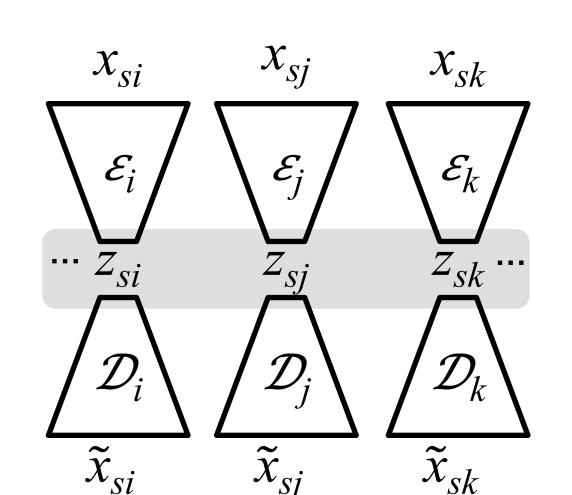


Rohan Gala, Nathan Gouwens, Zizhen Yao, Agata Budzillo, Osnat Penn, Bosiljka Tasic, Gabe Murphy, Hongkui Zeng, Uygar Sümbül, Allen Institute for Brain Science.

#### 1. Abstract

Recent developments in high throughput profiling of individual neurons have spurred data driven exploration of the idea that there exist natural groupings of neurons referred to as cell types. The promise of this idea is that the immense complexity of brain circuits can be reduced, and effectively studied by means of interactions between cell types. While clustering of neuron populations based on a particular data modality can be used to define cell types, such definitions are often inconsistent across different characterization modalities. We pose this issue of cross-modal alignment as an optimization problem and develop an approach based on coupled training of autoencoders as a framework for such analyses. We apply this framework to a Patch-seq dataset consisting of transcriptomic and electrophysiological profiles for the same set of neurons to study consistency of representations across modalities, and evaluate cross-modal data prediction ability. We explore the problem where only a subset of neurons is characterized with more than one modality, and demonstrate that representations learned by coupled autoencoders can be used to identify types sampled only by a single modality.

## 2. Coupled autoencoders to learn consistent representations



 $L = \sum_{i} \alpha_{i} R_{i} + \sum_{i} \lambda_{ij} C_{ij}$ Reconstruction cost

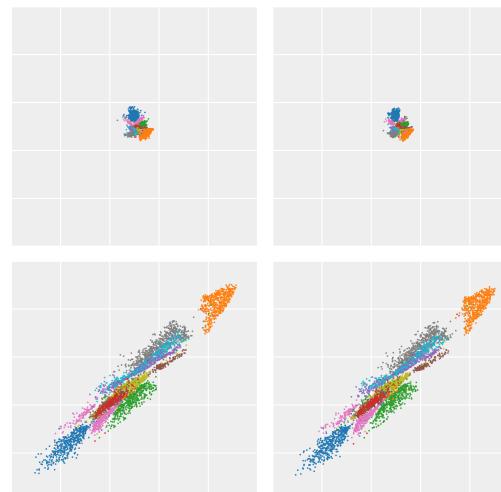
 $R_i = \left\| X_i - \widetilde{X}_i \right\|^2$ 

Coupling cost function  $C_{ij} \sim \|z_i - z_i\|^2$ 

 $\triangleright \alpha_i$  is fixed based on dataset specific noise considerations.

> $\rightarrow \lambda_i$  is controls the tradeoff between reconstruction and coupling costs.

