# Electrophysiological and Transcriptomic Features Reveal a Circular Taxonomy of Cortical Neurons - *Supplementary Material*

The Supplementary Material is divided into four sections. The first section describes briefly the Machine Learning supervised methods considered in the paper, whereas the second section provides boxplots of the FMM model parameters by Cre line. The third section contains details about the different principal component analysis conducted to obtain the proposed circular taxonomy. The final section shows graphically the agreement between the proposed taxonomy and the results in (Tasic et al., 2016).

## 1 MACHINE LEARNING SUPERVISED METHODS

Short descriptions of the methods are given below, based on (Hastie et al., 2009) and (Izenman, 2008).

- **Linear discriminant analysis (LDA)**: supervised method which seeks the linear combination of features that maximizes the between-class variance relative to the within-class variance. It assumes that the probability functions conditioned to the class are distributed as multivariate normals with homoscedasticity. This simple and interpretable method performs well if the classes can be linearly separated. However, it may underfit the data, leading to high bias models that fail to capture patterns.

- **Random forest (RF)**: supervised method in which various decision trees, induced with random subsets of the input features, are aggregated under bagging. Decision trees discriminate observations of different classes by splitting the data into subsets based on the features values with criteria like the Gini impurity. In RF, the final class prediction is decided by majority voting of the underlying trees. It has become a widely used method due to its good results in many fields without requiring deep knowledge of the methodology by the user. RF provides directly feature importance measurements and while, it normally achieves a good balance between bias and variance, it occasionally overfits the data.

- **Gradient boosting decision trees (GBDT)**: supervised method in which 'weak' decision trees are induced sequentially to eventually lead to a final prediction based on a weighted majority vote. Classifiers are considered 'weak' when they make decisions slightly better than a random choice, meaning that their bias is high. Should the user be capable of a fine hyperparameter tuning, the results achieved in many complex problems are excellent. However, the data is overfitted sometimes.

- **Support vector machines (SVM)**: supervised method in which the maximal separating hyperplane between the classes is constructed in the input feature space. The procedure becomes more flexible by means of the 'kernel trick', which enlarges the features space to find nonlinear boundaries between classes. In particular, the SVM implementation used searched for polynomial boundaries. SVM are considered black box methods as their predictions cannot be easily reasoned. Regularization is often used as these models can overfit the data, resulting in models with high variance.

- **Model averaged neural network (AvNNet)**: supervised method in which multiple multi-layer perceptrons with different random weight initializations are trained to discriminate the classes in the data. Multi-layer perceptrons are neural networks in which data is fed forward through multiple hidden layers to extract useful and discriminating features. The relative lack of knowledge about the inherently extracted features as well as the reasoning behind their predictions make neural networks the most notorious black box method. Nevertheless, the accuracy of their predictions is often unrivaled. Data is frequently overfitted by these models, making essential regularization techniques such as weight decay.

## 2    BOXPLOTS OF THE FMM PARAMETERS BY CRE LINE

In this section, Figures S1-S5 depict boxplots of the main FMM parameters by Cre line. In these figures, the parameters corresponding to the representative neurons of the Cre line have been highlighted as stars.
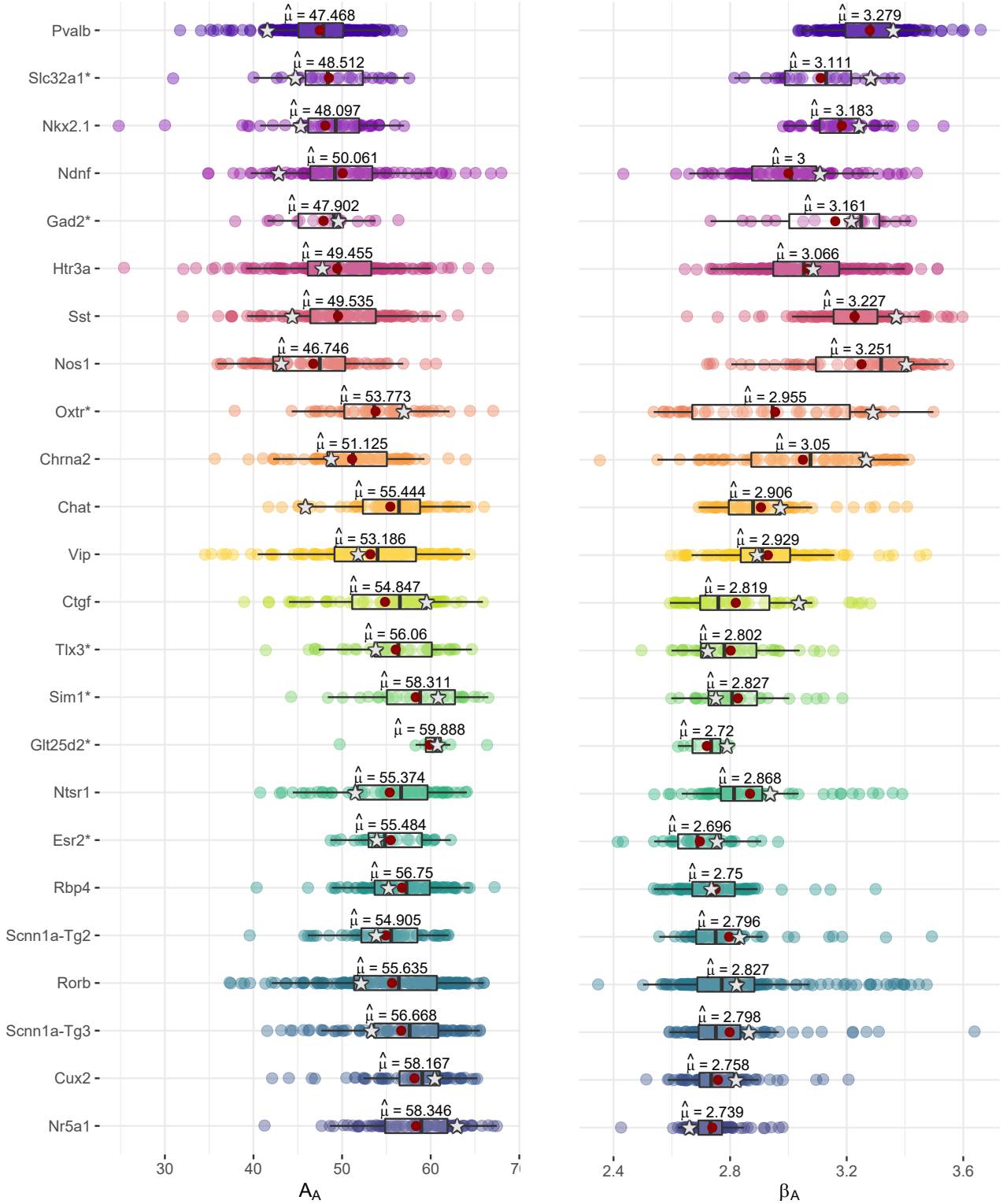


