

Investigating Neuronal Cell Classes and their Role in Cognition

by

Emily Huang

B.S. Computation and Cognition, MIT, 2022

Submitted to the
Department of Electrical Engineering and Computer Science, Brain and Cognitive Sciences
in partial fulfillment of the requirements for the degree of

MASTER OF ENGINEERING IN COMPUTATION AND COGNITION

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

February 2024

© 2024 Emily Huang. All rights reserved.

The author hereby grants to MIT a nonexclusive, worldwide, irrevocable, royalty-free license to exercise any and all rights under copyright, including to reproduce, preserve, distribute and publicly display copies of the thesis, or release the thesis under an open-access license.

Authored by: Emily Huang
Department of Brain and Cognitive Sciences
February 7, 2024

Certified by: Earl Miller
Picower Professor of Neuroscience, Thesis Supervisor

Accepted by: Katrina LaCurts
Chair
Master of Engineering Thesis Committee

Investigating Neuronal Cell Classes and their Role in Cognition

by

Emily Huang

Submitted to the

Department of Electrical Engineering and Computer Science, Brain and Cognitive Sciences
on February 7, 2024 in partial fulfillment of the requirements for the degree of

MASTER OF ENGINEERING IN COMPUTATION AND COGNITION

ABSTRACT

Classifying neurons into different cell classes is both an idea that has existed since the origins of neuroscience, and one that is essential to understanding the complex interactions of the brain. While there has been a substantial effort to categorize neurons morphologically, molecularly and physiologically in *in vitro* studies, there is a gap in experiments performed on awake and behaving animals. Using data collected from macaque monkeys performing a working memory task, and employing an unsupervised Gaussian mixture model (GMM) clustering algorithm, a number of different cell classes and their defining features were distinguished in area 7A, the lateral intraparietal area (LIP), the dorsolateral and ventrolateral prefrontal cortex (PFC) and the extrastriate visual area (V4). While the number of cell classes found across areas differed, there were several classes across areas that appeared to be correlates. Classes in each area also showed functional differences in information encoding during predictable trials and distributional differences in depth. This signifies both the potential of functionally distinct cell classes involved in prediction, as well as the existence of universal cell classes across different areas.

Thesis supervisor: Earl Miller

Title: Picower Professor of Neuroscience

Acknowledgments

I would like to express my gratitude to Dr. Earl Miller for supervising my work and for welcoming me into his lab, despite the newness of this program. I would also like to express my gratitude to Scott Brincat for mentoring me on this project, and for always being so patient and open. His insight was crucial in propelling this project forward.

A big thank you to my previous research mentors, Leo Kozachkov, Katya Tsimring and Michael Liu-Happ for starting me out on this path, and giving me guidance in not just research, but in all areas of life. Thank you for being someone for me to look up to when I was still a bright-eyed, bushy-tailed undergrad.

Another big thank you to Dr. Karchmer for taking me in as a TA for 6.1210 and giving me what was genuinely one of the best experiences of my life, as well as allowing me to fund this pursuit. And to all my students, thank you for being the reason I woke up in the morning. You made my weeks brighter.

Lastly, I would like to deeply thank my family and my close friends for being a source of unwavering support and love throughout both my time in college and this masters program. Even through the ups and downs of this process, I was always able to turn to someone. The only reason I'm able to keep beating against the current is because of you – thank you.

Contents

Title page	1
Abstract	3
Acknowledgments	5
List of Figures	9
List of Tables	11
1 Introduction	13
2 Methods	16
2.1 Datasets	16
2.1.1 Task Design	16
2.2 Spynal Library	17
2.3 Data Pre-Processing	18
2.3.1 Spike and Waveform Data	18
2.3.2 LFP Data	20
2.4 Feature Extraction	20
2.4.1 Spike Waveform Feature Extraction	20
2.4.2 Spike Rate Feature Extraction	21
2.4.3 Spike Timing Feature Extraction	21
2.5 Clustering	21
2.5.1 Optimal Cluster and Model Selection	22
2.6 Depth Distribution	22
2.7 Information Encoding	23
2.8 Spike Field Coupling	23
3 Results	25
3.1 Preliminary Clustering	25
3.2 Cell Classes	26
3.2.1 Class Features	27
3.2.2 Average Waveforms	28
3.2.3 Cluster Distribution	29
3.2.4 Depth Distribution	29

3.3	Cell Class & Firing Rate Variation	29
3.4	Cell Class & Spike Field Coupling	30
4	Discussion	48
4.1	Distinct cell classes	48
4.2	Functional and distributional differences between cell classes	49
4.3	Correlates between classification methods	50
	References	51
A	Code listing	54
A.1	Loading Data	54
A.2	Pre-processing	59
A.3	Clustering	68
A.4	Analysis	71
A.5	Plotting	72
A.6	Utilities	80