### Pathological representations



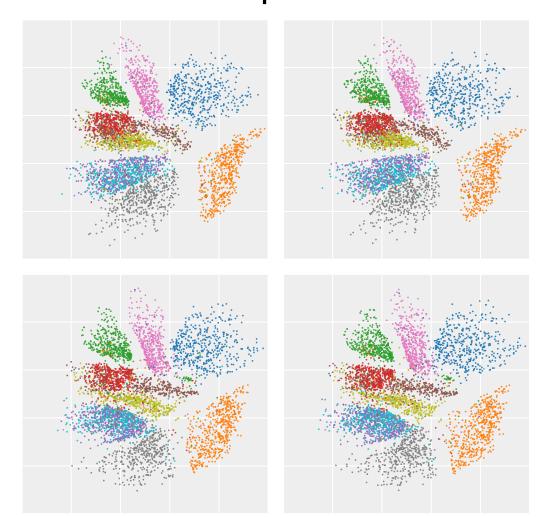
Batch normalized z $z_{id}$   $\longrightarrow$   $rac{z_{id}$  -  $\mu_{\scriptscriptstyle B}$ 

Full covariance normalized z

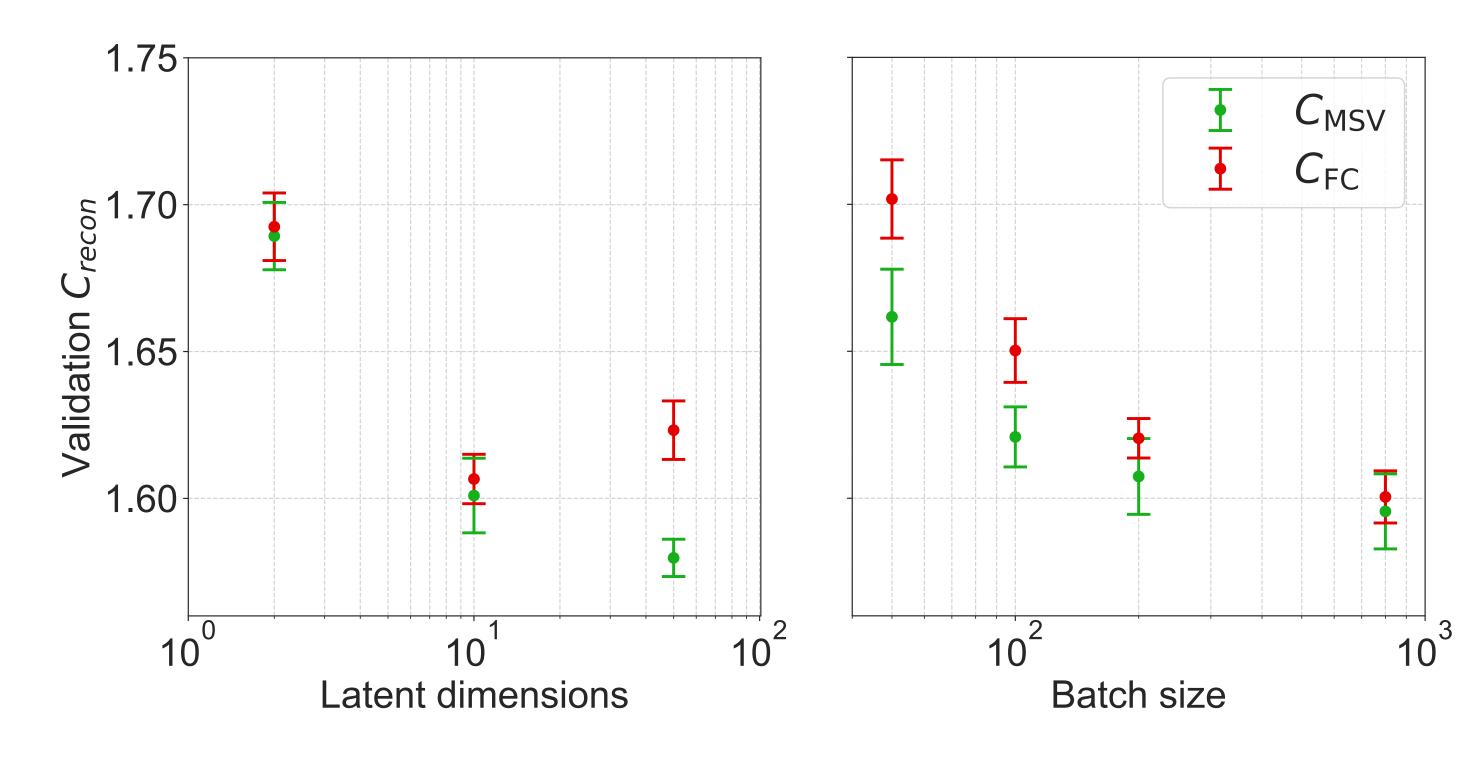
Minimum batch singular value normalized z

