

Consistent cross-modal identification of cortical neurons with coupled autoencoders

Rohan Gala¹, Agata Budzillo¹, Fahimeh Baftizadeh¹, Jeremy Miller¹, Nathan Gouwens¹, Anton Arkhipov¹, Bosiljka Tasic¹, Gabe Murphy¹, Hongkui Zeng¹, Michael Hawrylycz¹, and Uygur Smbl¹

¹Allen Institute, 615 Westlake Ave N, Seattle WA, USA

Abstract

Consistent identification of neurons and neuronal cell types across different observation modalities is an important problem in neuroscience. Here, we present an optimization framework to learn coordinated representations of multimodal data, and apply it to a large Patch-seq dataset of mouse cortical interneurons. Our approach reveals strong alignment between transcriptomic and electrophysiological profiles of neurons, enables accurate cross-modal data prediction, and identifies cell types that are consistent across modalities.

Keywords

neuronal cell type, Patch-seq, multimodal, cross-modal, coupled autoencoder

Highlights

Coupled autoencoders for multimodal assignment, Analysis of Patch-seq data consisting of more than 3000 cells

20 The characterization of cell types in the brain is an ongoing challenge in contemporary neu-
21 roscience. Describing and analyzing neuronal circuits using cell types can help simplify their
22 complexity and unravel their role in healthy and pathological brain function.^[1,6] However,
23 the effectiveness of such approaches rests on the existence of cellular identities that manifest
24 consistently across different observation modalities, and our ability to identify them. Recent
25 single cell RNA sequencing (scRNA-seq) experiments have provided a detailed window into
26 the transcriptomic organization of cortical cells in the mouse brain.^[7,8] Technological develop-
27 ments have enabled collection of large Patch-seq datasets that include electrophysiological and
28 transcriptomic properties for the same set of neurons.^[9,10] The problem of aligning multimodal
29 data for cell type research is challenging due to complexity of biological relationships between
30 modalities, difficulties in measuring signal and quantifying noise in each modality, and the high
31 dimensional nature of these datasets. Recent works to align single cell -omic measurements
32 have largely focused on removing experimental batch effects, or on estimating correspondences
33 between individual samples across unpaired modalities.^[11,12] For Patchseq-data, there are neither
34 overlapping features nor known associations across the modalities. However the same samples
35 are measured in each modality, and our goal is to formulate consistent cell identities. We
36 present a new deep neural network based methodology referred to as *coupled autoencoders* that
37 addresses the issue of data alignment, and demonstrate its utility for the multimodal cell type
38 identification problem using a Patch-seq dataset with transcriptomic and electrophysiological
39 profiles of 3,411 mouse cortical interneurons.^[9]

40 Coupled autoencoders consist of multiple autoencoder networks, each of which consists of
41 encoder and decoder subnetworks. These subnetworks are nonlinear transformations that
42 project input data into a low dimensional representation, and back to the input data space
43 respectively, Figure [1a](#). In learning these transformations, the goal is to simultaneously maximize
44 reconstruction accuracy for each data modality as well as similarity across representations for
45 the different modalities. In particular, hyper-parameter λ controls the relative importance
46 of achieving accurate reconstructions versus learning representations that are similar across
47 modalities.

48 We find that low-dimensional representations of transcriptomics and electrophysiological mea-
49 surements can be aligned to a high degree, while capturing salient characteristics of neurons
50 in the individual data modalities. This strongly supports the hypothesis that molecular and
51 electrophysiological properties of individual neurons are closely related, reflecting attributes
52 of a common cell type, albeit through a complicated mapping. Importantly, although linear
53 transformations^[13,14] can align the major cell classes, a more detailed alignment of features
54 and cell types is revealed only through non-linear transformations that avoid pathological

55 representations.

56 Using the aligned representations, we show that unsupervised clustering can identify ~ 33 classes
57 of GABAergic interneurons in the mouse visual cortex that are consistent across transcriptomic
58 and electrophysiological characterizations of this neuron population. Additionally, these classes
59 are in agreement with a reference transcriptomic taxonomy of cortical cell types.^[7] Our method
60 is general and can be extended to accommodate additional modalities of interest such as
61 morphology and connectivity, as the datasets mature. We further demonstrate how coupled
62 autoencoders trained on a reference dataset such as the one in this study can serve as a
63 dictionary for smaller, single modality datasets to accurately identify cell types as well as
64 predict data for unobserved modalities.

65 Aligned 3-d representations z_t and z_e for the transcriptomic and electrophysiological profiles for
66 the high-dimensional observation vectors X_t and X_e obtained with coupled autoencoders are
67 shown in Figure 1b-c. Cells labeled according to the reference taxonomy (see Figure S1) cluster
68 together in representations of both observation modalities. Moreover, the representations
69 largely preserve hierarchical relationships between cell types of the reference taxonomy. For
70 example, in Figure 1b-c various cell types of the Sst class appear close together, while remaining
71 well-separated from cell types of other classes such as Pvalb, Vip, and Lamp5.

72 Representations obtained with coupled autoencoders may be used to perform a variety of
73 downstream analyses on complex datasets. We considered supervised classification accuracy
74 in predicting cell type labels at different resolutions (Methods) of the reference taxonomy
75 from z_t and z_e in Figures 1d-e, and data reconstruction performance in Figure 1f. First, we
76 orient the reader with results for the uncoupled setting ($\lambda_{te} = 0.0$) at each of these tasks. In
77 Figure 1d, we note that the representations based on the transcriptomic data alone are best
78 suited for supervised cell type classification using QDA, leading to $>70\%$ accuracy for leaf
79 node cell type labels. This is not surprising, since the reference taxonomy was derived from
80 analyses of gene expression alone. Electrophysiological profiles are expected to be noisy, and of
81 lower resolution compared to transcriptomic profiles.^[15] Nevertheless in Figure 1e, classifiers
82 based on representations of electrophysiology alone predict leaf node cell type labels with
83 $\sim 30\%$ accuracy (chance level is $\sim 3\%$). Lastly, the within-modality reconstruction accuracy
84 of uncoupled representations in Figure 1f provides an upper limit for both, within- and cross-
85 modal reconstructions that may be achieved with 3-d representations obtained with coupled
86 autoencoders.