List of Figures

2.1	(a) DMTS task setup; (b) Recording areas illustrated. <i>NOTE:</i> FEF is labelled, but was not used in any analysis. Both figures taken from laminarPharm documentation.	17
3.1	Average waveforms per preliminary clustering area. Outlier waveforms shown in red. Average waveform shown in black. (a) 7A ; (b) PFC ; (c) V4	33
3.2	Mean BIC values and SEM values (shown as blue highlight) for each area (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	34
3.3	Mean AIC values and SEM values (shown as blue highlight) for each area (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	35
3.4	Five parameter values defining each cell class in areas (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	37
3.5	Visualization of clusters in the 2D space of each pair of features. Marginal distribution of each feature on the diagonal. (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	39
3.6	Average waveforms per preliminary clustering area. Outlier waveforms shown in light grey. Average waveform shown in black. (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	41
3.7	Percentage of data points in each cluster for areas (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	42
3.8	Depth distribution for clusters in areas (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	43
3.9	Percent explained variance of trial predictability on cluster firing rate (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	44
3.10	(a) PEV of each cluster on sample type in LIP ; (b) Raster plot of neuronal activity over different sample trials in LIP	45
3.11	Percent explained variance of trial sample type on cluster firing rate, with predictable trial blocks removed (a) LIP ; (b) 7A ; (c) PFC ; (d) V4	46
3.12	Spike-field coupling synchrony of select PFC units with superficial LFP signals (a) Cluster 0 unit (Group N1) ; (b) Cluster 5 unit (Group L1) ; (c) Cluster 7 unit (Group B1)	47

List of Tables

3.1	Clustering Information, laminarPharm Only	26
3.2	Clustering Information, Full Dataset	26

Chapter 1

Introduction

The identification and classification of neurons into different neuronal cell types is essential for understanding the complex behaviors and the local and distal circuit dynamics of the brain. While the foundations of neuroscience had its origins in attempting this classification [1], fascinated by all the diversity the brain has to offer, attempts were often hindered by the lack of computational methods available at the time, resulting in a laborious, and often biased approach to the problem [2]. Focus in the field shifted towards more hypothesis driven studies. However, with the recent development of reliable technological methods that significantly reduce human bias and offer the ability to perform powerful computations on large datasets, there has been a resurgence of this effort to identify different neuronal sub-classes, their properties, and their functional roles [3], [4].

There are several different approaches to the classification of cell types, generally categorized by morphological, molecular or physiological features. Morphological methods aim to classify cells based on their structure, such as their dendritic and axonal shapes, and their branching patterns, while molecular methods separate classes through properties such as protein and mRNA composition [2]. Physiological methods, on the other hand, distinguish neurons through properties such as firing rate, inter-spike intervals, and waveform shape. Although a large number of different cell classes have been identified using these methods

in *in vitro* studies, until recently, there had been a lack of classes identified using data collected from awake animals [3], as well as limited studies done with data from the primate pre-frontal cortex, critical for higher level cognition and decision making [5]. It is essential to observe the function of different cell types within a natural framework, with an animal engaging directly in tasks. Only then, can the full function of a cell class, as well as its interactive dynamics, both local and distal, be understood.

Classification from extracellular recordings has shown a bimodal distribution of spike-waveform widths, with each peak characterizing "broad-spiking and "narrow-spiking" cells. The broad-spiking cells are thought to be a correlate to excitatory pyramidal cells, while the narrow-spiking cells are thought to be GABAergic, inhibitory cells [6]. However, more recent studies have diversified and broken down the two categories even more, suggesting that there are more features than just spike waveform width that can be used to identify cell classes [3]–[5], [7]–[9]. One such study found four clusters, two of which mirrored the traditional broad and narrow-spiking cells, and the other two of which mirrored potential morphologically identified classes, non-fast spiking inter-neurons, and intrinsically bursting neurons [4]. Others identified many more cell types - up to seven or eight, with even more diversity beyond broad and narrow spiking [3], [5], [8]. These studies included more than just waveform features in their clustering analysis; they introduced other features of the neuron, such as firing rate features, and inter-spike interval features, and it made clear that there was more to cell classification than just the shape of the waveform itself; a neuron had other intrinsic properties that allowed it to be differentiated.