 $z_i \longrightarrow (B_i^T B_i)^{-1/2} z_i$ 

## Desirable representations



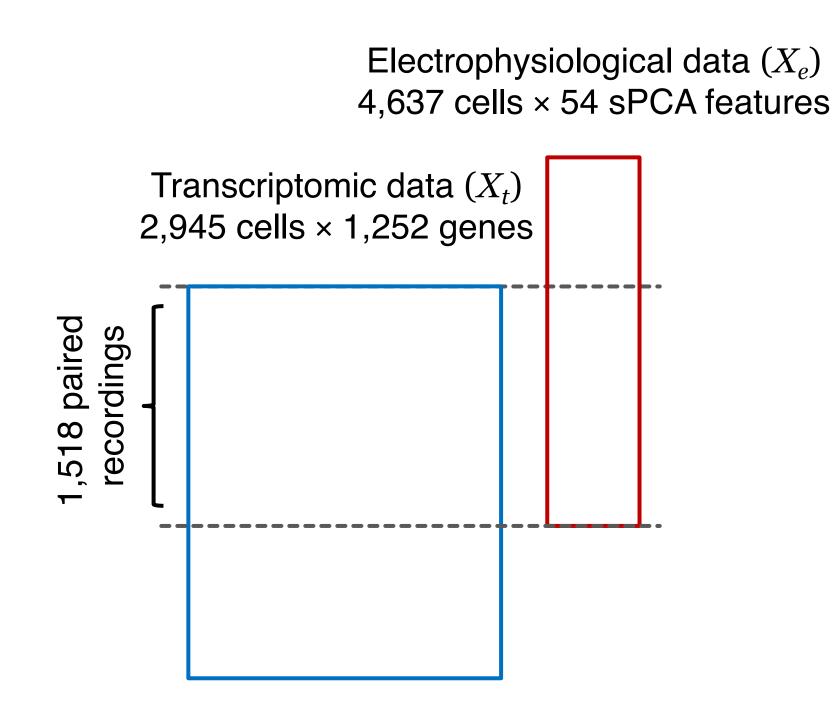
➤ Minimum singular value based normalization outperforms full covariance based normalization in regimes where the covariance estimates are unreliable.



