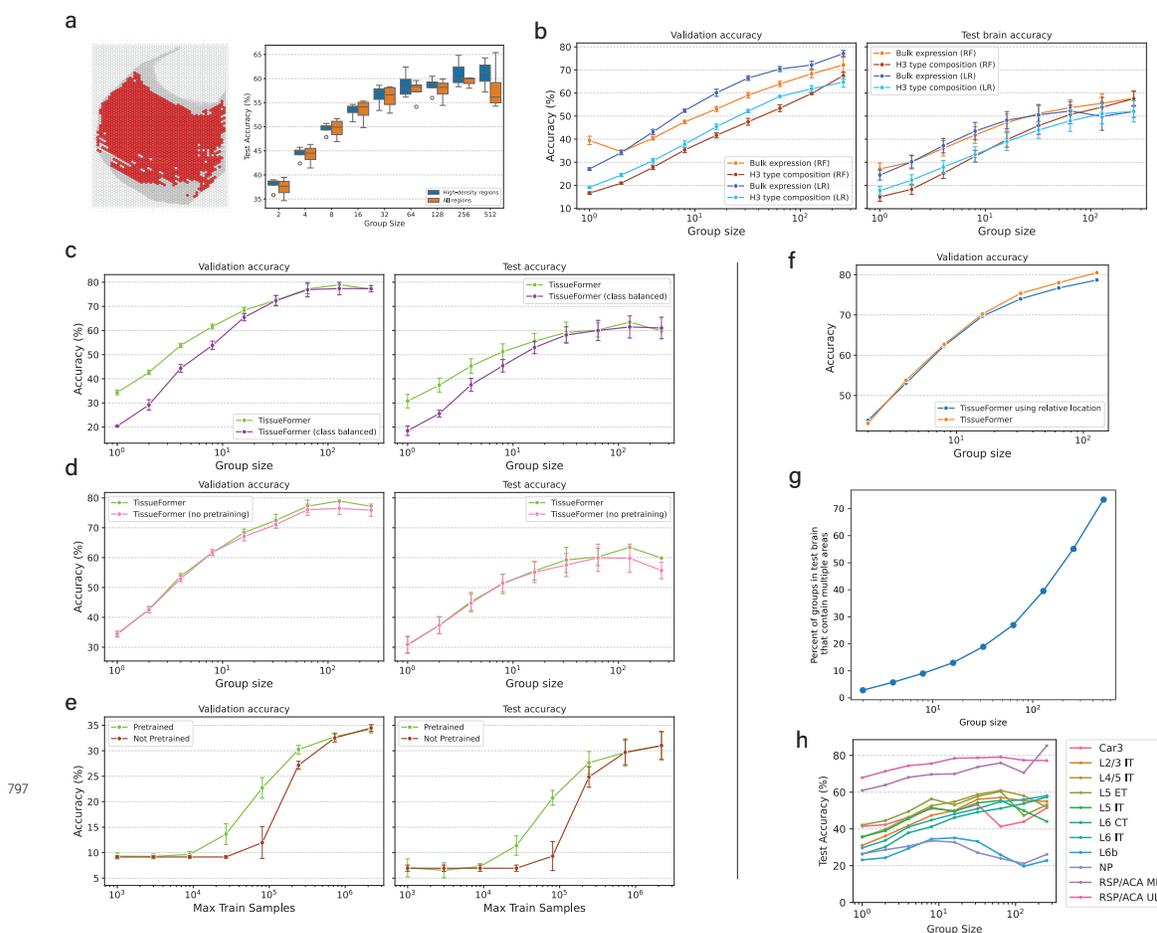
- 595
596 **Arvaniti E**, Claassen M. Sensitive detection of rare disease-associated cell subsets via representation learn-
597 ing. *Nature Communications*. 2017 Apr; 8(1):14825. <https://www.nature.com/articles/ncomms14825>, doi:
598 10.1038/ncomms14825, publisher: Nature Publishing Group.
- 599 **Ballouz S**, Verleyen W, Gillis J. Guidance for RNA-seq co-expression network construction and analysis: safety in
600 numbers. *Bioinformatics*. 2015; 31(13):2123–2130. Publisher: Oxford University Press.
- 601 **Benjamin A**, ColorMyCells: A Python package for biologically faithful colormaps for cell type visualization. Zenodo;
602 2025. <https://doi.org/10.5281/zenodo.15595324>, doi: 10.5281/zenodo.15595324.
- 603 **Biancalani T**, Scalia G, Buffoni L, Avasthi R, Lu Z, Sanger A, Tokcan N, Vanderburg CR, Segerstolpe A, Zhang M,
604 Avraham-Davidi I, Vickovic S, Nitzan M, Ma S, Subramanian A, Lipinski M, Buenrostro J, Brown NB, Fanelli D, Zhuang
605 X, et al. Deep learning and alignment of spatially resolved single-cell transcriptomes with Tangram. *Nature Methods*.
606 2021 Nov; 18(11):1352–1362. <https://www.nature.com/articles/s41592-021-01264-7>, doi: 10.1038/s41592-021-
607 01264-7, publisher: Nature Publishing Group.