87 To evaluate whether complicated, non-linear transformations underlie the relationship between
88 the transcriptomic and electrophysiological features of neurons, we considered the performance

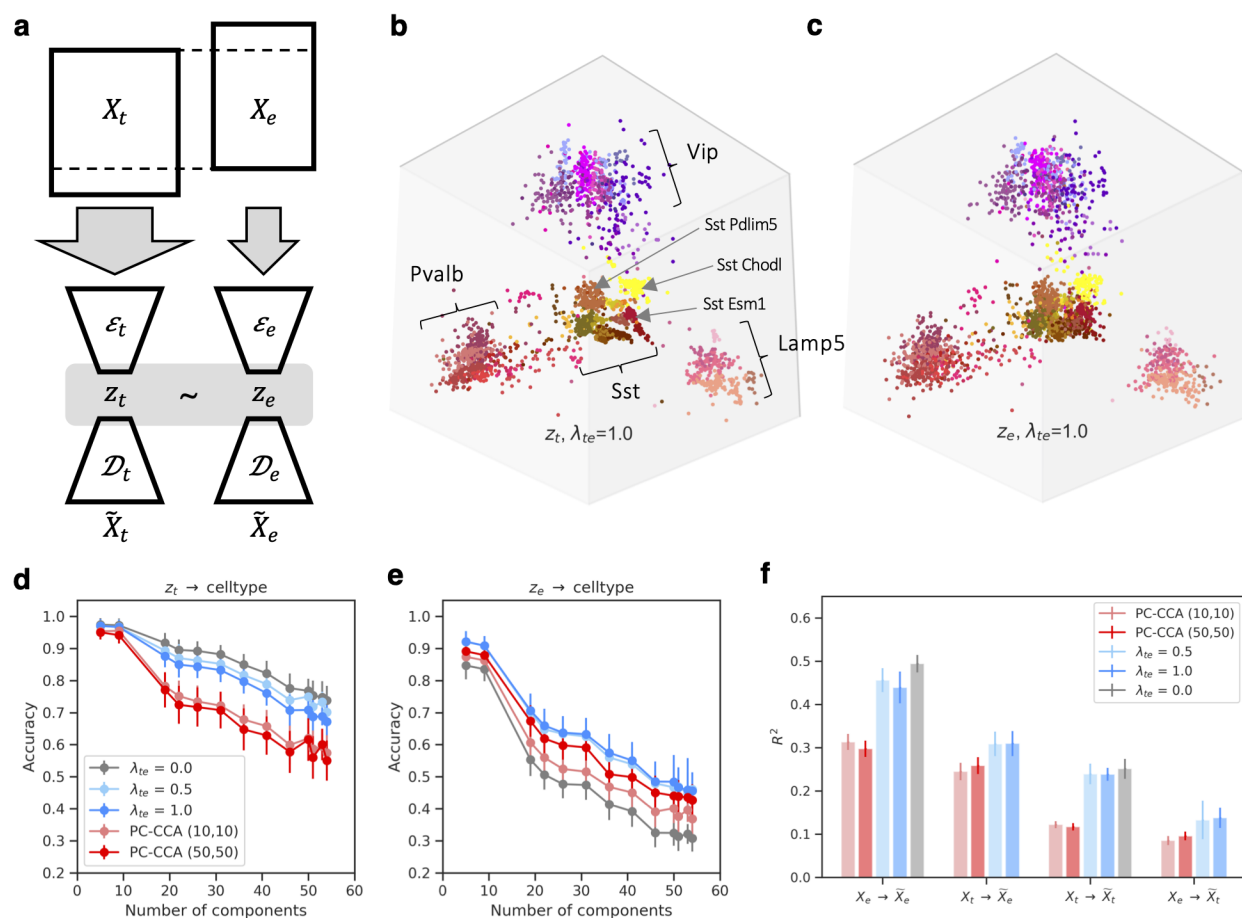


Figure 1: Coordinated representations of transcriptomic and electrophysiological profiles with coupled autoencoders (a) Schematic showing the coupled autoencoder architecture for Patch-seq data. Encoders (\mathcal{E}) compress input data (X) into low dimensional representations (z). Decoders (\mathcal{D}) reconstruct data (\tilde{X}) from representations. The coupling penalty in the loss function encourages representations to be similar across the transcriptomic (t) and electrophysiology (e) modalities. (b-c) 3-d coordinated representations of the transcriptomic and electrophysiological datasets. Each point represents a single cell, colored by the reference hierarchy leaf node to which the cell was mapped to. (d-e) Supervised cell type classification with QDA at different resolutions of the reference hierarchy that were based on 3-d representations obtained with coupled autoencoders and with linear methods. (f) Reconstruction performance as measured with coefficient of determination in the within-modality ($X_e \rightarrow \tilde{X}_e$ and $X_t \rightarrow \tilde{X}_t$), and in the cross-modality ($X_t \rightarrow \tilde{X}_e$ and $X_e \rightarrow \tilde{X}_t$) cases. Error bars show (mean \pm SD, 43-fold cross validation) for panels (d-f)

89 of linear methods (PC-CCA), and coupled autoencoders with $\lambda_{te} \in \{0.5, 1.0\}$ at these tasks,
90 with the representation dimensionality set to 3. We note that the Patch-seq experiment provides
91 perfect knowledge of *anchors* between the modalities by virtue of paired recordings. In this
92 setting, the popular tool Seurat¹⁶ uses a variant of linear CCA to achieve alignment, for
93 which the performance is expected to be comparable to baselines considered here. Results in
94 Figure 1d-f show that coupled autoencoders learn well-aligned representations of transcriptomic
95 and electrophysiology data, such that cell type labels can be predicted with better accuracy, and
96 the cross-modal data can be inferred more reliably compared to linear methods. Importantly,
97 the within-modality reconstruction error is comparable to that obtained in the uncoupled
98 setting, suggesting that the representations compress the individual data modalities with high
99 fidelity.