Once these features of interest are identified, a computational approach to grouping different neurons based on these features needs to be chosen. Since the problem of classification is an inherently unlabeled problem, an unsupervised clustering algorithm such as K-means clustering [3], [5], [8], [9] or a Gaussian mixture model [4] is chosen. Both models have advantages in different situations. K-means clustering uses a measure of squared euclidean distance and attempts to minimize the sum of the averages of each cluster [10]. On the other

hand, a Gaussian mixture model uses a model that maximizes expectation, and data points are assumed to have a normal distribution. This often makes it more sophisticated than K-means, and more able to handle non-linear data [11].

This project isolated five features of interest – trough to peak time and repolarization time, two markers of waveform width, the mean firing rate of the neuron, and its coefficient of variance (CV) and local variance (LV), two measures of how regular the neuron’s firing patterns are. A Gaussian Mixture Model was chosen over a k-means model due to its lack of bias in the perceived shape of the clusters, since there was no prior that suggested what geometry the data would be.

Chapter 2

Methods

2.1 Datasets

Two datasets collected in the Miller Lab were used for this project – the laminarPharm dataset, collected by Alex Major, and the wmPredict dataset, collected André Bastos. The data was collected from two macaque monkeys, Selma and Lucky. Data was formatted into the standard Miller Lab data format. Preliminary clustering and plotting was done using only the laminarPharm dataset; the wmPredict dataset was then added so that the clustering algorithm could have more data points to produce a more accurate clustering. Since both datasets were formatted using the standard Miller lab format, most functions were easily extended – any discrepancies were dealt with manually in the code.

2.1.1 Task Design

Both datasets contain data from working memory tasks, with manipulations on sample predictability and laminar recording. The laminarPharm dataset contained additional data with manipulation of laminar pharmacology, though that data was not used for this project. The task was a delayed match to sample task (DMTS), which consisted of two different trial blocks – blocks in which the sample was predictable, and blocks in which the sample was

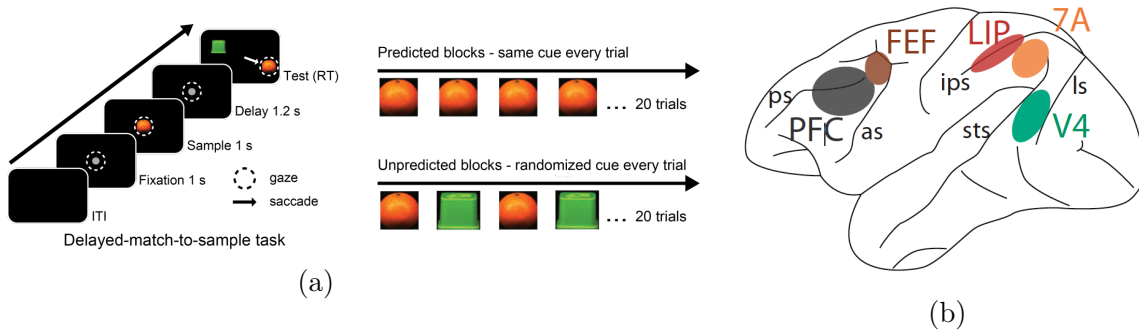


Figure 2.1: (a) DMTS task setup; (b) Recording areas illustrated.

NOTE: FEF is labelled, but was not used in any analysis.

Both figures taken from laminarPharm documentation.

randomly selected from a set of three images (unpredictable). Samples were chosen from a set of 12 used across all sessions. In the wmPredict dataset, blocks were 50 trials long. In the laminarPharm dataset, they were 20 trials long.

Each trial consisted of a fixation period (1 second), followed by a sample image (1 second), a delay period (variable) and finally, a test array containing 1-2 distractor objects. Subjects had to saccade to the image matching the sample shown previously. The task is illustrated in Fig. 2.1a.

Data in the laminarPharm datasets was recorded from the dorsolateral and ventrolateral prefrontal cortex (dlPFC and vlPFC respectively), area 7a of the posterior parietal cortex (7A) and the extrastriate visual area (V4). In addition to those areas, data in the wmPredict dataset recorded from the lateral intraparietal area of posterior parietal cortex (LIP) was also used. Data from dlPFC and vlPFC was combined for the analysis, and compiled into datasets labeled 'PFC'. Areas are highlighted in 2.1b.

2.2 Spynal Library

The analysis in this project was heavily supported by Spynal, a neural data analysis library coded by Scott Brincat [12]. Various functions from Spynal were used, mainly in pre-processing data. These functions included loading data, extracting spike rates, inter-spike

intervals, waveform statistics, concatenating arrays, calculating percent explained variance, and interpolating waveform data. Additionally, plotting functions from Spynal was used to create the raster plots shown later on.