**Figure S1.** Distribution of the $A_A$ and $\beta_A$ parameters by Cre line. Parameters of representative neurons used throughout the manuscript are highlighted as stars.
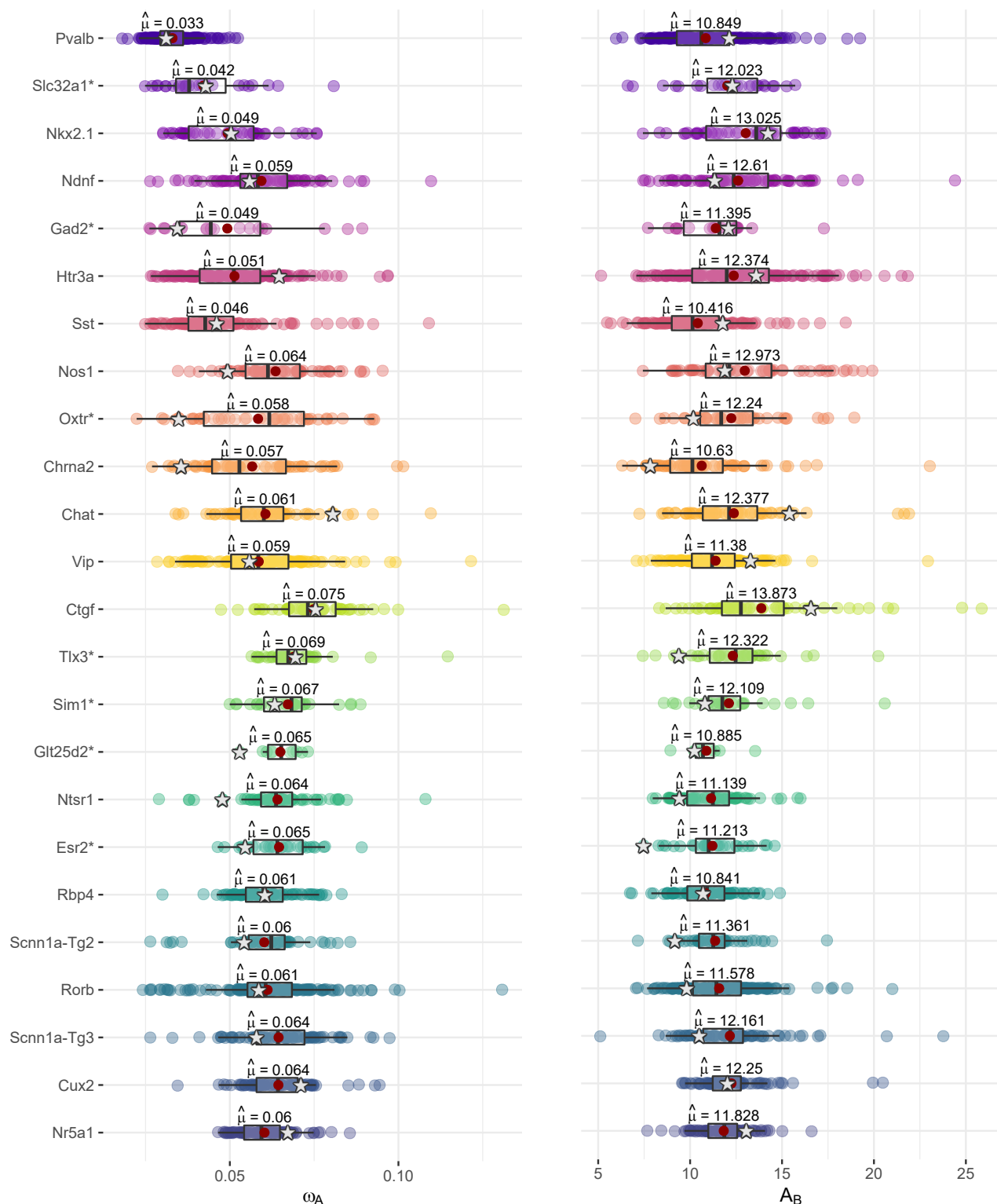
**Figure S2.** Distribution of the $\omega_A$ and $A_B$ parameters by Cre line. Parameters of