100 Cross-modal data prediction is a key computational tool for identifying corresponding properties
101 of cell types, and in the design of new experiments. Non-linear transformations to align single
102 cell modalities directly in the data domain have been explored before,¹⁷ but crucially did not
103 provide low dimensional co-ordinated representations. We considered a subset of genes that
104 underlie recently discovered cell type specific paracrine signaling pathways in the cortex.¹⁸ The
105 Patch-seq transcriptomic data shows these cell type specific gene expression patterns, Figure 2a.
106 We used only electrophysiology features to infer the expression patterns for all genes in the
107 cross-modal setting, and show results for the same subset of genes as before in Figure 2b. The
108 striking similarity of these expression patterns (Pearson's $r=0.89\pm 0.10$, mean \pm SD over cell
109 types) not only demonstrates the effectiveness of coupled autoencoders at the cross-modal
110 prediction task at a granular level, but also suggests that intrinsic electrophysiology contains
111 information regarding neuropeptide communication networks.

112 We considered cross-modal prediction of electrophysiological features in an analogous manner,
113 pooling values of the features on a per cell type basis. We considered electrophysiological
114 features that are captured by the compressed representation well (within-modality reconstruction
115 $R^2 > 0.25$, Figure S6). While results of Figure 1d-e already suggest that the electrophysiology
116 features are not as specific to transcriptomic cell types, we can nevertheless identify cell type
117 specific patterns, Figure 2c. The cross-modal reconstruction of these features also matches
118 the data (Pearson's $r=0.99\pm 0.01$, mean \pm SD over cell types), reinforcing the idea that gene
119 expression can explain many intrinsic electrophysiological features accurately, and that coupled
120 autoencoders are a powerful starting point to unravel such non-linear relationships.

121 We directly tested the idea that pre-trained coupled autoencoders can be used to predict
122 unobserved cross-modal features in smaller independent experiments by using the Patch-seq

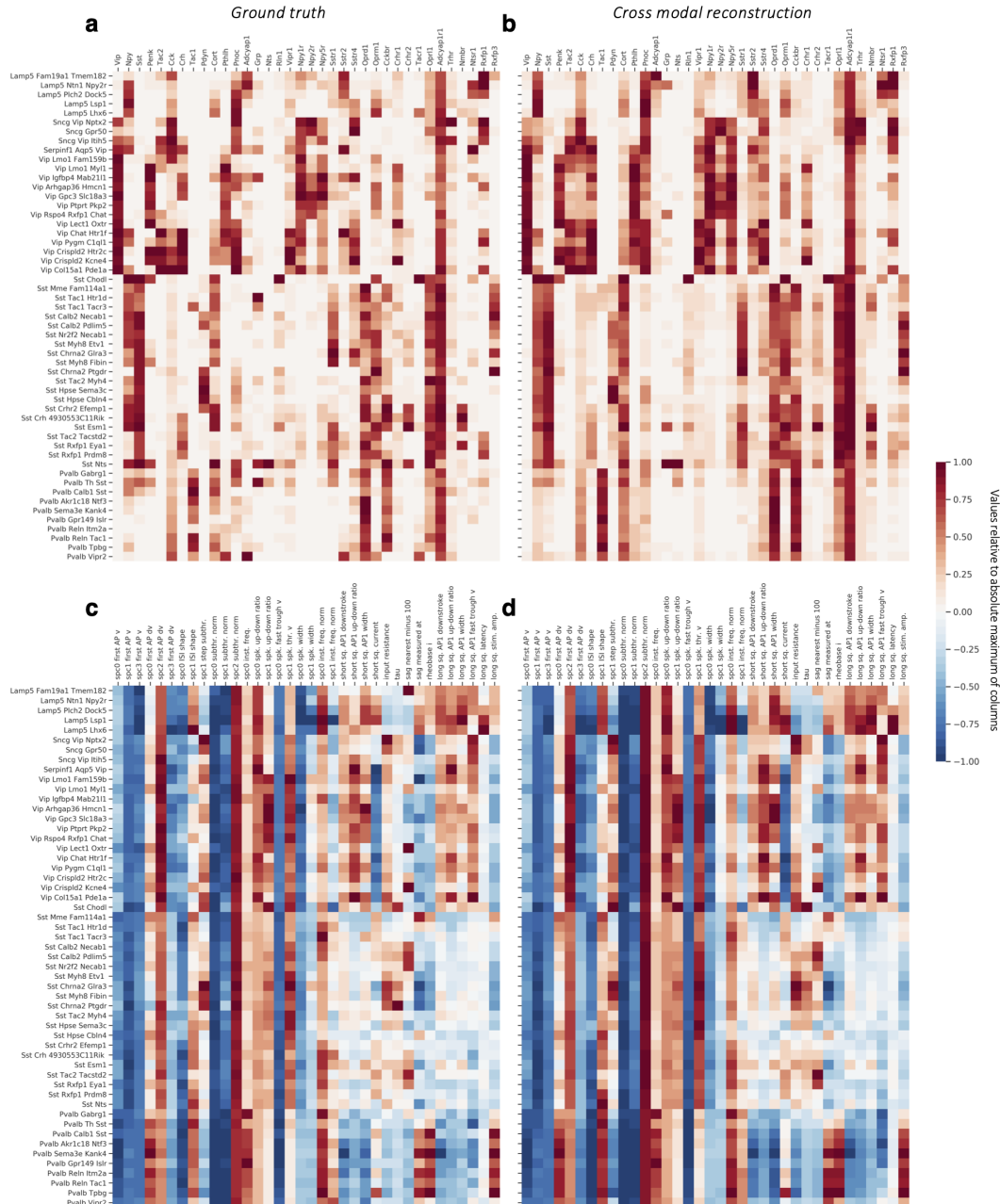


Figure 2: Cross-modal reconstructions capture cell type specific gene expression patterns and electrophysiological features. (a) Gene expression levels averaged over samples of individual cell types of the reference taxonomy, normalized per gene by the maximum value of each column. (b) Cell type specificity of different genes is captured well by cross-modal prediction of gene expression profiles from electrophysiological features. (c) A subset of electrophysiological features pooled by cell types shows analogous cell type specificity. (d) Cross-modal reconstructions of the electrophysiology features from gene expression profiles match the measured electrophysiology features.