2.3 Data Pre-Processing

2.3.1 Spike and Waveform Data

Data pre-processing consisted of several steps: (1) Loading and reshaping the data, (2) Area selection, (3) Trial Selection, (4) Unit selection, (5) Wave Alignment and (6) Concatenation. Code for all pre-processing, feature extraction, clustering and plotting can be found in Appendix A.

Loading Data

All relevant data for each session was loaded using the matIO module from Spynal. Spike and waveform data was reshaped such that all individual trains were at least 1-dimensional arrays, and all individual waveforms were at least 2-dimensional. The data for the proper areas was selected by isolating the corresponding indices for the area from the 'unitInfo' dataframe. Separation by area was chosen so that area would not be a potential hidden parameter used in clustering. Additionally, it's possible that not every area has the same distribution or prevalence of varying cell-types, so the data from different areas was separated to avoid inaccurate grouping.

Trial Selection

Valid trials were defined as trials that were (1) correct and (2) ≤ 5 trials before drug onset (in the 'laminarPharm' dataset). Trials after drug onset were chosen to not be used in analysis due to previous analysis suggesting an almost complete suppression of activity after drug onset, and thus would not provide much use for clustering.

Unit Selection

Valid units were defined as units that (1) had a mean spike rate of ≥ 1 spike/ second and (2) $< 0.1\%$ of inter-spike intervals that were within 1 millisecond. These criteria weeded out silent cells and poorly isolated cells respectively. Additionally, a few outlier units with inverted waveforms due to an error in the data recording process were also not used for the analysis.

Wave Alignment

Wave alignment served to align all waveforms to the mean trough of all the waveforms in a session. Firstly, waves were interpolated by 10x using the `interpolate` Spynal function. Then, the mean waveform over all units was calculated. The trough of that waveform was extracted to be the mean trough that all the waveforms would be aligned to. All waveforms were aligned to the trough by shifting the waveform either left or right; the ends were padded with NaN values as necessary. Zero-padding was avoided, so that zero values would not affect any means taken over the waveforms. Then, the mean aligned waveform was found by taking the mean over all of the aligned waveforms, ignoring NaN values.

Concatenation

Data from across all sessions was concatenated into large dataframes. Only sessions with more than two valid trials and two valid units were used – otherwise, they were not concatenated. All dataframes were then saved for ease of future access, so that they would not have to be recalculated every time the clustering algorithm ran.

Merging vIPFC and dIPFC

In order to create the PFC dataset, data for vIPFC and dIPFC were pre-processed, concatenated and saved separately, and then manually merged. dIPFC data was always concatenated to the vIPFC data, for consistency.

2.3.2 LFP Data

LFP data was not heavily processed due to the data having already been pre-processed into the Miller lab data format. Similar to spike and waveform data, the units from the area of interest were isolated in the LFP data. Then, LFP data was pooled into two groups – data recorded from electrodes above to Layer 4, and data recorded from electrodes deeper than Layer 4. Depth information was taken from the corresponding area electrodes in the `electrodeInfo` dataframe. Depths were recorded based on the estimated relative depth of the electrode to layer 4, such that a label of 0 corresponded to layer 4, a label of > 0 corresponded to deep layers, and a label of < 0 corresponded to superficial layers.

2.4 Feature Extraction

Five features of interest were identified. Two were spike waveform features – the time in milliseconds from a waveform’s trough to its peak, and repolarization time. Two were spike rate features – firing rate, coefficient of variation (CV). The last feature was a spike timing (ISI) feature, local variation (LV). The five features of each unit was recorded in a Pandas dataframe.

2.4.1 Spike Waveform Feature Extraction

Spike waveform features for each unit were extracted from the mean aligned waveform for that unit, which was found through the wave alignment section of pre-processing. Both trough to peak time and repolarization time were calculated using the `waveform_stats` Spynal function. An array of the trough to peak time and repolarization time for each unit was returned, and added as columns in the feature dataframe.

2.4.2 Spike Rate Feature Extraction

The mean firing rate for each unit was taken from the firing rates found using the `rate` Spynal function. Firing rate was averaged over all time points, as well as over all trials for a unit. An Anscombe transform was then performed on the mean rates, so that the data would be normalized.

The coefficient of variation for each unit was calculated using the `rate_stats` Spynal function, which took the mean firing rate of each trial per unit as input, and returned $\frac{SD(rates)}{\text{mean}(rates)}$ for each unit.

Arrays with the mean firing rate and CV were added as columns in the feature dataframe.