- 608 **Brodmann K**. Vergleichende Lokalisationslehre der Grosshirnrinde in ihren Prinzipien dargestellt auf Grund des
609 Zellenbaues. Barth; 1909.
- 610 **Callaway EM**. Local circuits in primary visual cortex of the macaque monkey. *Annual review of neuroscience*. 1998;
611 21(1):47–74. Publisher: Annual Reviews 4139 El Camino Way, PO Box 10139, Palo Alto, CA 94303-0139, USA.
- 612 **Cao Y**, Lin Y, Patrick E, Yang P, Yang JYH. scFeatures: multi-view representations of single-cell and spatial data for
613 disease outcome prediction. *Bioinformatics*. 2022 Oct; 38(20):4745–4753. [https://doi.org/10.1093/bioinformatics/
614 btac590](https://doi.org/10.1093/bioinformatics/btac590), doi: 10.1093/bioinformatics/btac590.
- 615 **Cederquist GY**, Azim E, Shnider SJ, Padmanabhan H, Macklis JD. Lmo4 Establishes Rostral Motor Cortex Projection
616 Neuron Subtype Diversity. *The Journal of Neuroscience*. 2013 Apr; 33(15):6321–6332. [https://www.ncbi.nlm.nih.
617 gov/pmc/articles/PMC3698850/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3698850/), doi: 10.1523/JNEUROSCI.5140-12.2013.
- 618 **Chen SQ**, Chen CH, Xiang XJ, Zhang SY, Ding SL. Chemoarchitecture of area prostriata in adult and developing
619 mice: Comparison with presubiculum and parasubiculum. *The Journal of Comparative Neurology*. 2022 Oct;
620 530(14):2486–2517. doi: 10.1002/cne.25346.
- 621 **Chen X**, Fischer S, Rue MCP, Zhang A, Mukherjee D, Kanold PO, Gillis J, Zador AM. Whole-cortex in situ sequencing
622 reveals input-dependent area identity. *Nature*. 2024 Apr; <https://www.nature.com/articles/s41586-024-07221-6>,
623 doi: 10.1038/s41586-024-07221-6.
- 624 **Chitra U**, Arnold BJ, Sarkar H, Sanno K, Ma C, Lopez-Darwin S, Raphael BJ. Mapping the topography of spatial gene
625 expression with interpretable deep learning. *Nature Methods*. 2025 Feb; 22(2):298–309. [https://doi.org/10.1038/
626 s41592-024-02503-3](https://doi.org/10.1038/s41592-024-02503-3), doi: 10.1038/s41592-024-02503-3.
- 627 **Connell W**, Khan U, Keiser MJ, A single-cell gene expression language model. arXiv; 2022. [http://arxiv.org/abs/2210.
628 14330](http://arxiv.org/abs/2210.14330), doi: 10.48550/arXiv.2210.14330, arXiv:2210.14330 [cs, q-bio] version: 1.
- 629 **Crowell HL**, Soneson C, Germain PL, Calini D, Collin L, Raposo C, Malhotra D, Robinson MD. muscat detects
630 subpopulation-specific state transitions from multi-sample multi-condition single-cell transcriptomics data.
631 *Nature Communications*. 2020 Nov; 11(1):6077. <https://www.nature.com/articles/s41467-020-19894-4>, doi:
632 10.1038/s41467-020-19894-4, publisher: Nature Publishing Group.
- 633 **Cui H**, Wang C, Maan H, Pang K, Luo F, Wang B, scGPT: Towards Building a Foundation Model for Single-Cell
634 Multi-omics Using Generative AI. bioRxiv; 2023. <https://www.biorxiv.org/content/10.1101/2023.04.30.538439v2>,
635 doi: 10.1101/2023.04.30.538439, pages: 2023.04.30.538439 Section: New Results.

- 636 **Dann E**, Cujba AM, Oliver AJ, Meyer KB, Teichmann SA, Marioni JC. Precise identification of cell states altered
637 in disease using healthy single-cell references. *Nature genetics*. 2023; 55(11):1998–2008. Publisher: Nature
638 Publishing Group US New York.
- 639 **Devlin J**, Chang MW, Lee K, Toutanova K, BERT: Pre-training of Deep Bidirectional Transformers for Language
640 Understanding. arXiv; 2019. <http://arxiv.org/abs/1810.04805>, doi: 10.48550/arXiv.1810.04805, arXiv:1810.04805
641 [cs].
- 642 **Dong K**, Zhang S. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph atten-
643 tion auto-encoder. *Nature Communications*. 2022 Apr; 13(1):1739. <https://doi.org/10.1038/s41467-022-29439-6>,
644 doi: 10.1038/s41467-022-29439-6.
- 645 **Dries R**, Zhu Q, Dong R, Eng CHL, Li H, Liu K, Fu Y, Zhao T, Sarkar A, Bao F, George RE, Pierson N, Cai L, Yuan GC.