#### 3. Dataset

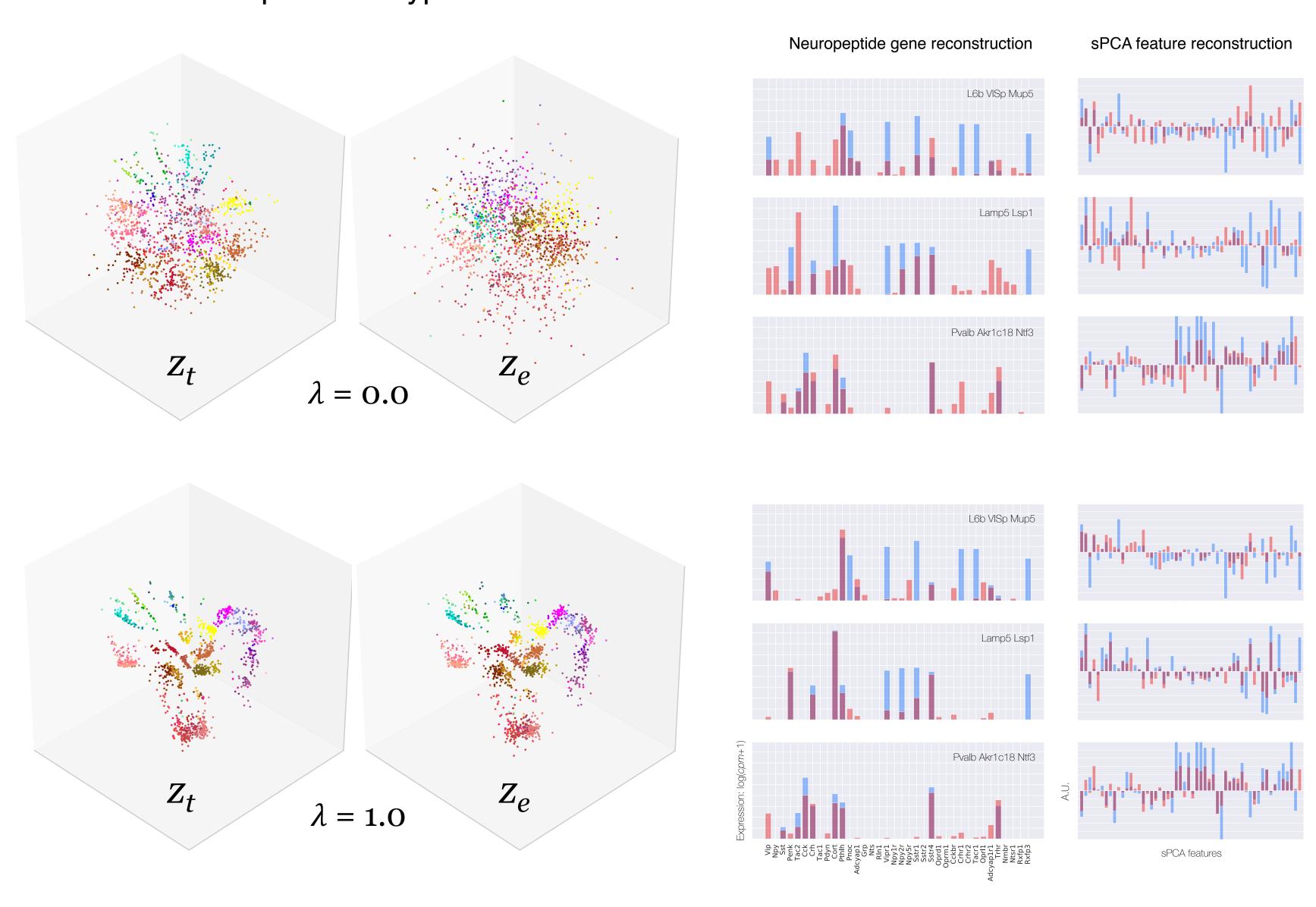
➤ Patch-seq dataset analyzed here consists of expression profiles of 1,252 genes across 2,945 neurons, and 54 sPCA features of electrophysiological recordings from 4,637 neurons.

➤ A subset of 1,518 neurons were profiled with both data modalities, and mapped to an established taxonomy [Tasic et al. 2018] based on differential expression of marker genes.