2.4.3 Spike Timing Feature Extraction

The local variation of each unit was found by calculating the LV of each trial of the unit using the `isi_stats` Spynal function. If the trial had ≤ 1 ISI, meaning there were no intervals to compare against each other, the LV was recorded as NaN. After all the LV of all trials was calculated, the mean LV was found by taking a mean over all LVs, ignoring NaN values. The final array of the mean LV per unit was added as a column to the feature dataframe.

2.5 Clustering

The optimal number of clusters for each area was found by performing Gaussian Mixture Model (GMM) clustering on the given feature dataframe. Before feeding the data into the GMM, the data was standardized by using sklearn's `StandardScaler` and `fit_transform` functions, which produces a scaler and then fits the scaler to the data and transforms it, respectively. GMM clustering was performed using sklearn's `GaussianMixture` function from the `mixture` package.

2.5.1 Optimal Cluster and Model Selection

500 repetitions of the GMM algorithm was performed, each repetition fitting the model on component sizes in the range 2-27. For each repetition, every time a model was found for a certain component size, the model was then fit to the features, and cluster labels for each unit (which cluster each unit belonged to based on its five features) were predicted.

The Bayesian information criterion (BIC) and Akaike information criterion (AIC) values were then calculated by using sklearn's `bic` and `aic` functions from the `mixture` package. The BIC value, AIC value, the current model, and labels were then recorded. The best model per repetition was determined by finding the model with the minimum BIC value; the best component size per repetition then recorded as the component size of the model that produced that minimum. AIC value was not used to find the best fitting model since AIC doesn't penalize for larger component sizes. BIC penalizes for larger component sizes, so that overfitting doesn't occur.

The best component size (c) over all repetitions was determined by calculating the mode of the array of best component sizes per repetition. The best c -component model was then found by taking the c -component model with the lowest BIC value. The model itself, the labels produced by that model, and c was returned for each area. Cluster labels were saved for ease of access when performing analyses.

2.6 Depth Distribution

In order to investigate the distribution of each cluster across different depths, depth data was isolated from the `unit_info` dataframes from both datasets. Both datasets took a measure of depth relative to Layer 4 that was estimated by the crossing point of relative power in beta and gamma bands. In the `wmPredict` dataset, this was called `betaGammaDepth`, and in the `laminarPharm` dataset, this was called `laminarDepth`. In both datasets, negative numbers meant more superficial depths, while positive numbers meant deeper depths. Zeros

were considered to be within layer 4.

In the `wmPredict` dataset, depths that could not be reliably estimated were labelled as NaN values. These values were converted to a value of -3, which was lower than the minimum depth value recorded, to allow the investigation of whether certain clusters had more depth information missing than others.

In addition to the original depth data, which was discrete, a "jittered" version of the depth data was generated to aid in plotting the distributions. Jittered values helped prevent the points from being plotted on top of each other due to the confines of the figure. The array of jittered values was created by randomly selecting from a uniform distribution in the range of `[-0.02, 0.02]` using numpy's random package, `np.random.uniform`. The random values were then added to the original depth values.

2.7 Information Encoding

The information encoding of each cluster was measured by finding the percent explained variance of each unit's firing rate over time by the variation in sample (three groups) or block predictability (two groups). The percent explained variance was calculated using the Spynal function `neural_info`.

In addition to using the full sample data, a truncated version of the sample labels was created by filtering out the trials that occurred in predictable blocks, due to the presence of units that would indiscriminately fire during those predictable blocks. The percent explained variance calculated using the truncated sample data was also calculated using the same Spynal function.

2.8 Spike Field Coupling

Spike field coupling/ synchrony was calculated using the Spynal function `spike_field_coupling`, which uses the default method of measuring the magnitude of pairwise phase consistency

(PPC). The LFP data was truncated to the three seconds surrounding sample onset ($[-1, 2]$) so that its time dimensions would match the spike time dimensions used for previous analyses. In each area, the unique probes used for recording in that area were identified. Only linear probes were used – single electrodes (only present in the `wmPredict` dataset) were excluded from this analysis, due to the lack of comparable depth data. To reduce processing burden, the lfp data for a specific probe was grouped into "deep" lfp data, measured with the "deep" electrodes (relative to Layer 4), and "superficial" lfp data, measured with the "superficial" electrodes. The spike data for that probe was then isolated by referencing the unique ID of the probe.

Spike times were converted to boolean spike trains using the Spynal function `times_to_bool`. Then, for every unit, the PPC between the spikes recorded by that unit, and the deep and superficial lfp data was calculated, yielding two datasets. The data was then visualized using the Spynal `plot_spectrogram` function.