representative neurons used throughout the manuscript are highlighted as stars.
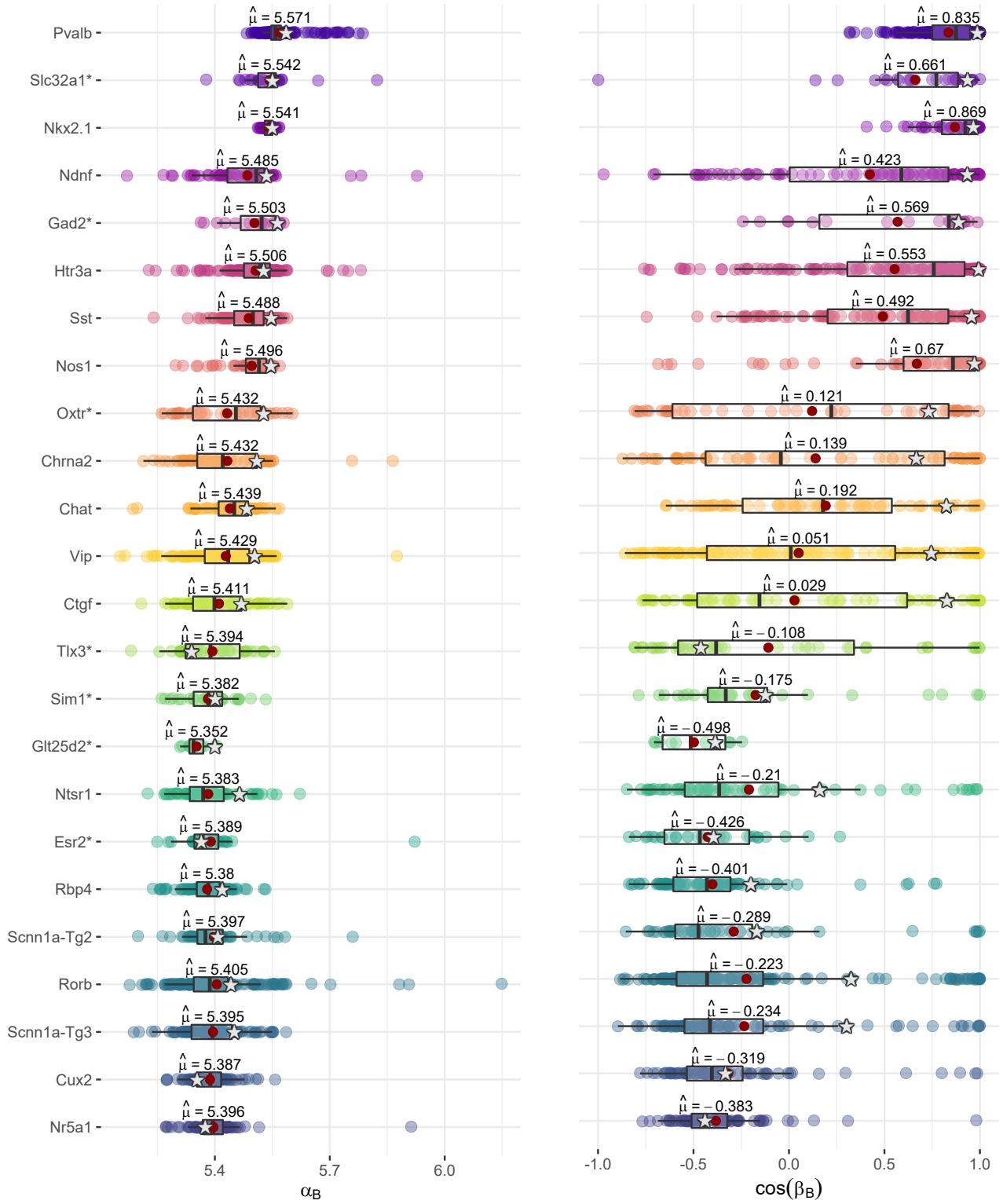
**Figure S3.** Distribution of the $\alpha_B$ and $\cos(\beta_B)$ parameters by Cre line. Parameters of representative neurons used throughout the manuscript are highlighted as stars.