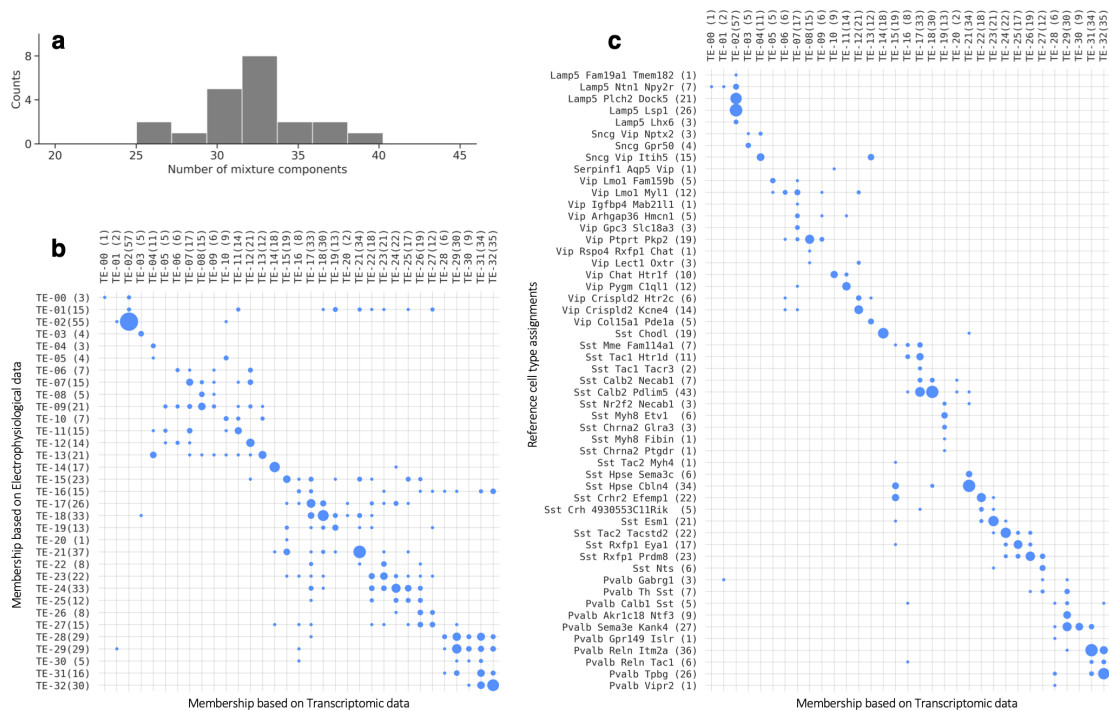


Figure 3: A proposal for consensus clusters (a) Unsupervised clustering using Gaussian Mixtures on the coordinated representation z_t and BIC based model selection suggests ~ 33 consensus clusters. (b) Contingency matrix for cluster assignments based on independent, unsupervised clustering of the transcriptomic and electrophysiology representations shows that the clusters are highly consistent. (c) Contingency matrix for the leaf node cell type labels of the reference hierarchy compared to unsupervised cluster assignments show that these unsupervised clusters have substantial overlap with known transcriptomic cell types. Number of test cells for each label are indicated within parentheses next to the label, and area of the dots is proportional to the number cells in panels (b) and (c).

123 dataset of Scala *et al.*,^[19] which includes 107 inhibitory neurons from mouse motor cortex.
124 We applied a coupled autoencoder without additional training to predict the transcriptomic
125 labels and electrophysiological properties of the 107 neurons from their transcriptomic profiles.
126 Results in Figures S8 and S7 show that this approach yields accurate prediction of cell type
127 labels and certain electrophysiological properties, despite $\sim 5\%$ mismatch between the gene lists
128 and significant differences in electrophysiology protocols.

129 While clustering of individual modalities into cell type taxonomies shows general correspondence,
130 a strategy for consensus clustering is less clear. The notion of a consensus set of cell types can be
131 formalized as a statistical mixture model. Accordingly, the observation for each cell is explained
132 by a combination of its membership to one of a discrete number of types, and continuous
133 variability around the type representative. Encouraged by the clustering of cells belonging

134 to similar transcriptomic types in Figure 1b-c, we explored the extent to which such a model
135 can explain the data consistently across modalities. Specifically, we performed unsupervised
136 clustering by fitting a Gaussian mixture model on coordinated representations obtained with
137 the coupled autoencoder to explain both modalities. Figure 3a shows the distribution of optimal
138 number of mixture components over representations obtained with different coupled autoencoder
139 initializations. This plot suggests that the number of clusters that can be consistently defined
140 with coordinated representations has a tight distribution around ~ 33 . We refer to this *de*
141 *novo* clustering of the data as consensus clusters. Figure 3b demonstrates that the same
142 consensus cluster can be assigned to neurons not used during training with high frequency,
143 based on observing either the transcriptomic or electrophysiological (but not both) modality.
144 While the dominant diagonal of this contingency matrix indicates the success of this notion of
145 consistent, multimodal cortical cell types, the off-diagonal entries point to imperfections of this
146 view, either due to experimental noise and limitations of experimental characterization, or due
147 to imperfection of the model itself.

148 Lastly, the consensus clusters are also consistent with the reference transcriptomic taxonomy,
149 Figure 3c. This might suggest over-splitting in the transcriptomic taxonomy and help identify
150 transcriptomic “super-clusters” of GABAergic neurons, as well as point towards the limitations
151 of the dataset, such as having too few samples for certain transcriptomic labels (see Figure S2)
152 to support a mixture component.