646 Giotto: a toolbox for integrative analysis and visualization of spatial expression data. *Genome Biology*. 2021 Mar;
647 22(1):78. <https://doi.org/10.1186/s13059-021-02286-2>, doi: 10.1186/s13059-021-02286-2.
- 648 **Govek KW**, Chen S, Sgourdou P, Yao Y, Woodhouse S, Chen T, Fuccillo MV, Epstein DJ, Camara PG. Developmental
649 trajectories of thalamic progenitors revealed by single-cell transcriptome profiling and *Shh* perturbation. *Cell*
650 reports. 2022; 41(10). Publisher: Elsevier.
- 651 **Harrison PW**, Amode MR, Austine-Orimoloye O, Azov AG, Barba M, Barnes I, Becker A, Bennett R, Berry A, Bhai J,
652 others. Ensembl 2024. *Nucleic acids research*. 2024; 52(D1):D891–D899. Publisher: Oxford University Press.
- 653 **Hashikawa Y**, Hashikawa K, Rossi MA, Basiri ML, Liu Y, Johnston NL, Ahmad OR, Stuber GD. Transcriptional and
654 Spatial Resolution of Cell Types in the Mammalian Habenula. *Neuron*. 2020 Jun; 106(5):743–758.e5. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7285796/>, doi: 10.1016/j.neuron.2020.03.011.
655
- 656 **Hsieh KL**, Chu Y, Li X, Pilié PG, Dai Y, scEMB: Learning context representation of genes based on large-scale
657 single-cell transcriptomics. bioRxiv; 2024. <https://www.biorxiv.org/content/10.1101/2024.09.24.614685v1>, doi:
658 10.1101/2024.09.24.614685, pages: 2024.09.24.614685 Section: New Results.
- 659 **Hu J**, Li X, Coleman K, Schroeder A, Ma N, Irwin DJ, Lee EB, Shinohara RT, Li M. SpaGCN: Integrating gene expres-
660 sion, spatial location and histology to identify spatial domains and spatially variable genes by graph convolu-
661 tional network. *Nature Methods*. 2021 Nov; 18(11):1342–1351. <https://doi.org/10.1038/s41592-021-01255-8>, doi:
662 10.1038/s41592-021-01255-8.
- 663 **Ito K**, Hirakawa T, Shigenobu S, Fujiyoshi H, Yamashita T. Mouse-Geneformer: A deep learning model for mouse single-
664 cell transcriptome and its cross-species utility. *PLOS Genetics*. 2025 Mar; 21(3):e1011420. <https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1011420>, doi: 10.1371/journal.pgen.1011420, publisher: Public
665 Library of Science.
- 666
- 667 **Jackson KC**, Boeshaghi AS, Galvez-Merchan A, Moses L, Chari T, Kim A, Pachter L, Identification of spatial homoge-
668 neous regions in tissues with concordex. bioRxiv; 2024. <https://www.biorxiv.org/content/10.1101/2023.06.28.546949v2>, doi: 10.1101/2023.06.28.546949, pages: 2023.06.28.546949 Section: New Results.
669
- 670 **Kalatsky VA**, Stryker MP. New paradigm for optical imaging: temporally encoded maps of intrinsic signal. *Neuron*.
671 2003 May; 38(4):529–545. doi: 10.1016/s0896-6273(03)00286-1.
- 672 **Kozareva V**, Martin C, Osorno T, Rudolph S, Guo C, Vanderburg C, Nadaf N, Regev A, Regehr WG, Macosko E. A tran-
673 scriptomic atlas of mouse cerebellar cortex comprehensively defines cell types. *Nature*. 2021 Oct; 598(7879):214–
674 219. <https://www.nature.com/articles/s41586-021-03220-z>, doi: 10.1038/s41586-021-03220-z, number: 7879
675 Publisher: Nature Publishing Group.
- 676 **Lee AJ**, Dubuc A, Kunst M, Yao S, Lusk N, Ng L, Zeng H, Tasic B, Abbasi-Asl R. Data-driven fine-grained region
677 discovery in the mouse brain with transformers. bioRxiv. 2025; <https://www.biorxiv.org/content/early/2025/02/26/2024.05.05.592608>, doi: 10.1101/2024.05.05.592608, publisher: Cold Spring Harbor Laboratory _eprint:
678 <https://www.biorxiv.org/content/early/2025/02/26/2024.05.05.592608.full.pdf>.
679

- 680 **Lopez R**, Regier J, Cole MB, Jordan MI, Yosef N. Deep generative modeling for single-cell transcriptomics. *Nature Meth-*
681 *ods*. 2018 Dec; 15(12):1053–1058. <https://www.nature.com/articles/s41592-018-0229-2>, doi: 10.1038/s41592-
682 018-0229-2, publisher: Nature Publishing Group.