**Figure S4.** Distribution of the $\omega_B$ and $\alpha_C$ parameters by Cre line. Parameters of representative neurons used throughout the manuscript are highlighted as stars.
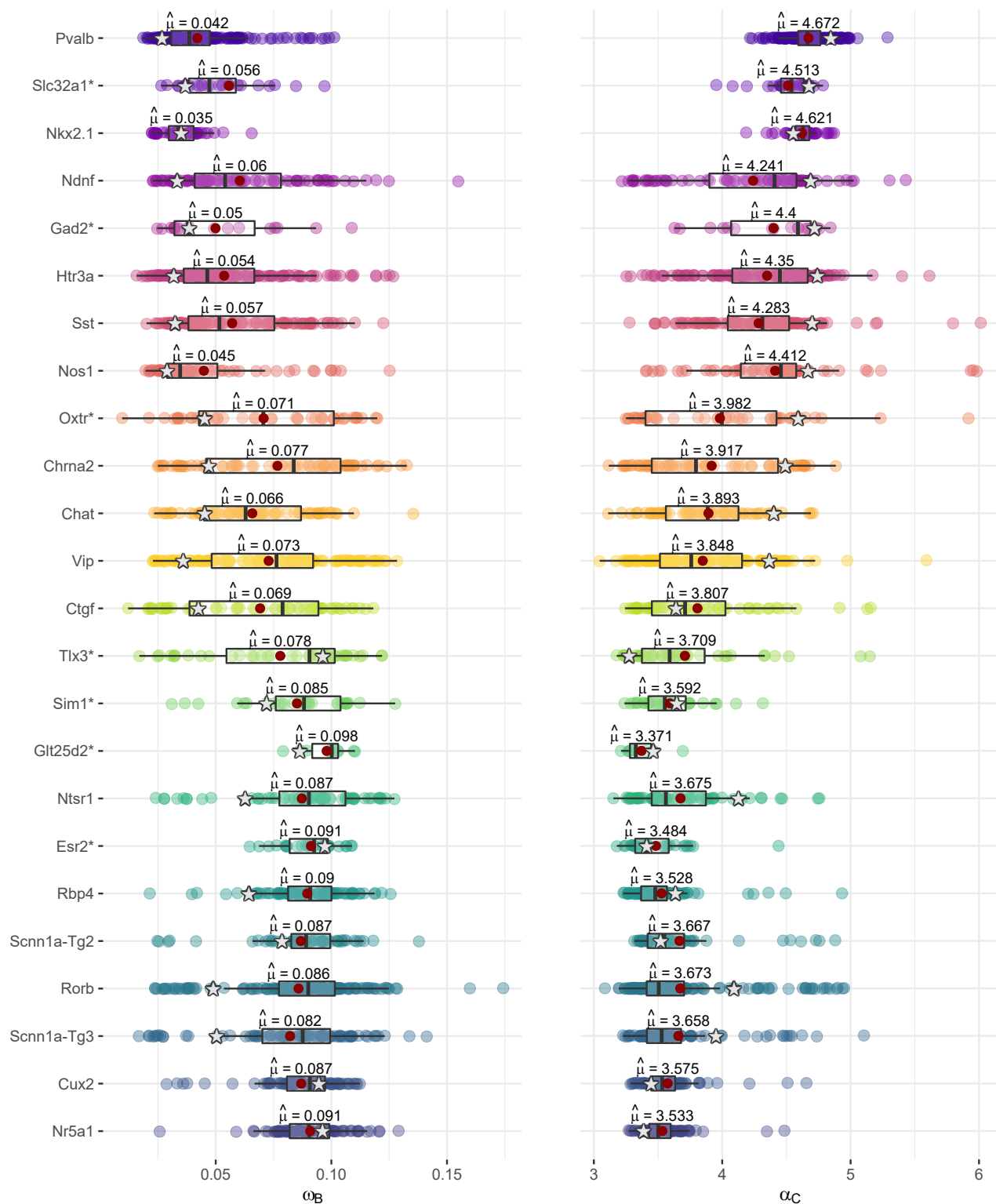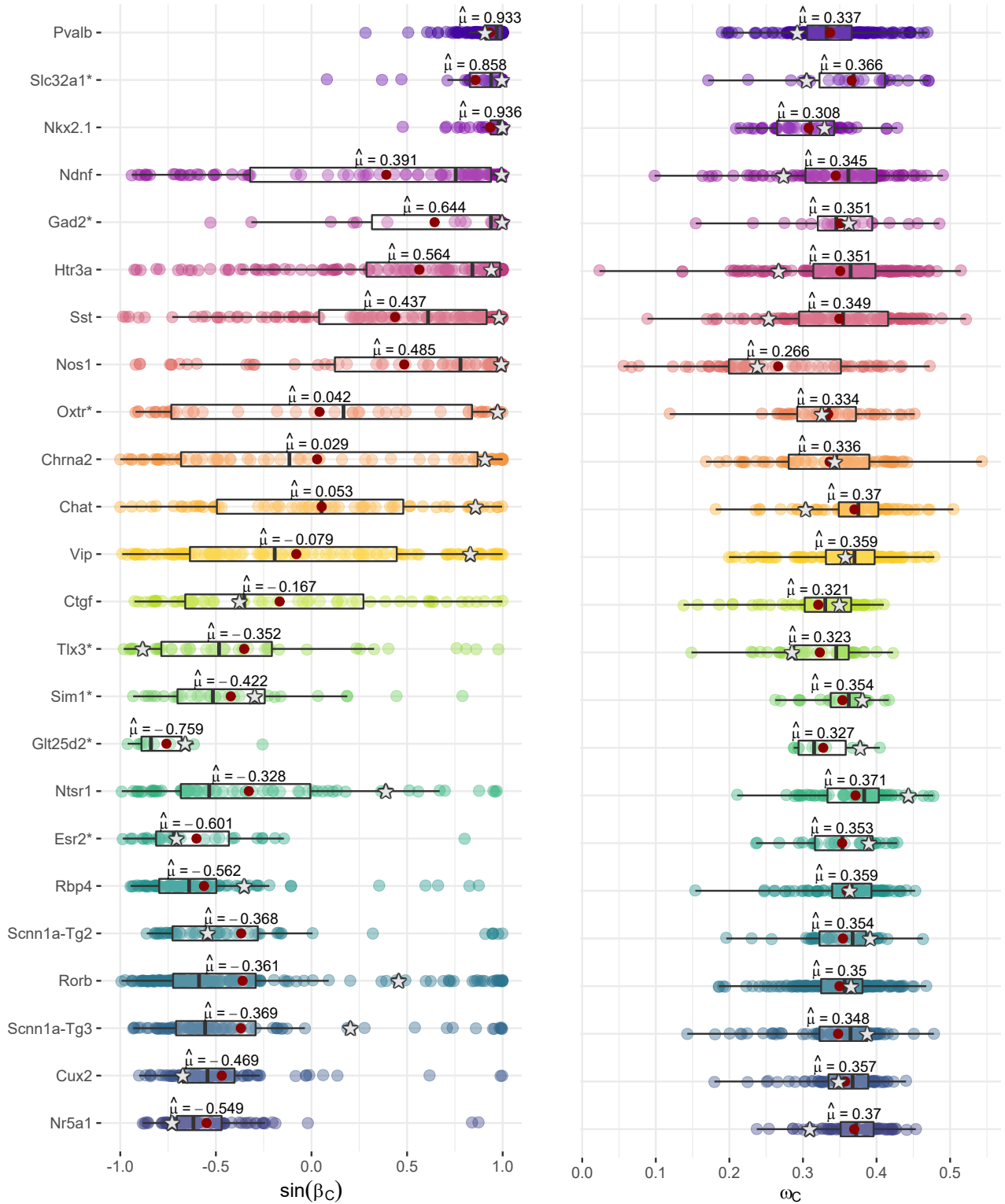
**Figure S5.** Distribution of the $\sin(\beta_C)$ and $\omega_C$ parameters by Cre line. Parameters of representative neurons used throughout the manuscript are highlighted as stars.