Chapter 3

Results

3.1 Preliminary Clustering

The preliminary clustering done with only the laminarPharm dataset, and for only datasets recorded from 7A, PFC, and V4, revealed an optimal component size of 4, 3, and 5 for 7A ($n_pts = 354$), PFC ($n_pts = 121$), and V4 ($n_pts = 261$) respectively. Table 3.1 illustrates this data. The wmPredict dataset was added after preliminary results to increase the number of data points for the model to be able to predict on, as the original dataset was small, especially for the PFC.

In the preliminary clustering, some clusters had high variability in waveform shape (Figure 3.1). Most notably, the last cluster of area 7A has an increased number of outlier waveforms [3.1a], categorized by those with repolarization times or trough to peak times that were more than 2.5 standard deviations outside of the mean waveform, shown in black. Other clusters, such as the last two clusters in V4[3.1b], and the last cluster of PFC[3.1c], appear to have two different waveform shapes – one with a narrower trough to peak, and one with a broader trough to peak.

Table 3.1: Clustering Information, laminarPharm Only

	7A	PFC	V4
Number Datapoints	354	121	261
Optimal Component Size	4	3	5

Table 3.2: Clustering Information, Full Dataset

	LIP	7A	PFC	V4
Number Datapoints	689	765	1083	1663
Optimal Component Size	7	5	8	8

3.2 Cell Classes

The minimum mean BIC value for each component size revealed an optimal component size of 7, 5, 8, and 8 for areas LIP ($n_{pts} = 689$), 7A ($n_{pts} = 765$), PFC ($n_{pts} = 1083$), and V4 ($n_{pts} = 1663$) respectively. Table 3.2 illustrates this data.

All four areas showed a distinct "elbow" in the mean BIC curve for component sizes, indicating a clear optimum in choosing component size [3.2]. On the other hand, the mean AIC curves don't show a distinct elbow for any areas [3.3] – this is due to AIC allowing for overfitting by not punishing larger component sizes, which is why only the BIC values were used to determine the optimal component size.

The size of the data appeared to have no correlation with the optimal number of clusters found (i.e. more data points fed into the model did not imply a larger component size). This can most notably be seen in LIP having an optimal component size of 7, while having the smallest number of data points out of all the areas (689), while 7A, having the second smallest number of data points out of the areas (765), and closer in size to the LIP dataset than the other two areas, had an optimal component size of 5, smaller than that found for LIP.

3.2.1 Class Features

The mean value of each parameter for each cluster is shown in Figure 3.4. While each area had a differing number of cell classes, there were some trends in parameter values that seemed to remain constant across areas. Each area had a cell class characterized by a broad waveform ($TTP \cong .0065$; $RPT \cong 0.0030$), a lower mean firing rate ($3.5 < MR < 5$), and a semi-regular firing pattern ($0.5 < CV < 1$) and medium "burstiness" ($LV \cong 1.0$), such as cluster 0 in area 7A [3.4b], cluster 7 in area PFC [3.4c], and cluster 7 in area V4 [3.4d]. Clusters 0 and 4 in area LIP seem to have similar properties, only being split up by cluster 0 having a smaller repolarization time, a smaller mean rate, as well as a less regular firing pattern [3.4a]. For convenience, this cluster will be referred to as B1. Inspection of the marginal distribution of each feature for each cluster in Figure 3.5 suggests that B1 is actually a group of neurons with highly variable trough to peak and repolarization times, which is why their mean waveform trough to peak and repolarization times are larger than other clusters.

The trough to peak and repolarization marginal distributions of each area clearly shows a bimodal distribution of waveform width, which would be the distinction between broad-spiking and narrow-spiking neurons. While there are some clusters that fit into either side of the bimodal distribution, there are others, like cluster 3 in 7A which has peaks on both sides. This suggests that what characterizes the cell class is not the waveform shape, but rather the rate or ISI features [3.5b].

In addition to the broad-spiking B1 class, each area also appears to have a narrow waveform cluster ($0.0002 < TTP < 0.0003$; $0.00005 < RPT < 0.00015$) with a very high mean rate ($MR > 6$), and low variability in both pattern and burstiness ($CV < 0.6$; $LV < 0.95$). This pattern corresponds to cluster 3 in area LIP, cluster 0 in area PFC, and cluster 0 in area V4. The cluster in 7A that seems to correspond to this pattern (cluster 3) seems to have a slightly waveform than the others ($TTP \cong 0.00045$; $RPT \cong 0.00018$). For convenience, this

cluster will be referred to as N1.