- 683 **Madhu H**, Rocha JF, Huang T, Viswanath S, Krishnaswamy S, Ying R. HEIST: A Graph Foundation Model for Spatial
684 Transcriptomics and Proteomics Data. *arXiv preprint arXiv:250611152*. 2025; .
- 685 **Mao Y**, Lin YY, Wong NKY, Volik S, Sar F, Collins C, Ester M. Phenotype prediction from single-cell RNA-seq data using
686 attention-based neural networks. *Bioinformatics*. 2024 Feb; 40(2):btæ067. [https://doi.org/10.1093/bioinformatics/
687 btæ067](https://doi.org/10.1093/bioinformatics/btæ067), doi: 10.1093/bioinformatics/btæ067.
- 688 **Mountcastle V**. The columnar organization of the neocortex. *Brain*. 1997 Apr; 120(4):701–722. [https://academic.
689 oup.com/brain/article-lookup/doi/10.1093/brain/120.4.701](https://academic.oup.com/brain/article-lookup/doi/10.1093/brain/120.4.701), doi: 10.1093/brain/120.4.701.
- 690 **Ortiz C**, Navarro JF, Jurek A, Martín A, Lundeberg J, Meletis K. Molecular atlas of the adult mouse brain. *Science*
691 *Advances*. 2020; 6(26):eabb3446. <https://www.science.org/doi/abs/10.1126/sciadv.abb3446>, doi: 10.1126/sci-
692 adv.abb3446, eprint: <https://www.science.org/doi/pdf/10.1126/sciadv.abb3446>.
- 693 **Partel G**, Hilscher MM, Milli G, Solorzano L, Klemm AH, Nilsson M, Wählby C. Automated identification of the
694 mouse brain's spatial compartments from in situ sequencing data. *BMC Biology*. 2020 Dec; 18(1):144. [https:
695 //bmcbiol.biomedcentral.com/articles/10.1186/s12915-020-00874-5](https://bmcbiol.biomedcentral.com/articles/10.1186/s12915-020-00874-5), doi: 10.1186/s12915-020-00874-5.
- 696 **Paszke A**, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison A,
697 Kopf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, et al. PyTorch: An Imperative
698 Style, High-Performance Deep Learning Library. In: Wallach H, Larochelle H, Beygelzimer A, Alché-Buc Fd, Fox
699 E, Garnett R, editors. *Advances in Neural Information Processing Systems*, vol. 32 Curran Associates, Inc.; 2019.
700 https://proceedings.neurips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.
- 701 **Pedregosa F**, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg
702 V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine Learning in
703 Python. *Journal of Machine Learning Research*. 2011; 12:2825–2830.
- 704 **Program CCS**, Abdulla S, Aevermann B, Assis P, Badajoz S, Bell SM, Bezzi E, Cakir B, Chaffer J, Chambers S, others.
705 CZ CELLxGENE Discover: a single-cell data platform for scalable exploration, analysis and modeling of aggregated
706 data. *Nucleic Acids Research*. 2025; 53(D1):D886–D900. Publisher: Oxford University Press.
- 707 **Raschka S**, Patterson J, Nolet C. Machine Learning in Python: Main developments and technology trends in data
708 science, machine learning, and artificial intelligence. *arXiv preprint arXiv:200204803*. 2020; .
- 709 **Rosen Y**, Roohani Y, Agarwal A, Samotorčan L, Tabula Sapiens Consortium, Quake SR, Leskovec J, Universal Cell
710 Embeddings: A Foundation Model for Cell Biology; 2023. <http://biorxiv.org/lookup/doi/10.1101/2023.11.28.568918>,
711 doi: 10.1101/2023.11.28.568918.
- 712 **Schaar AC**, Tejada-Lapuerta A, Palla G, Gutgesell R, Halle L, Minaeva M, Vornholz L, Dony L, Drummer F, Bahrami M,
713 Theis FJ, Nicheformer: a foundation model for single-cell and spatial omics. *bioRxiv*; 2024. [https://www.biorxiv.
714 org/content/10.1101/2024.04.15.589472v1](https://www.biorxiv.org/content/10.1101/2024.04.15.589472v1), doi: 10.1101/2024.04.15.589472, pages: 2024.04.15.589472 Section:
715 New Results.
- 716 **Singhal V**, Chou N, Lee J, Yue Y, Liu J, Chock WK, Lin L, Chang YC, Teo EML, Aow J, Lee HK, Chen KH, Prabhakar S.
717 BANKSY unifies cell typing and tissue domain segmentation for scalable spatial omics data analysis. *Nature*
718 *Genetics*. 2024 Mar; 56(3):431–441. <https://doi.org/10.1038/s41588-024-01664-3>, doi: 10.1038/s41588-024-
719 01664-3.
- 720 **Steuernagel L**, Lam BYH, Klemm P, Dowsett GKC, Bauder CA, Tadross JA, Hitschfeld TS, del Rio Martin A, Chen W,
721 de Solis AJ, Fenselau H, Davidsen P, Cimino I, Kohnke SN, Rimmington D, Coll AP, Beyer A, Yeo GSH, Brüning JC.