153 In this study, we have presented a principled way to align multimodal observations of neuronal
154 data and define clusters that are consistent across data modalities. Our analysis of the
155 largest multimodal Patch-seq dataset to date with an unsupervised clustering on coordinated
156 representations reveals ~ 33 clusters that can be defined consistently with transcriptomic and
157 electrophysiological measurements of cortical GABAergic neurons. We demonstrated that
158 coupled autoencoders trained on reference datasets can serve as efficient look up tables for
159 smaller, single modality neuron characterization to not only infer cell types, but also their
160 properties in other modalities. Refining this ability will enable the design of new kinds of
161 experiments.

162 An intriguing and essential issue regarding cell types is whether they should be considered as
163 discrete entities or as a continuum.²⁰ Here, we tested a mixture model view on multimodal
164 data, which allows for types to overlap each other in the representation space so long as the
165 cluster centers are more dominant than the peripheries. With this model, mouse visual cortex
166 interneuron Patch-seq data suggests the existence of ~ 33 clusters, more than the ~ 5 well-known
167 subclasses but less than the > 50 partitions suggested by scRNA-seq data alone.

168 Finally, dataset size plays an important role in all our results. More samples can allow the
169 use of larger representation space dimensionality and improve cross-modal data prediction.
170 Similarly, clustering is ill-defined for cell types with too few samples. Therefore, the number of
171 cortical GABAergic interneuron types is likely to grow, and the number of consensus clusters
172 in Figure 3 more likely represents an under-count of the diversity when the notion of cell types
173 is considered as a mixture model.

174 Methods

175 Coupled autoencoders

176 Approaches to discover and extract relationships in multimodal datasets are discussed in litera-
177 ture as cross-modal retrieval, multimodal alignment, multi-view representation learning.²¹⁻²³
178 Deep learning methods such as DeepCCA^{24,25} and correspondence autoencoders²⁶ are promising
179 approaches to achieve multimodal data alignment, but have had limited success in associating
180 complex neural datasets. Our coupled autoencoder networks are related architectures with key
181 improvements to scaling of representations that are critical for the overall quality of learned
182 representations.²⁷

183 We first describe the general coupled autoencoder framework. Then, we show its application to
184 the Patch-seq dataset. For K observation modalities, we represent the coupled autoencoder by

$$\Phi = (\{(\mathcal{E}_i, \mathcal{D}_i, \alpha_i)\}_{1 \leq i \leq K}, \lambda), \quad (1)$$

185 where \mathcal{E}_i and \mathcal{D}_i denote the encoding and decoding networks for the i -th observation modality,
186 α_i sets the relative importance of the different modalities, and $\lambda \geq 0$ sets the relative impor-
187 tance of representation fidelity within observation modalities versus the alignment of different
188 representations.

189 For a set of paired observations $X = \{(x_{s1}, x_{s2}, \dots, x_{sK}), s \in S\}$, we define the loss due to Φ as

$$L(X; \Phi) = \sum_{s \in S} \left[\sum_{i=1}^K \alpha_i \|x_{si} - \mathcal{D}_i(\mathcal{E}_i(x_{si}))\|_2^2 + \lambda \sum_{\substack{i, j \in K, \\ i < j}} \frac{\|\mathcal{E}_i(x_{si}) - \mathcal{E}_j(x_{sj})\|_2^2}{f_{ij}(X)} \right]. \quad (2)$$

190 That is, each autoencoding agent (Figure 1a) within the coupled architecture processes a
191 separate data modality and optimizes a loss function that consists of penalties for (1) the

192 discrepancies between the actual input X and reconstructed input \tilde{X} (2) mismatches between
193 the representations learned by the different agents. (A slightly more general treatment can be
194 found in Ref.^[27])

195 In Eq. [2](#), the functional form of the denominator f_{ij} that scales the mean squared difference
196 between representations of the same sample based on the different data modalities, is crucial to
197 learn good quality representations. Common choices for f_{ij} lead to pathological solutions, i.e.
198 the latent representations collapses into a zero- or one-dimensional space (see Propositions in
199 Supplementary Methods). To avoid such pathological solutions, we propose using:

$$f_{ij}(X) = \min(\sigma_{\min}^2(Z_i), \sigma_{\min}^2(Z_j)) \quad (3)$$

200 where $\sigma_{\min}(Z_i)$ denotes the minimum singular value of the matrix Z_i , which consists of rows
201 $Z_i(s, :) = z_{si}$ where $z_{si} = \mathcal{E}_i(x_{si})$. In practice, we perform stochastic gradient descent and
202 calculate f_{ij} by its mini-batch approximation. Scaling the coupling loss term in this manner
203 approximates whitening by the full covariance matrix well, and also is practically important
204 when the batch size is small or representation dimensionality is large, regimes where calculating
205 the full covariance matrix would be unreliable and computationally expensive.

206 Application to the Patch-seq dataset

207 We use the fact that the same neurons were profiled with both modalities to obtain aligned,
208 low-dimensional representations of the gene expression profiles and electrophysiological features.
209 In the case of just these two data modalities, transcriptomics (t) and electrophysiology (e),
210 the loss function according to Eq. [2](#) consists of two reconstruction error terms, and a single
211 coupling error term. For a single sample s ,

$$L((x_{st}, x_{se})) = \alpha_t \|x_{st} - \mathcal{D}_t(\mathcal{E}_t(x_{st}))\|_2^2 + \alpha_e \|x_{se} - \mathcal{D}_e(\mathcal{E}_e(x_{se}))\|_2^2 + \lambda_{te} \frac{\|z_{st} - z_{se}\|_2^2}{f_{te}(X)}, \quad (4)$$

212 where $z_{st} = \mathcal{E}_t(x_{st})$ and $z_{se} = \mathcal{E}_e(x_{se})$. Here x_{st} denotes gene expression vector for sample s
213 and x_{se} denotes the concatenated sPC and physiological feature measurement vectors for the
214 same sample. The interplay between the accuracy with which the representations capture the
215 individual data modality, versus how well the representations are aligned is a fundamental trade-
216 off that any attempt to define consistent multimodal cell types must resolve (see Supplementary
217 Material for an equivalent formulation in the probabilistic setting). The hyper-parameters α_t , α_e
218 and λ_{te} explicitly control this trade-off in our formulation (Figure [S3](#)).