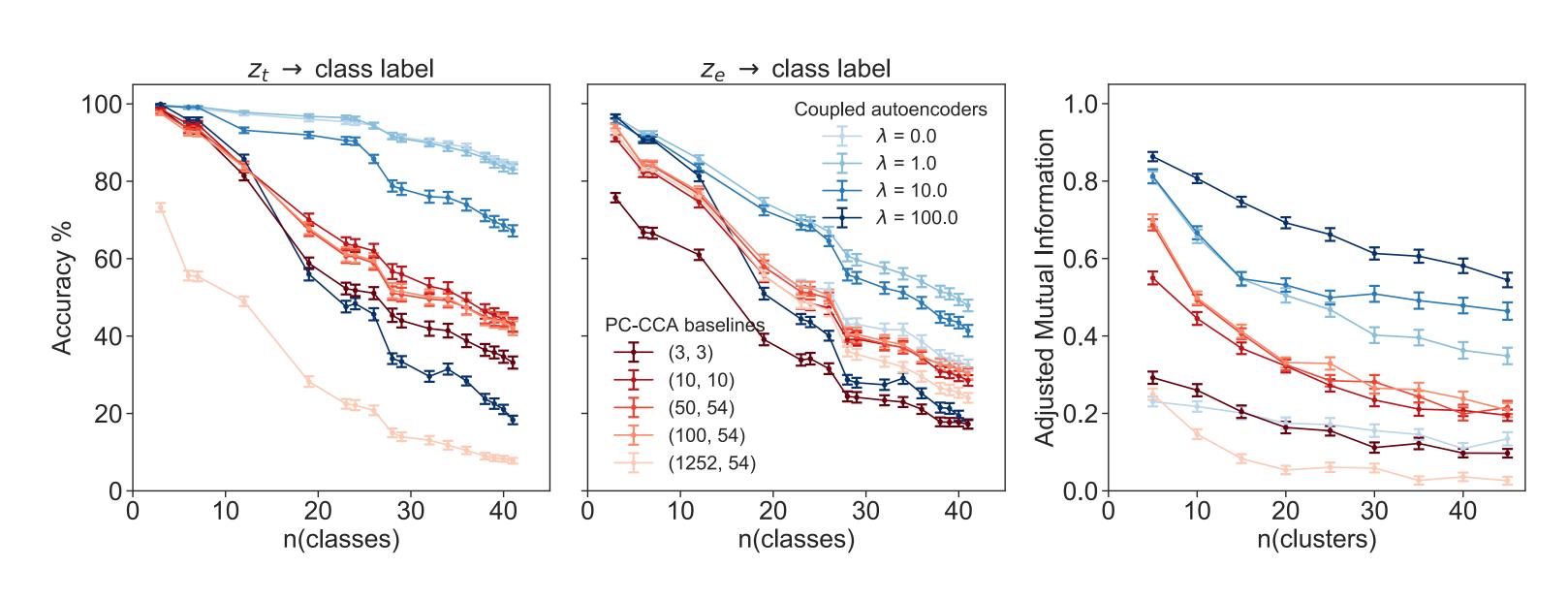
# 4. Results with patch-seq dataset

➤ Known cell types cluster as distinct islands in low dimensional (3d) representations. These representations also appear to preserve hierarchical relationships of cell types.



There is a trade-off between how similar the representations are, and accuracy of the reconstructions the can be inferred from the representations. Coupling constant  $\lambda$  explicitly controls this tradeoff.

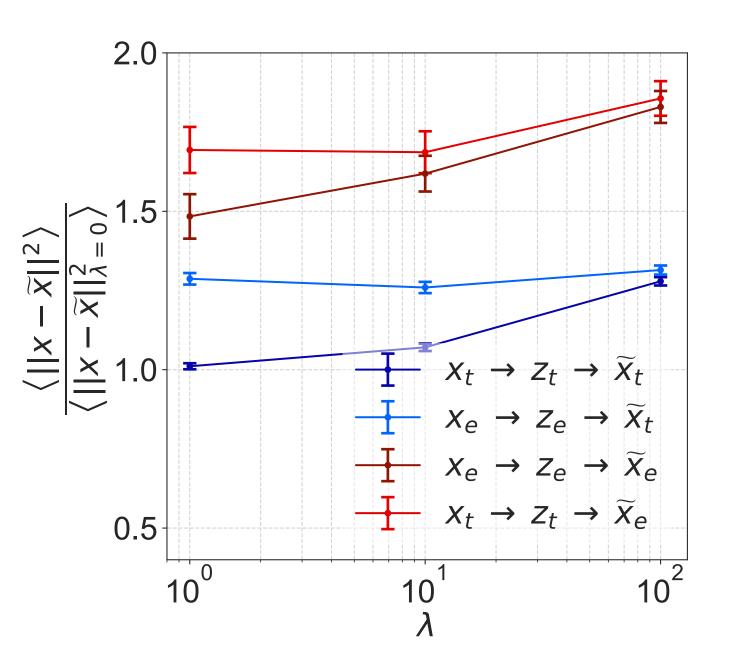
# 5. Cross modal cell type classification