722 HypoMap—a unified single-cell gene expression atlas of the murine hypothalamus. *Nature Metabolism*. 2022
723 Oct; 4(10):1402–1419. <https://www.nature.com/articles/s42255-022-00657-y>, doi: 10.1038/s42255-022-00657-y,
724 number: 10 Publisher: Nature Publishing Group.

- 725 **Sun X**, Axelrod ML, Waks AG, Fu J, DiLullo M, Van Allen EM, Tolaney SM, Mittendorf EA, Xu Y, Balko JM. Dynamic
726 single-cell systemic immune responses in immunotherapy-treated early-stage HR+ breast cancer patients. *npj*
727 *Breast Cancer*. 2025 Jul; 11(1):65. <https://doi.org/10.1038/s41523-025-00776-1>, doi: 10.1038/s41523-025-00776-1.
- 728 **Svensson V**, Vento-Tormo R, Teichmann SA. Exponential scaling of single-cell RNA-seq in the past decade. *Nature*
729 *Protocols*. 2018 Apr; 13(4):599–604. <https://www.nature.com/articles/nprot.2017.149>, doi: 10.1038/nprot.2017.149,
730 publisher: Nature Publishing Group.
- 731 **Szalata A**, Hrovatin K, Becker S, Tejada-Lapueta A, Cui H, Wang B, Theis FJ. Transformers in single-cell omics: a
732 review and new perspectives. *Nature Methods*. 2024 Aug; 21(8):1430–1443. <https://www.nature.com/articles/s41592-024-02353-z>, doi: 10.1038/s41592-024-02353-z, publisher: Nature Publishing Group.
- 734 **Tang Z**, Luo S, Zeng H, Huang J, Sui X, Wu M, Wang X. Search and match across spatial omics samples at single-cell
735 resolution. *Nature Methods*. 2024 Oct; 21(10):1818–1829. <https://www.nature.com/articles/s41592-024-02410-7>,
736 doi: 10.1038/s41592-024-02410-7, publisher: Nature Publishing Group.
- 737 **Tasic B**, Yao Z, Graybuck LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J, Garren E, Economo MN, Viswanathan S,
738 others. Shared and distinct transcriptomic cell types across neocortical areas. *Nature*. 2018; 563(7729):72–78.
739 Publisher: Nature Publishing Group UK London.
- 740 **Templeton AJ**, McNamara MG, Šeruga B, Vera-Badillo FE, Aneja P, Ocaña A, Leibowitz-Amit R, Sonpavde G, Knox JJ,
741 Tran B, Tannock IF, Amir E. Prognostic role of neutrophil-to-lymphocyte ratio in solid tumors: a systematic review
742 and meta-analysis. *Journal of the National Cancer Institute*. 2014 Jun; 106(6):dju124. doi: 10.1093/jnci/dju124.
- 743 **Theodoris CV**, Xiao L, Chopra A, Chaffin MD, Al Sayed ZR, Hill MC, Mantineo H, Brydon EM, Zeng Z, Liu XS, Ellinor
744 PT. Transfer learning enables predictions in network biology. *Nature*. 2023 Jun; 618(7965):616–624. <https://www.nature.com/articles/s41586-023-06139-9>, doi: 10.1038/s41586-023-06139-9.
- 746 **Vaswani A**, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin I. Attention is all you need.
747 *Advances in neural information processing systems*. 2017; 30.
- 748 **Verlaan T**, Bouland G, Mahfouz A, Reinders MJT. scAGG: Sample-level embedding and classification of
749 Alzheimer’s disease from single-nucleus data. *bioRxiv*. 2025; <https://www.biorxiv.org/content/early/2025/01/30/2025.01.28.635240>, doi: 10.1101/2025.01.28.635240, publisher: Cold Spring Harbor Laboratory _eprint:
750 <https://www.biorxiv.org/content/early/2025/01/30/2025.01.28.635240.full.pdf>.
- 752 **Virshup I**, Bredikhin D, Heumos L, Palla G, Sturm G, Gayoso A, Kats I, Koutrouli M, Berger B, others. The scverse
753 project provides a computational ecosystem for single-cell omics data analysis. *Nature biotechnology*. 2023;
754 41(5):604–606. Publisher: Nature Publishing Group US New York.
- 755 **Wang C**, Cui H, Zhang A, Xie R, Goodarzi H, Wang B. scGPT-spatial: Continual Pretraining of Single-Cell Foun-
756 dation Model for Spatial Transcriptomics. *bioRxiv*. 2025; <https://www.biorxiv.org/content/early/2025/02/08/2025.02.05.636714>, doi: 10.1101/2025.02.05.636714, publisher: Cold Spring Harbor Laboratory _eprint:
757 <https://www.biorxiv.org/content/early/2025/02/08/2025.02.05.636714.full.pdf>.
- 759 **Wang Q**, Ding SL, Li Y, Royall J, Feng D, Lesnar P, Graddis N, Naemi M, Facer B, Ho A, others. The Allen mouse brain
760 common coordinate framework: a 3D reference atlas. *Cell*. 2020; 181(4):936–953. Publisher: Elsevier.