Besides clusters that can be further categorized within the broad/ narrow spectrum, there appear to be clusters defined by very high/ low variability, or a high/ low mean rate. Cluster 2 in LIP, cluster 2 in PFC and cluster 6 in area V4 can all be characterized by a very high CV, implying a much more irregular firing pattern than the other cell classes in the area. This cell class is also characterized by a low mean rate ($3 < MR < 4$), a narrower waveform ($0.0003 < TTP < 0.0005$; $RPT \cong 0.00015$), and are not particularly bursty ($0.85 < LV < 1.0$). Area 7A doesn't appear to have a direct correlate, although cluster 2 exhibits a similar pattern, with a higher burstiness ($LV \cong 1.15$). For convenience, this cluster will be referred to as C1.

Similarly, cluster 6 in LIP, cluster 5 in PFC, and cluster 4 in V4 seem to be characterized by a high LV, or being more bursty neurons. It also appears to be characterized by a more regular firing pattern, a lower mean rate, and narrow to medium sized waveforms. For convenience, this cluster will be called L1.

3.2.2 Average Waveforms

Cell classes that demonstrated similar parameter values, especially those strongly characterized by their waveform features, appeared to have similar average waveform shapes, and little outliers [3.6]. The broad spiking cluster, B1, had larger trough to peak durations ($\tilde{150}$ timepoints) compared to the narrow spiking cluster, N1, which had a trough to peak duration of less than 100 timepoints, as expected. This can also be seen in other broad spiking clusters, such as cluster 4 of LIP and cluster 3 of V4 [3.6a, 3.6d].

Notably, C1, which was characterized by a high CV, appeared to have two distinct waveform groupings, as can be seen by the average waveform plots for cluster 2 in LIP, cluster 2 in PFC, and cluster 6 in V4 [3.6a, 3.6c, 3.6d]. Similarly, L1, which was characterized by a high LV, shows the same waveform groupings in the average waveform plots for cluster 6 in LIP, cluster 5 in PFC and cluster 4 in V4 [3.6a, 3.6c, 3.6d].

3.2.3 Cluster Distribution

Figure 3.7 shows the percentage of data which fell into each cluster in each area. The distribution of cells into clusters appears to differ by area – for group B1, it comprises 11.0% of units in 7A, 15.3% of units in LIP, 7.7% of units in PFC, and 7.0% of units in V4. N1 comprised of much more units in 7A (29.5%) than the other areas (11.3, 15.6 and 11.8%). C1 comprised of more units in V4 (23.8%) than other areas (12.0, 6.8, 5.4%). L1 seemed to have the largest overall proportion of units, with 22.1% in LIP, 18.4% in PFC, and 16.4% in V4.

3.2.4 Depth Distribution

The distribution of units in each cluster was visually represented in swarmplots, shown in Figure 3.8. Due to lack of data, as well as missing depth data, the shapes of the distribution appears unclear, especially for clusters that had proportionally fewer units in them. Most clusters are missing a proportional number of depth information – the bigger they are, the more unlabeled clusters (plotted at $d = -3$) there are.

From the distributions with more units, clusters 1 and 4 in LIP, clusters 3 and 4 in 7A, clusters 0 and 1 in PFC, and clusters 0 and 5 in V4 appear "pinched" around layer 4 ($d = 4$), and bottom heavy (with a higher distribution of neurons in deep layers). Conversely, cluster 3 in LIP and cluster 3 in PFC are also pinched around layer 4, but appear top heavy (higher distribution in superficial layers). Cluster 5 in PFC and clusters 1 and 2 in V4 appear fatter around layer 4.

3.3 Cell Class & Firing Rate Variation

The percent explained variance (PEV) of the predictability of the block on each cluster's firing rate can be seen in Figure 3.9. Areas LIP and 7A had little to no effect on PEV

between clusters, but PFC showed a small peak in PEV for cluster 6 at sample onset. V4 showed a peak as well, with cluster 0. Cluster 6 had unusually high PEV throughout the time period before and after sample onset. Cluster 0 in V4 was part of group N1, characterized by a narrow waveform, high mean firing rate, and low CV and LV [3.4d]. Cluster 6 in PFC showed similar patterns, albeit with a less narrow waveform [3.4c]. This implies that in these clusters, information about the predictability of the trials in a block is encoded through a high firing rate.

The PEV of the sample image on each cluster’s firing rate was also calculated, but was found to have unusually high elevated PEV throughout the trial [3.10a]. Investigating the activity of individual neurons revealed that some neurons would fire indiscriminately during the entire period only in certain sample blocks [3.10b]. After removing trials from predictable blocks from the PEV calculation, the elevated activity disappeared [3.11].