# 3 PRINCIPAL COMPONENT ANALYSIS (PCA)

This section illustrates in Figures S6-S10 different aspects of the three principal component analyses conducted in order to obtain the proposed circular taxonomy.
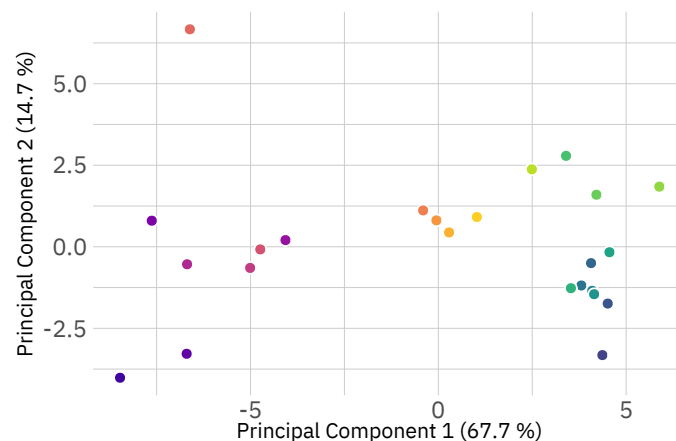
**(A)**

| | M | $A_A$ | $A_B$ | $A_C$ | $\alpha_A$ | $\alpha_B$ | $\alpha_C$ | $\beta_A$ | $\omega_A$ | $\omega_B$ | $\omega_C$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Component 1 | 0.98 | 0.95 | | | -0.74 | -0.99 | -0.99 | -0.95 | 0.8 | 0.98 | |
| Component 2 | | | | | | | | | | | |

| | $\sin(\beta_B)$ | $\cos(\beta_B)$ | $\sin(\beta_C)$ | $\cos(\beta_C)$ | $d_{AB}$ | $d_{AC}$ | $d_{BC}$ | $Var_A$ | $Var_B$ | $Var_C$ | $t_A^U$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Component 1 | -0.85 | -0.99 | -0.99 | 0.83 | 0.96 | 0.99 | 0.99 | | 0.88 | | 0.87 |
| Component 2 | | | | | | | | 0.9 | | -0.83 | |

| | $t_B^U$ | $t_C^U$ | $t_A^L$ | $t_B^L$ | $t_C^L$ | $f(t_A^U)$ | $f(t_B^U)$ | $f(t_C^U)$ | $f(t_A^L)$ | $f(t_B^L)$ | $f(t_C^L)$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Component 1 | | 0.99 | -0.77 | 0.84 | -0.97 | 0.95 | | -0.73 | 0.84 | 0.91 | |
| Component 2 | | | | | | | | | | | 0.79 |

**(B)**



Cre Line: Pvalb, Ndnf, Sst, Chrna2, Ctgf, Tlx3*, Rbp4, Scnn1a-Tg3, Slc32a1*, Gad2*, Nos1, Chat, Glt25d2*, Ntsr1, Scnn1a-Tg2, Cux2, Nkx2.1, Htr3a, Oxtr*, Vip, Sim1*, Esr2*, Rorb, Nr5a1

**Figure S6.** (A) Electrophysiological features correlation with extracted principal components. (B) Scores of the Cre lines with the electrophysiological features PCA.
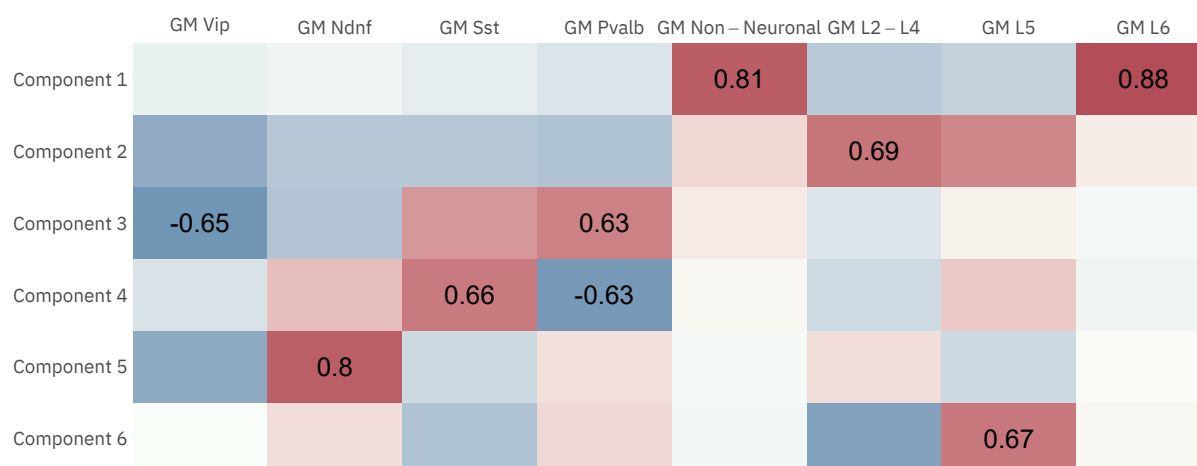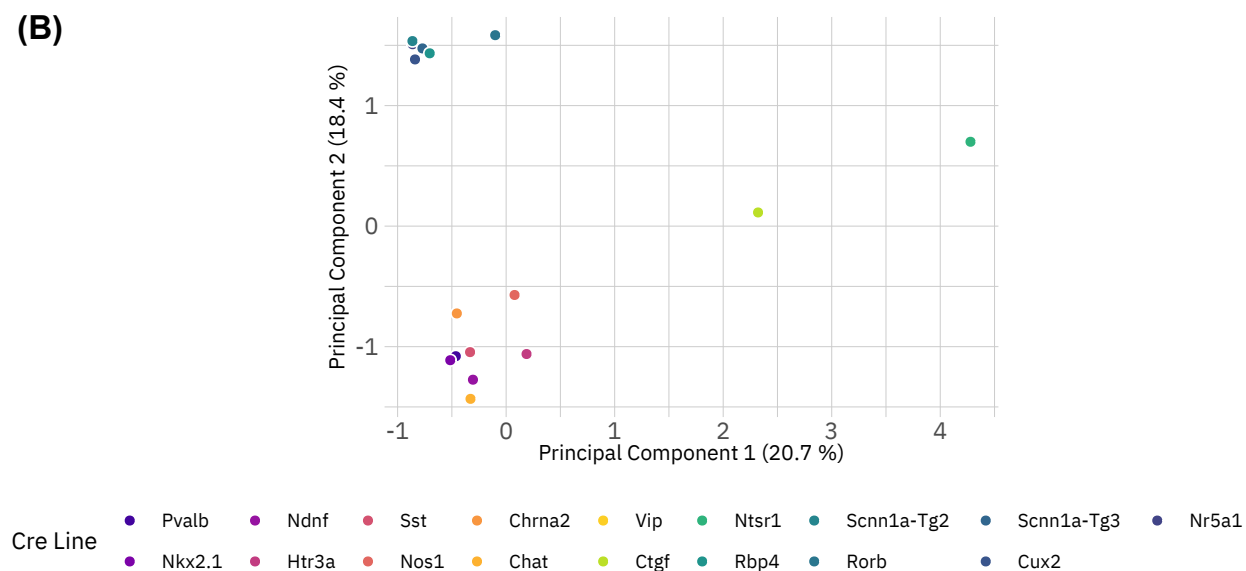
**(A)**