219 **Data augmentation**

220 Data augmentation is important to regularize the networks and alleviate overfitting, particularly
221 when the dataset size is small. We mimicked the biological dropout phenomenon²⁸ and used
222 Bernoulli noise (i.e., Dropout²⁹) to augment repeated presentations of the transcriptomic
223 vectors while training. This strategy also renders the network robust to partial mismatches
224 in gene lists, and reduces dependence of the representations and reconstructions on specific
225 marker genes. The individual electrophysiological features have unequal variances, since the
226 total variance in the sPC is normalized on a per-experiment basis. We therefore used additive
227 Gaussian noise with variance proportional to that of the individual features to augment the
228 electrophysiological vectors while training the network. The reconstruction loss for the decoders
229 was calculated with both, the representation obtained by the encoder network of the same
230 modality, and that obtained by the encoder for the other modality. This was done to improve
231 performance of cross-modal prediction. We view this way of calculating the reconstruction loss
232 function as an augmentation strategy for the decoder networks.

233 **Linear baselines**

234 Canonical correlation analysis (CCA) is a standard linear method to align low dimensional
235 representations.¹³ To optimize the performance with linear methods, we first used principle
236 component analysis (PCA) to reduce the dimensionality of individual data modalities, followed
237 by CCA to achieve aligned representations across the modalities. The number of dimensions to
238 which the transcriptomic and electrophysiology data were reduced to with PCA is indicated as
239 a tuple in the legends of Figure **I**. The dimensionality of CCA representations was chosen to
240 match the dimensionality obtained with coupled autoencoders (dim=3). The inverse CCA and
241 PCA transformations were used to reconstruct data from the representations both, for the the
242 within- and across- modality cases in Figure **I**f.

243 **Supervised cell type classification**

244 Label sets obtained at different resolutions of the reference taxonomy were used as ground truth
245 to evaluate representations. The different resolutions correspond to different horizontal levels
246 of the reference taxonomy hierarchy in Figure **SI**. Starting from the leaf node cell type labels,
247 each cell is assigned the parent node label based on the set of labels that remains at a given
248 level of the hierarchy. Quadratic Discriminant Analysis (QDA)¹³ was used to train classifiers on

249 the representations obtained with coupled autoencoders or CCA, and used to predict the cell
250 type labels for all such label sets. Cells that were not used to train the coupled autoencoder
251 were used to obtain accuracy values shown in Figure 1(d-e) using a $k=43$ fold cross validation
252 approach. Validation folds were obtained such the class distribution in each fold was similar to
253 that for the overall dataset. Classes with $n \leq 10$ samples in the dataset were discarded from
254 the analysis. Similarly, classes for which there were less than $n=6$ samples in the training set
255 of any fold were discarded from evaluation for only that fold, since QDA classifier parameters
256 for those poorly represented classes would be unreliable. The results were pooled across the
257 folds for the remaining number of classes (i.e. QDA components) in Figure 1(d-e).

258 Unsupervised clustering and consensus clusters

259 Gaussian mixture models with a different number of components (15 to 45 in steps of 1) were
260 fit on z_t obtained with coupled autoencoders ($\lambda_{te} = 1.0$) for 21 different network initializations
261 trained on the same 80% of the dataset. The remaining 20% of cells serve as the test set
262 for this analysis. The training and test sets had similar distributions of the cell type labels
263 based on the reference taxonomy. Each mixture model fit was initialized 50 times and fit
264 until convergence. For the representation from each network initialization, we used Bayesian
265 Information Criterion^[3] (BIC) to perform model selection. The distribution for optimal number
266 of mixture components across the 21 different representations was binned using the Freedman-
267 Diaconis rule,^[30] Figure 3a. Based on this distribution we estimated the number of clusters that
268 can be consistently defined with coordinated representations to be 33. We picked the model
269 with the lowest reconstruction error, and refer to the mixture model with 33 components fitted
270 on z_t as consensus clusters. The fitted mixture model was then used to assign consensus cluster
271 labels to test cells based on z_t , as well as based on z_e . The consensus cluster assignments
272 obtained in this manner are compared in Figure 3b. We used the Hungarian algorithm to
273 match the consensus clusters with leaf node cell types of the reference taxonomy, using the
274 negative of the contingency matrix based on training cells as the cost function. The order of
275 the consensus clusters in Figure 3b-c reflects this optimal match.

276 Patch-seq dataset

277 We used the transcriptomic and electrophysiological profiles of 3,411 GABAergic interneurons
278 from mouse visual cortex of a recent Patch-seq dataset.^[9] The dataset includes cell type labels
279 that were obtained by mapping the gene expression profiles to a reference taxonomy.^[7] The

280 relevant taxonomy, and abundances of cells per type are shown in Figure S1 and Figure S2.
281 There are 59 cell types at the highest resolution (i.e. leaf nodes) of this reference taxonomy.
282 A set of 1,252 genes after removing genes related to mitochondria and sex were used as
283 input for the analyses in this study. Gene expression values were CPM normalized, and then
284 $\log_e(\bullet + 1)$ transformed. 44 sparse principle components (sPC) were extracted to summarize
285 the time series data from different portions of the electrophysiology measurement protocol.⁶
286 Additionally 24 measurements of intrinsic physiology features were obtained using the IPFX
287 library <https://ipfx.readthedocs.io/>. The sPC values were scaled to have unit variance
288 per experiment. The remaining features were individually normalized to have zero mean and
289 unit norm. Data was divided into $k=43$ folds for cross validation experiments. For the consensus
290 cluster experiments, 20% of the cells were set aside as the test set. Different random seeds were
291 used to train networks 21 times on the remaining 80% of the cells.