- 761 **Weinkauff R**, Estey EH, Starostik P, Hayes K, Huh YO, Hirsch-Ginsber C, Andreeff M, Keating M, Kantarjian HM,
762 Freireich EJ, others. Use of peripheral blood blasts vs bone marrow blasts for diagnosis of acute leukemia.
763 *American journal of clinical pathology*. 1999; 111(6):733–740. Publisher: Oxford University Press Oxford, UK.
- 764 **Wolf T**, Debut L, Sanh V, Chaumond J, Delangue C, Moi A, Cistac P, Rault T, Louf R, Funtowicz M, Davison J,
765 Shleifer S, Platen Pv, Ma C, Jernite Y, Plu J, Xu C, Scao TL, Gugger S, Drame M, et al. Transformers: State-of-
766 the-Art Natural Language Processing. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural*
767 *Language Processing: System Demonstrations* Online: Association for Computational Linguistics; 2020. p. 38–45.
768 <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.

- 769 **Xu X**, Holmes TC, Luo MH, Beier KT, Horwitz GD, Zhao F, Zeng W, Hui M, Semler BL, Sandri-Goldin RM.
770 Viral Vectors for Neural Circuit Mapping and Recent Advances in Trans-synaptic Anterograde Tracers.
771 *Neuron*. 2020 Sep; 107(6):1029–1047. [https://www.cell.com/neuron/abstract/S0896-6273\(20\)30527-4](https://www.cell.com/neuron/abstract/S0896-6273(20)30527-4), doi:
772 10.1016/j.neuron.2020.07.010, publisher: Elsevier.
- 773 **Yadan O**, Hydra - A framework for elegantly configuring complex applications; 2019. [https://github.com/](https://github.com/facebookresearch/hydra)
774 [facebookresearch/hydra](https://github.com/facebookresearch/hydra).
- 775 **Yao Z**, Liu H, Xie F, Fischer S, Adkins RS, Aldridge AI, Ament SA, Bartlett A, Behrens MM, Van den Berge K, Bertagnolli
776 D, de Bézieux HR, Biancalani T, Boeshaghi AS, Bravo HC, Casper T, Colantuoni C, Crabtree J, Creasy H, Crichton
777 K, et al. A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex. *Nature*. 2021 Oct;
778 598(7879):103–110. <https://www.nature.com/articles/s41586-021-03500-8>, doi: 10.1038/s41586-021-03500-8,
779 number: 7879 Publisher: Nature Publishing Group.
- 780 **Yao Z**, van Velthoven CTJ, Nguyen TN, Goldy J, Seden-Cortes AE, Baftizadeh F, Bertagnolli D, Casper T, Chiang M,
781 Crichton K, Ding SL, Fong O, Garren E, Glandon A, Gouwens NW, Gray J, Graybuck LT, Hawrylycz MJ, Hirschstein
782 D, Kroll M, et al. A taxonomy of transcriptomic cell types across the isocortex and hippocampal formation.
783 *Cell*. 2021 Jun; 184(12):3222–3241.e26. <https://linkinghub.elsevier.com/retrieve/pii/S0092867421005018>, doi:
784 10.1016/j.cell.2021.04.021.
- 785 **You Y**, Wang ZJ, Fleisher K, Liu R, Thomson M. Building Foundation Models to Characterize Cellular Interactions
786 via Geometric Self-Supervised Learning on Spatial Genomics. *bioRxiv*. 2025; [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/early/2025/01/27/2025.01.25.634867)
787 [early/2025/01/27/2025.01.25.634867](https://www.biorxiv.org/content/early/2025/01/27/2025.01.25.634867), doi: 10.1101/2025.01.25.634867, publisher: Cold Spring Harbor Laboratory
788 _eprint: <https://www.biorxiv.org/content/early/2025/01/27/2025.01.25.634867.full.pdf>.
- 789 **Yuan X**, Zhan Z, Zhang Z, Zhou M, Zhao J, Han B, Li Y, Tang J, Cell-ontology guided transcriptome foundation model.
790 *arXiv*; 2024. <http://arxiv.org/abs/2408.12373>, doi: 10.48550/arXiv.2408.12373, arXiv:2408.12373 [cs].
- 791 **Zhao E**, Stone MR, Ren X, Guenthoer J, Smythe KS, Pulliam T, Williams SR, Uyttingco CR, Taylor SEB, Nghiem P, Bielas
792 JH, Gottardo R. Spatial transcriptomics at subspot resolution with BayesSpace. *Nature Biotechnology*. 2021 Nov;
793 39(11):1375–1384. <https://doi.org/10.1038/s41587-021-00935-2>, doi: 10.1038/s41587-021-00935-2.