| | GM Vip | GM Ndnf | GM Sst | GM Pvalb | GM Non−Neuronal | GM L2−L4 | GM L5 | GM L6 |
|---|---|---|---|---|---|---|---|---|
| Component 1 | | | | | 0.81 | | | 0.88 |
| Component 2 | | | | | | 0.69 | | |
| Component 3 | -0.65 | | | 0.63 | | | | |
| Component 4 | | | 0.66 | -0.63 | | | | |
| Component 5 | | 0.8 | | | | | | |
| Component 6 | | | | | | | 0.67 | |

**(B)**



**Figure S7.** (A) Transcriptomic features correlation with extracted principal components. (B) Scores of the Cre lines with the transcriptomic features PCA.
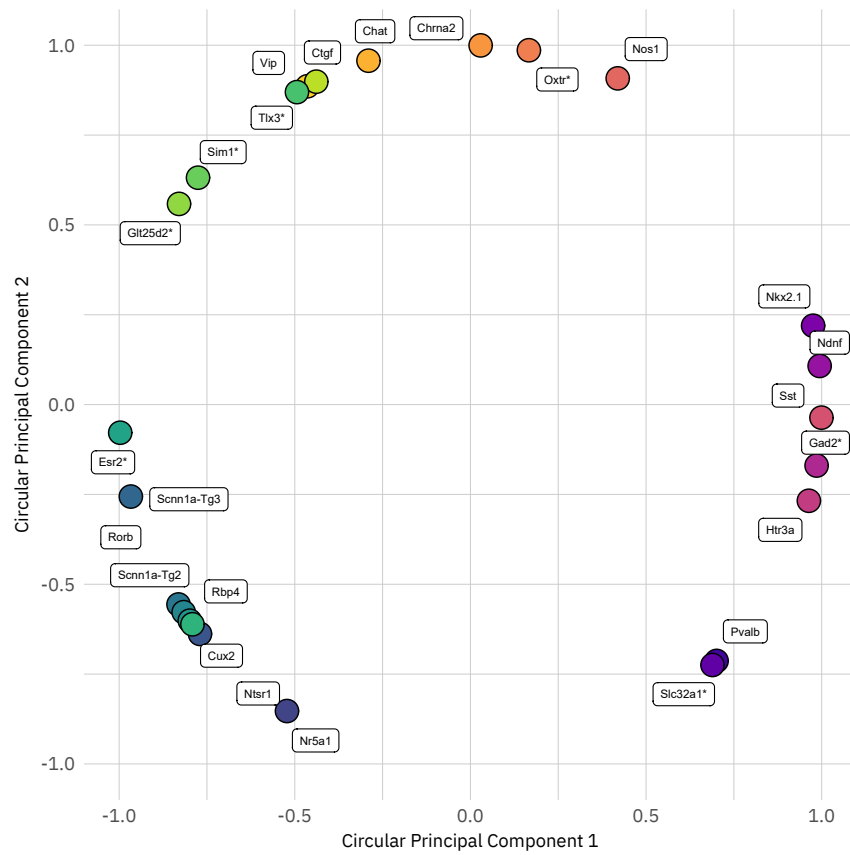
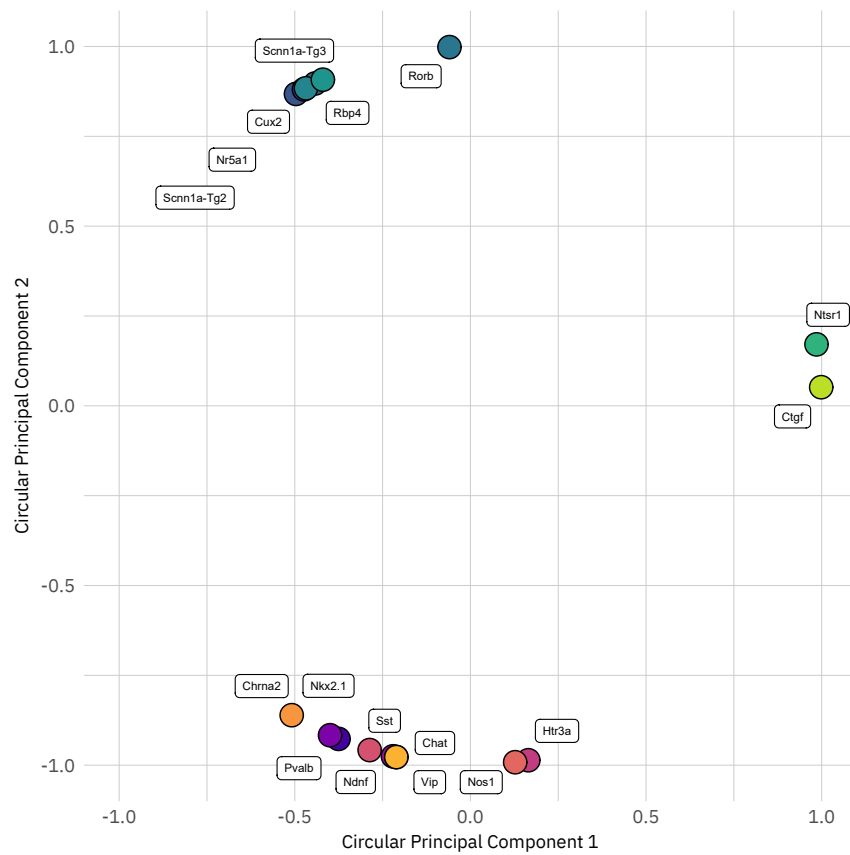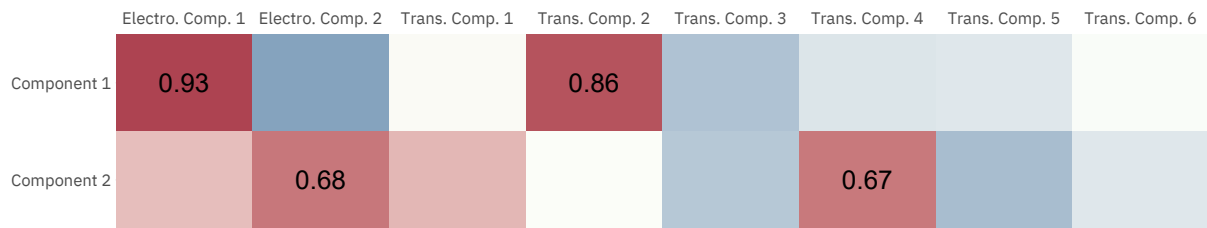**Figure S8.** CPCA of the electrophysiological features.