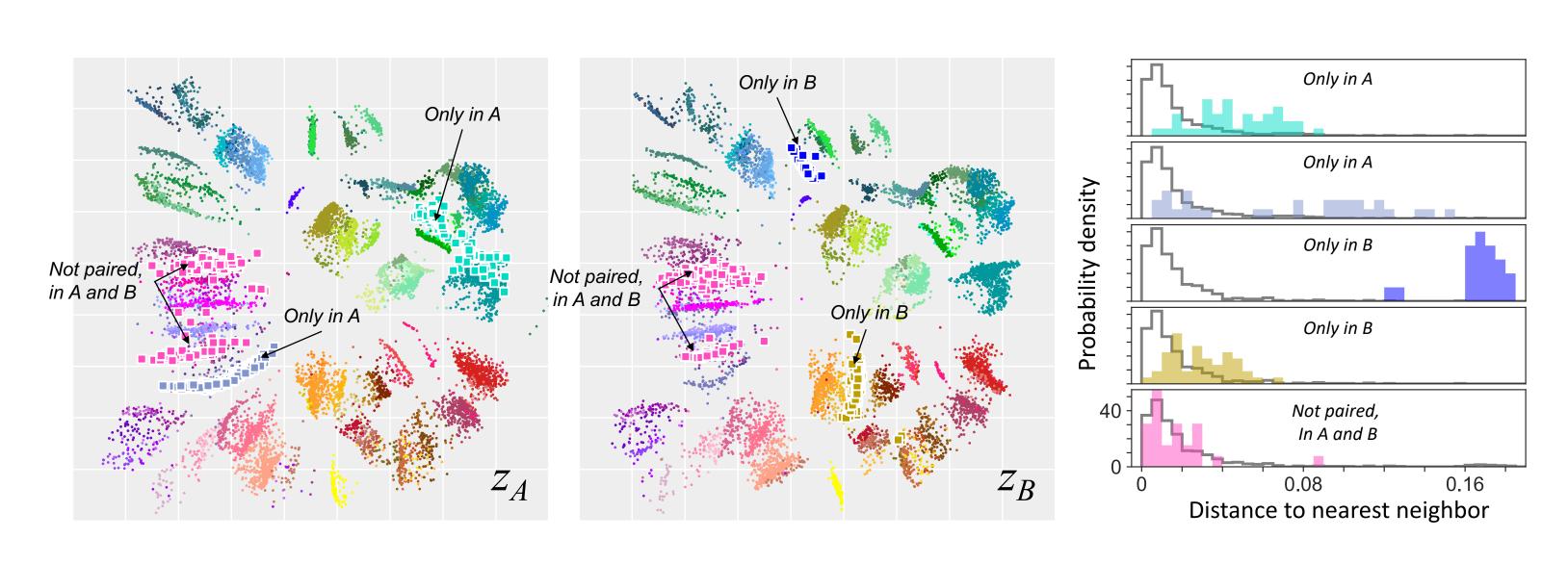
- > Joint representations agree with established transcriptomic hierarchy. Representations with coupled autoencoders reveal this agreement better than CCA based approaches.
- > Clusters in the transcriptomic data are consistent with those in the electrophysiology data. Thus, transcriptomic profiles of neurons seem to map better to functionally consequential properties better than previously thought.

## 6. Cross modal data prediction

- ➤ Jointly learned representations enable prediction of transcriptomic gene expression profiles from electrophysiology features and vice versa.
- ➤ Learning these relationships can minimize the number of measurements required to precisely determine a neurons' distinctive characteristics across modalities.



# 7. Discovery of shared and distinct cell types



➤ Cell types that are unique to a dataset, or are shared across datasets but the association is unknown can both be identified.

#### 8. References

- 1. Cadwell CR., et al. Electrophysiological, transcriptomic and morphologic profiling of single neurons using Patch-seq. *Nature biotechnology* (2016)
- 2. Gouwens NW., et al. Classification of electrophysiological and morphological neuron types in the mouse visual cortex. *Nature neuroscience* (2019)
- 3. Tasic B., et al. Shared and distinct transcriptomic cell types across neocortical areas. *Nature* (2018)

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