292 Code availability

293 Code for the coupled autoencoder implementation and analysis are available at <https://github.com/AllenInstitute/coupledAE-patchseq>. The coupled autoencoder was imple-
294 mented using Tensorflow 2.1. Scikit-learn³¹ version 0.22.2 implementations of PCA, CCA,
295 QDA and Gaussian Mixture Models, and Scipy version 1.4.1 implementation of the Hungarian
296 algorithm (linear sum assignment) were used to perform the analyses.
297

298 Acknowledgements

299 We wish to thank the Allen Institute for Brain Science founder, Paul G Allen, for his vision,
300 encouragement and support.

301 References

302 ¹ Robin Tremblay, Soohyun Lee, and Bernardo Rudy. Gabaergic interneurons in the neocortex:
303 from cellular properties to circuits. *Neuron*, 91(2):260–292, 2016.

304 ² Hongkui Zeng and Joshua R Sanes. Neuronal cell-type classification: challenges, opportunities
305 and the path forward. *Nature Reviews Neuroscience*, 18(9):530, 2017.

- 306 ³ Anirban Paul, Megan Crow, Ricardo Raudales, Miao He, Jesse Gillis, and Z Josh Huang.
307 Transcriptional architecture of synaptic communication delineates gabaergic neuron identity.
308 *Cell*, 171(3):522–539, 2017.
- 309 ⁴ Z Josh Huang and Anirban Paul. The diversity of gabaergic neurons and neural communication
310 elements. *Nature Reviews Neuroscience*, 20(9):563–572, 2019.
- 311 ⁵ Giorgio A Ascoli, Lidia Alonso-Nanclares, Stewart A Anderson, German Barrionuevo, Ruth
312 Benavides-Piccione, Andreas Burkhalter, György Buzsáki, Bruno Cauli, Javier DeFelipe,
313 Alfonso Fairén, et al. Petilla terminology: nomenclature of features of gabaergic interneurons
314 of the cerebral cortex. *Nature Reviews Neuroscience*, 9(7):557, 2008.
- 315 ⁶ Philipp Berens and Thomas Euler. Neuronal diversity in the retina. *e-Neuroforum*, 23(2):93–
316 101, 2017.
- 317 ⁷ Bosiljka Tasic, Zizhen Yao, Lucas T Graybuck, Kimberly A Smith, Thuc Nghi Nguyen, Darren
318 Bertagnolli, Jeff Goldy, Emma Garren, Michael N Economo, Sarada Viswanathan, et al.
319 Shared and distinct transcriptomic cell types across neocortical areas. *Nature*, 563(7729):72–78,
320 2018.
- 321 ⁸ A Zeisel, AB Muñoz-Manchado, S Codeluppi, P Lönnerberg, Manno G La, A Juréus,
322 S Marques, H Munguba, L He, C Betsholtz, C Rolny, G Castelo-Branco, J Hjerling-Leffler,
323 and S Linnarsson. Brain structure. Cell types in the mouse cortex and hippocampus revealed
324 by single-cell RNA-seq. *Science*, 347:1138–42, Mar 2015.
- 325 ⁹ Nathan W Gouwens, Staci A Sorensen, Fahimeh Baftizadeh, Agata Budzillo, Brian R Lee,
326 Tim Jarsky, Lauren Alfiler, Anton Arkhipov, Katherine Baker, Eliza Barkan, Kyla Berry,
327 Darren Bertagnolli, Kris Bickley, Jasmine Bomben, Thomas Braun, Krissy Brouner, Tamara
328 Casper, Kirsten Crichton, Tanya L Daigle, Rachel Dalley, Rebecca de Frates, Nick Dee,
329 Tsega Desta, Samuel D Lee, Nadezhda Dotson, Tom Egdorf, Lauren Ellingwood, Rachel
330 Enstrom, Luke Esposito, Colin Farrell, David Feng, Olivia Fong, Rohan Gala, Clare Gamlin,
331 Amanda Gary, Alexandra Glandon, Jeff Goldy, Melissa Gorham, Lucas Graybuck, Hong Gu,
332 Kristen Hadley, Michael J Hawrylycz, Alex M Henry, DiJon Hill, Madie Hupp, Sara Kebede,
333 Tae Kyung Kim, Lisa Kim, Matthew Kroll, Changkyu Lee, Katherine E Link, Matthew
334 Mallory, Rusty Mann, Michelle Maxwell, Medea McGraw, Delissa McMillen, Alice Mukora,
335 Lindsay Ng, Lydia Ng, Kiet Ngo, Philip R Nicovich, Aaron Oldre, Daniel Park, Hanchuan
336 Peng, Osnat Penn, Thanh Pham, Alice Pom, Lydia Potekhina, Ramkumar Rajanbabu,
337 Shea Ransford, David Reid, Christine Rimorin, Miranda Robertson, Kara Ronellenfitch,
338 Augustin Ruiz, David Sandman, Kimberly Smith, Josef Sulc, Susan M Sunkin, Aaron Szafer,
339 Michael Tieu, Amy Torkelson, Jessica Trinh, Herman Tung, Wayne Wakeman, Katelyn Ward,