**Figure S9.** CPCA of the transcriptomic features.

**(A)**

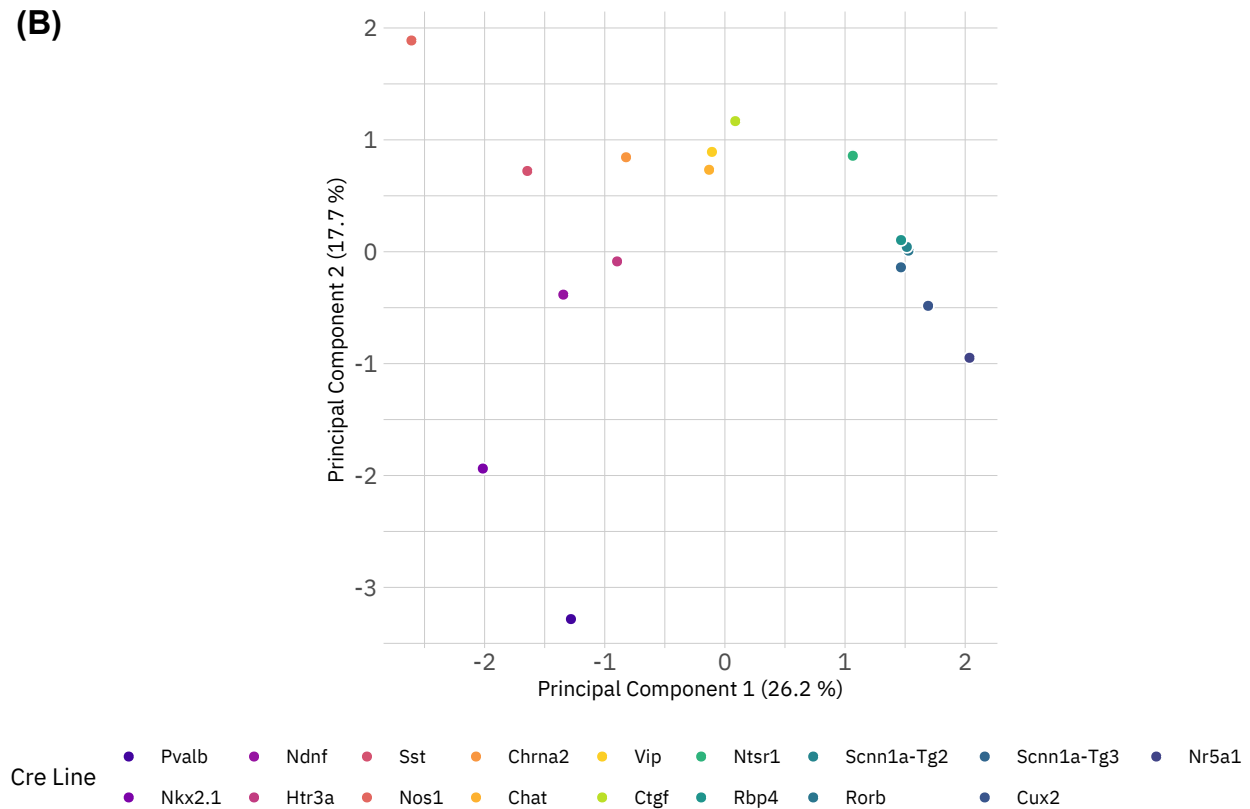| | Electro. Comp. 1 | Electro. Comp. 2 | Trans. Comp. 1 | Trans. Comp. 2 | Trans. Comp. 3 | Trans. Comp. 4 | Trans. Comp. 5 | Trans. Comp. 6 |
|---|---|---|---|---|---|---|---|---|
| Component 1 | 0.93 | | | 0.86 | | | | |
| Component 2 | | 0.68 | | | | 0.67 | | |