- 794 **Zilles K**, Amunts K. Centenary of Brodmann’s map — conception and fate. *Nature Reviews Neuroscience*. 2010
795 Feb; 11(2):139–145. <https://www.nature.com/articles/nrn2776>, doi: 10.1038/nrn2776, publisher: Nature Publishing
796 Group.



797

Figure 2—figure supplement 1. a) The drop in test brain accuracy relative to validation accuracy is only partially due to differences in sampling density. We selected hexes in the test brain which had high density in the 3 training brains (example at left), specifically, hexes with at least as many cells as the group size. (right) Test accuracy on data in those hexes was marginally higher than for all points, yet still much below validation accuracy. **b)** Performance of random forests (RF) compared against logistic regression (LR) given identical inputs. **c)** Performance of a class-balanced model which controls for some areas having more cells in training brains with a modified objective that discounts samples from frequent areas. **d)** Effect of initializing TissueFormer with a pretrained Murine Geneformer as single-cell module vs. randomly initialized. **e)** Varying the number of train samples, comparing random initializations to pretrained initializations, here for $N = 1$ group size. $N = 1$ was chosen to reduce ambiguity and ensure that the number of cells seen equaled the number of groups seen. **f)** Effect of giving TissueFormer the relative spatial location (relative to the center of mass) within each group of each cell. **g)** Number of groups in the test brain that contain an area boundary as a function of the group size. **h)** The accuracy curves of models trained while only ever able to observe groups containing a single cell type. Groups were constructed by selecting a single cell and choosing the nearest N cells of the same time. The high performance of some types, such as RSP/ACA, reflects the fact that these cells are localized in just one area of cortex (retrosplenial/anterior cingulate areas). Slopes of these curves are shown in Fig. 2.

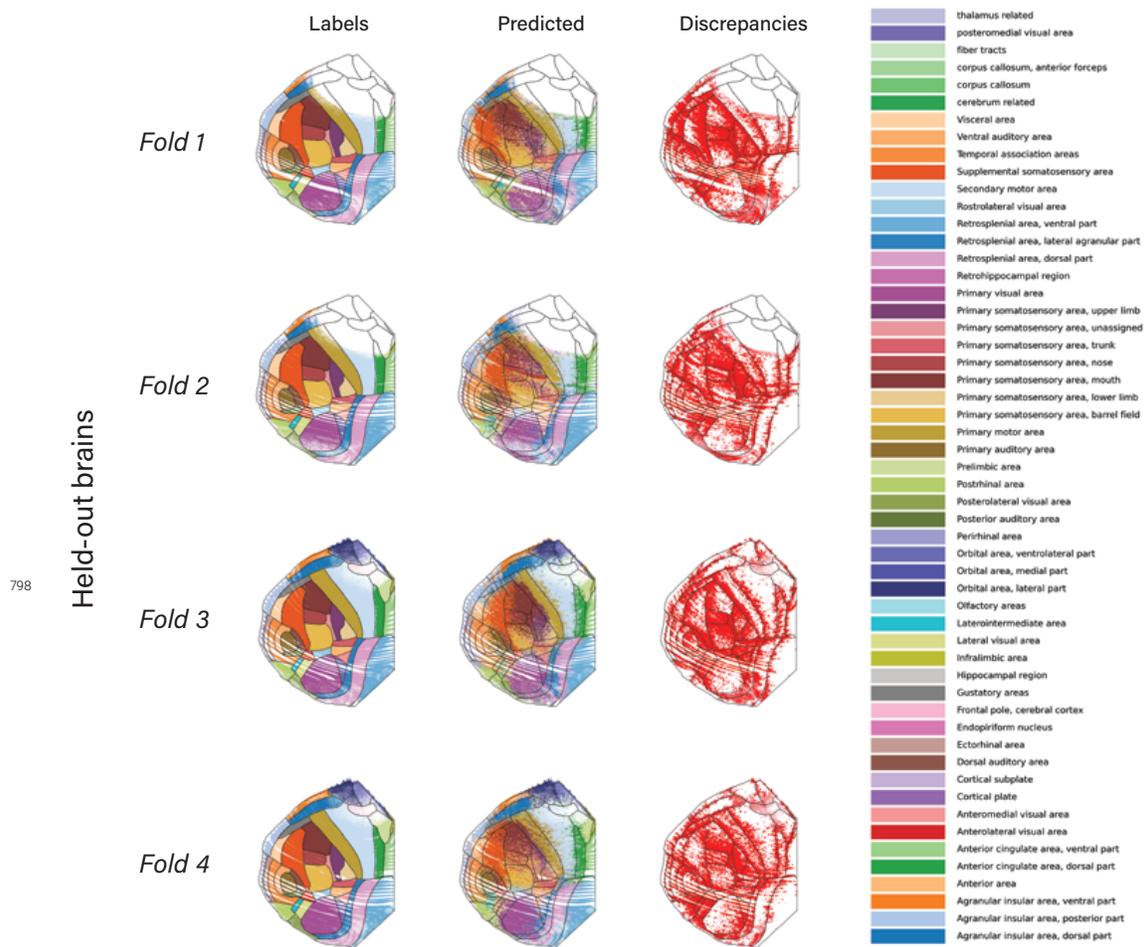


Figure 3—figure supplement 1. Across all 4 brains, we show the CCF labels on single cells (left), the predicted labels (middle) and the location of mismatches between the two. We call these discrepancies rather than errors as it is not clear which is closer to ground truth. Notable systematic differences include a gross shift of the Primary Visual Area (magenta), and the somatosensory-motor boundary. In both cases, a shift appears in the discrepancy map as a large amount of discrepancies on one side of an area but very few on the other side, and a large discontinuity at the CCF boundary.