- 340 Grace Williams, Zhi Zhou, Jonathan Ting, Uygur Sumbul, Ed Lein, Christof Koch, Zizhen
341 Yao, Bosiljka Tasic, Jim Berg, Gabe J Murphy, and Hongkui Zeng. Toward an integrated
342 classification of neuronal cell types: morphoelectric and transcriptomic characterization of
343 individual GABAergic cortical neurons. *bioRxiv*, 2020.
- 344 ¹⁰ Federico Scala, Dmitry Kobak, Matteo Bernabucci, Yves Bernaerts, Cathryn R Cadwell,
345 Jesus R Castro, Leonard Hartmanis, Xiaolong Jiang, Sophie R Laturus, Elanine Miranda,
346 et al. Phenotypic variation within and across transcriptomic cell types in mouse motor cortex.
347 *bioRxiv*, 2020.
- 348 ¹¹ Tim Stuart and Rahul Satija. Integrative single-cell analysis. *Nature Reviews Genetics*,
349 20(5):257–272, 2019.
- 350 ¹² Maria Colomé-Tatché and Fabian J Theis. Statistical single cell multi-omics integration.
351 *Current Opinion in Systems Biology*, 7:54–59, 2018.
- 352 ¹³ Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning:*
353 *data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- 354 ¹⁴ Dmitry Kobak, Yves Bernaerts, Marissa A Weis, Federico Scala, Andreas Tolias, and Philipp
355 Berens. Sparse reduced-rank regression for exploratory visualization of multimodal data sets.
356 *bioRxiv*, page 302208, 2019.
- 357 ¹⁵ NW Gouwens, SA Sorensen, J Berg, C Lee, T Jarsky, J Ting, SM Sunkin, D Feng, CA Anas-
358 tassiou, E Barkan, K Bickley, N Blesie, T Braun, K Brouner, A Budzillo, S Caldejon,
359 T Casper, D Castelli, P Chong, K Crichton, C Cuhaciyani, TL Daigle, R Dalley, N Dee,
360 T Desta, SL Ding, S Dingman, A Doperalski, N Dotson, T Egdorf, M Fisher, Frates RA
361 de, E Garren, M Garwood, A Gary, N Gaudreault, K Godfrey, M Gorham, H Gu, C Habel,
362 K Hadley, J Harrington, JA Harris, A Henry, D Hill, S Josephsen, S Kebede, L Kim, M Kroll,
363 B Lee, T Lemon, KE Link, X Liu, B Long, R Mann, M McGraw, S Mihalas, A Mukora,
364 GJ Murphy, L Ng, K Ngo, TN Nguyen, PR Nicovich, A Oldre, D Park, S Parry, J Perkins,
365 L Potekhina, D Reid, M Robertson, D Sandman, M Schroedter, C Slaughterbeck, G Soler-
366 Llavina, J Sulc, A Szafer, B Tasic, N Taskin, C Teeter, N Thatra, H Tung, W Wakeman,
367 G Williams, R Young, Z Zhou, C Farrell, H Peng, MJ Hawrylycz, E Lein, L Ng, A Arkhipov,
368 A Bernard, JW Phillips, H Zeng, and C Koch. Classification of electrophysiological and
369 morphological neuron types in the mouse visual cortex. *Nat Neurosci*, 22:1182–1195, Jul
370 2019.
- 371 ¹⁶ Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi,
372 William M Mauck III, Yuhao Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija.
373 Comprehensive integration of single-cell data. *Cell*, 177(7):1888–1902, 2019.

- 374 ¹⁷ Matthew Amodio and Smita Krishnaswamy. Magan: Aligning biological manifolds. *arXiv*
375 *preprint arXiv:1803.00385*, 2018.
- 376 ¹⁸ Stephen J Smith, Uygur Sümbül, Lucas T Graybuck, Forrest Collman, Sharmishta Seshamani,
377 Rohan Gala, Olga Gliko, Leila Elabbady, Jeremy A Miller, Trygve E Bakken, et al. Single-cell
378 transcriptomic evidence for dense intracortical neuropeptide networks. *Elife*, 8:e47889, 2019.
- 379 ¹⁹ F Scala, D Kobak, S Shan, Y Bernaerts, S Laturus, CR Cadwell, L Hartmanis, E Froudarakis,
380 JR Castro, ZH Tan, S Papadopoulos, SS Patel, R Sandberg, P Berens, X Jiang, and AS Tolias.
381 Layer 4 of mouse neocortex differs in cell types and circuit organization between sensory
382 areas. *Nat Commun*, 10:4174, Sep 2019.
- 383 ²⁰ Kenneth D Harris, Hannah Hochgerner, Nathan G Skene, Lorenza Magno, Linda Katona,
384 Carolina Bengtsson Gonzales, Peter Somogyi, Nicoletta Kessaris, Sten Linnarsson, and Jens
385 Hjerling-Leffler. Classes and continua of hippocampal ca1 inhibitory neurons revealed by
386 single-cell transcriptomics. *PLoS biology*, 16(6):e2006387, 2018.
- 387 ²¹ Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine
388 learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine*
389 *intelligence*, 41(2):423–443, 2018.
- 390 ²² Yingming Li, Ming Yang, and Zhongfei Zhang. A survey of multi-view representation learning.
391 *IEEE transactions on knowledge and data engineering*, 31(10):1863–1883, 2018.
- 392 ²³ Kaiye Wang, Qiyue Yin, Wei Wang, Shu Wu, and Liang Wang. A Comprehensive Survey on
393 Cross-modal Retrieval, 2016.
- 394 ²⁴ Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation
395 analysis. In *International conference on machine learning*, pages 1247–1255, 2013.
- 396 ²⁵ Weiran Wang, Raman Arora, Karen Livescu, and Jeff Bilmes. On deep multi-view repre-
397 sentation learning. In *International Conference on Machine Learning*, pages 1083–1092,
398 2015.
- 399 ²⁶ Fangxiang Feng, Xiaojie Wang, and Ruifan Li. Cross-modal retrieval with correspondence
400 autoencoder. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages
401 7–16, 2014.
- 402 ²⁷ Rohan Gala, Nathan Gouwens, Zizhen Yao, Agata Budzillo, Osnat Penn, Bosiljka Tasic,
403 Gabe Murphy, Hongkui Zeng, and Uygur Sümbül. A coupled autoencoder approach for
404 multi-modal analysis of cell types. In *Advances in Neural Information Processing Systems*,
405 pages 9263–9272, 2019.

- 406 ²⁸ Peter V Kharchenko, Lev Silberstein, and David T Scadden. Bayesian approach to single-cell
407 differential expression analysis. *Nature methods*, 11(7):740, 2014.
- 408 ²⁹ Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
409 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*
410 *learning research*, 15(1):1929–1958, 2014.
- 411 ³⁰ David Freedman and Persi Diaconis. On the histogram as a density estimator:L 2 theory.
412 *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 57(4):453–476, dec 1981.
- 413 ³¹ F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
414 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
415 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
416 *Learning Research*, 12:2825–2830, 2011.
- 417 ³² Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training
418 by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.