**(B)**



**Figure S10.** (A) Extracted electrophysiological and transcriptomic components correlation with extracted ensemble principal components. (B) Scores of the Cre lines with the ensemble PCA.

# 4 TRANSCRIPTOMIC CELL-TYPE MATRIX

Figure S11(A) shows the number of core cells by genetic cluster and Cre line from (Tasic et al., 2016). This Figure's rearrangement according to the proposed circular taxonomy order is shown on Figure S11(B).
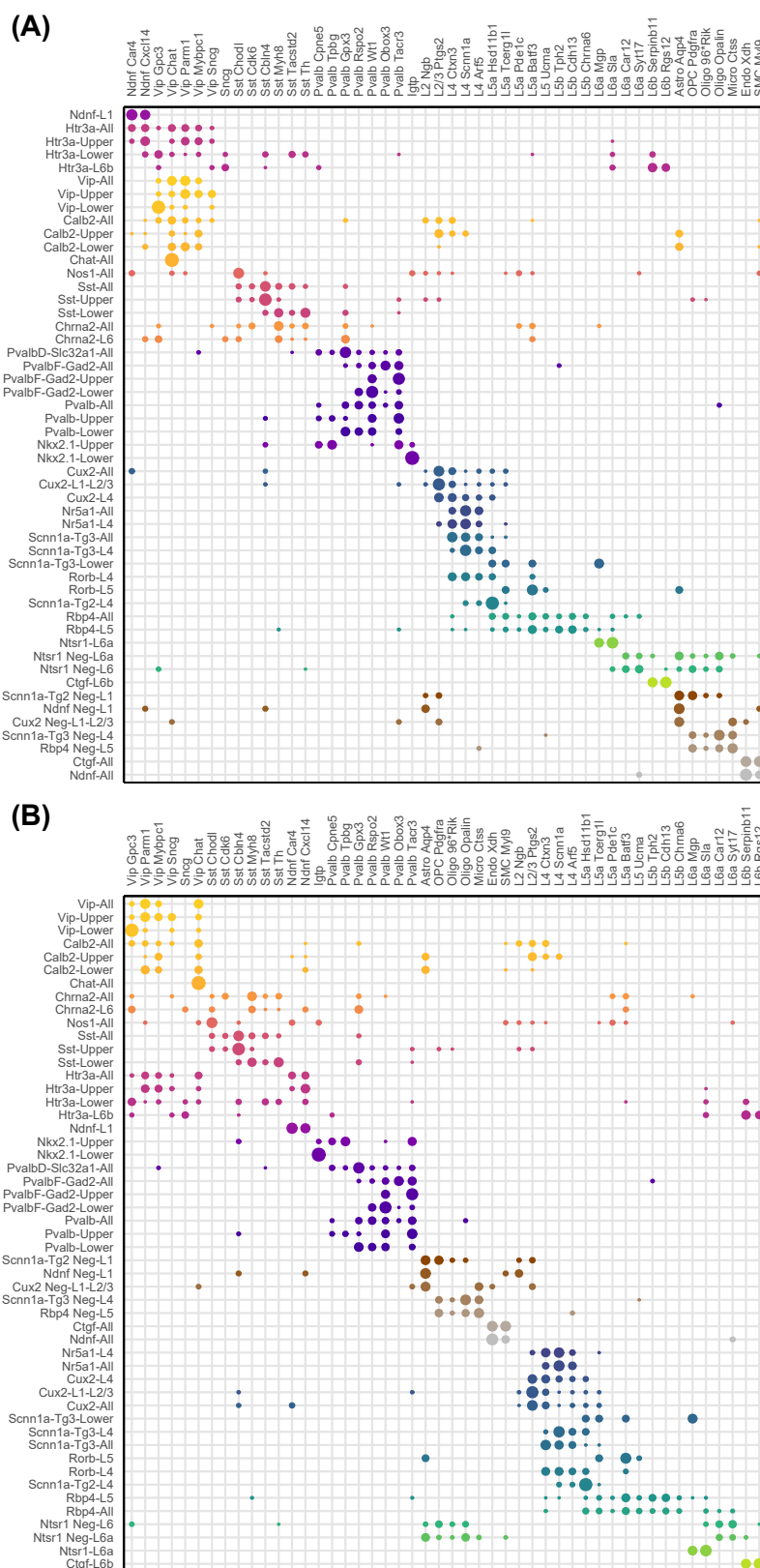


**Figure S11.** (A) Arrangement of the number of core cells by genetic cluster and Cre line according to (Tasic et al., 2016). (B) Rearrangement according to the proposed circular taxonomy.

# REFERENCES

Hastie, T., Tibshirani, R., and Jerome, F. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (Springer), second edn.

Izenman, A. J. (2008). *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning* (Springer), first edn.

Tasic, B., Menon, V., Nguyen, T. N., Kim, S., Jarsky, T., Yao, Z., et al. (2016). Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nature Neuroscience* 19, 335–346