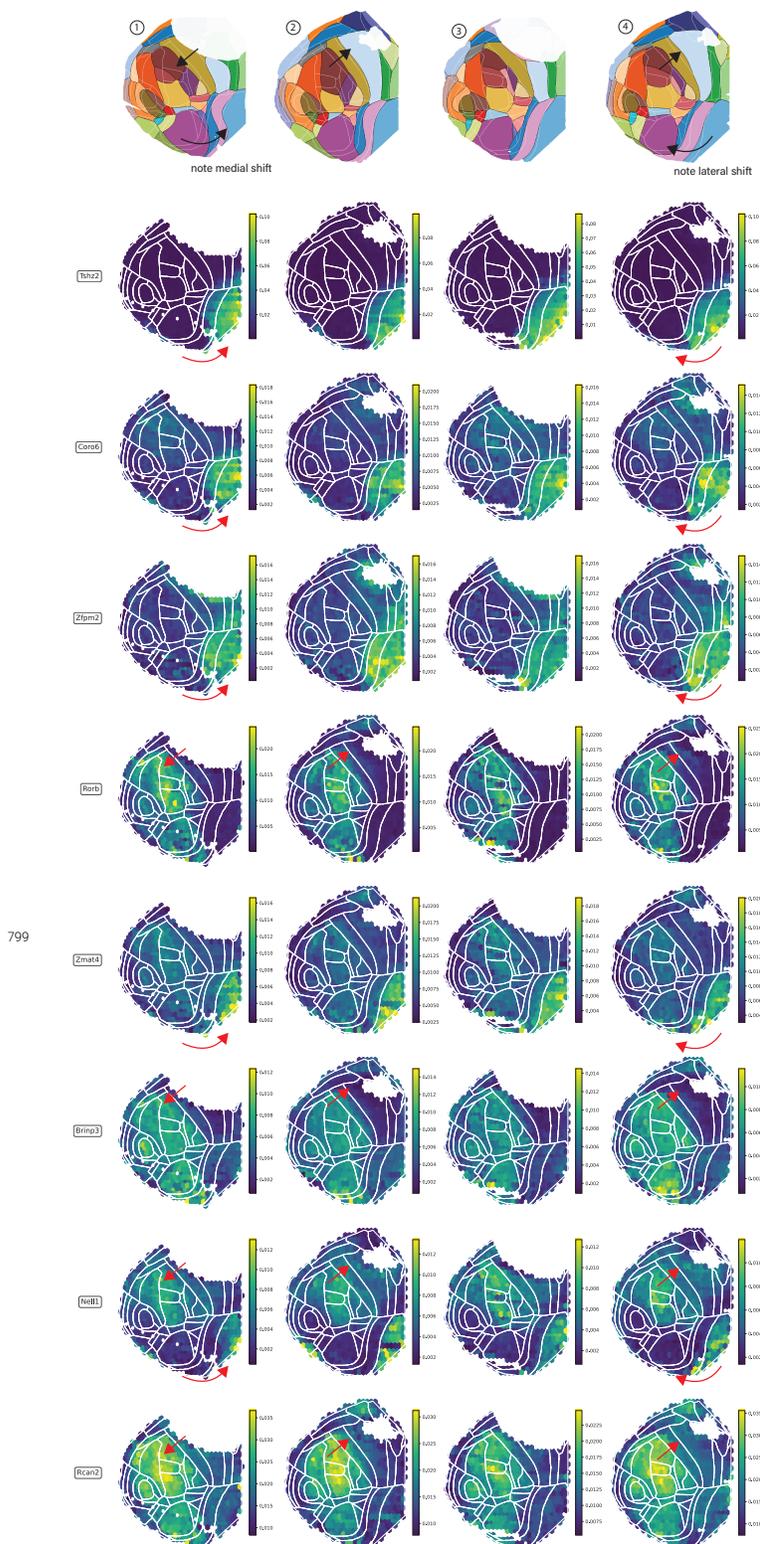


Figure 3—figure supplement 2. Normalized spatial density of select genes in each animal. Arrows highlight the animal-specific shifts of regional anatomy relative to the CCF (white boundaries), and how these are consistent between genes and with the predictions of Tissueformer (top row, reproduced from **Figure 3**). The top three rows (Tshz2, Coro6, and Zfmp2) recapitulate the medial shift in the boundary to the retrosplenial areas in Brain 1 and the lateral shift in this boundary in Brain 4. The bottom 5 rows recapitulate the shift in the somatosensory/motor boundary, the same as analyzed in cell type distributions in Figure 3. To construct normalized gene densities, we first normalized mRNA transcription levels within each cell. We then plot the average normalized transcription within each hex. Normalization was necessary to reduce confounds of unequal spatial sampling.