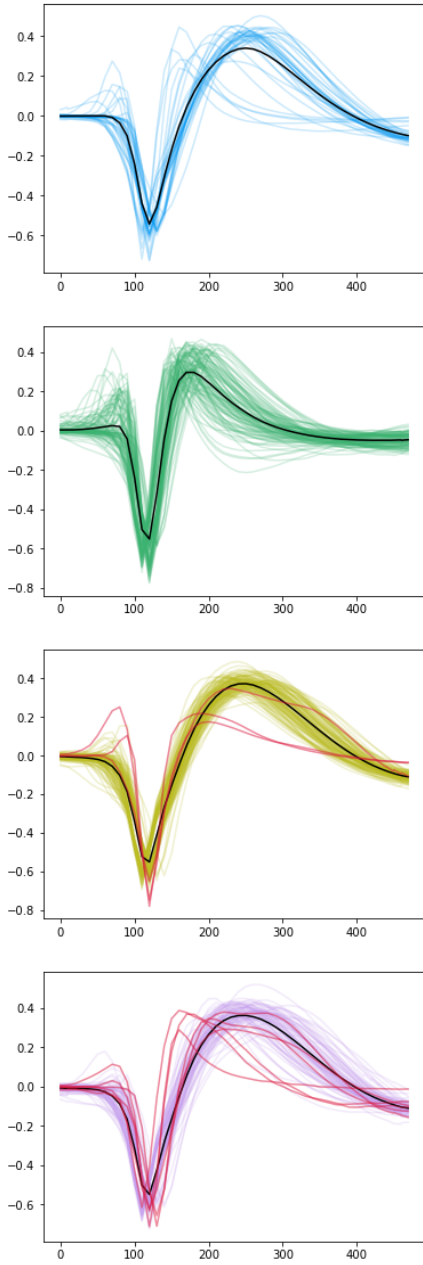
While clusters in 7A didn’t have much of an effect on PEV, cluster 3 in LIP, clusters 0 and 6 in PFC, and clusters 0 and 2 in V4 all showed a higher PEV at sample onset. All clusters had a high firing rate, with half belonging to group N1, that similarly implies that sample information is encoded through firing rate in these clusters [3.4]. Notably, cluster 4 in PFC also high firing rate, but had a lower PEV than clusters 0 and 6 [3.4c]. This could be due to it having a higher CV, and thus having a less regular firing process.

3.4 Cell Class & Spike Field Coupling

Preliminary data from spike-field coupling shows synchronization occurring at different times in the trial in different clusters in PFC. The unit from cluster 0 (part of group N1) exhibits highest synchronization with superficial LFP signals throughout trial start, sample onset, and the delay period [3.12a]. The unit from cluster 5 in group L1 exhibits highest synchronization only *before* sample onset [3.12b]. The unit from cluster 7 in group B1 exhibits highest synchronization before and after sample onset, but not during [3.12c]. While this shows

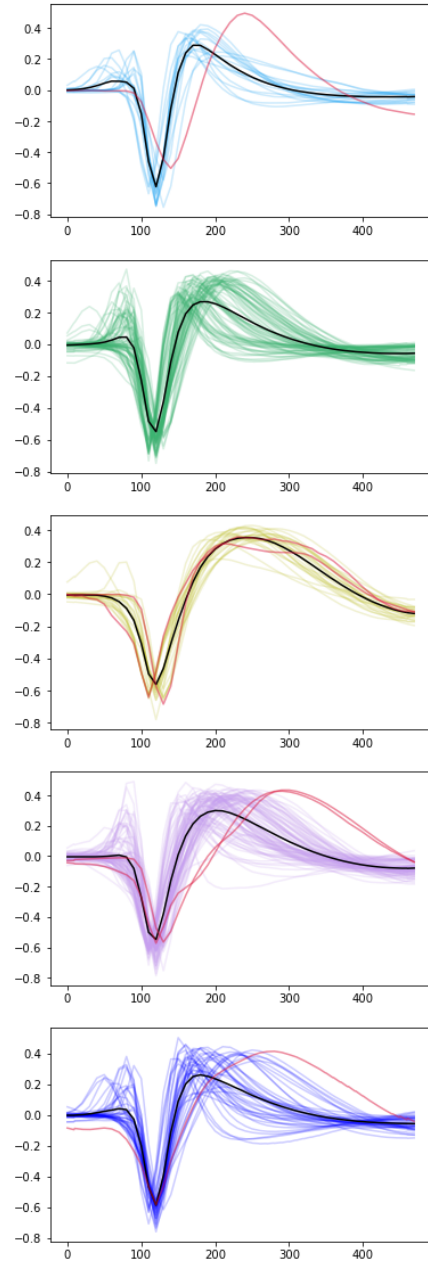
potential for different periods of synchrony in different clusters during prediction periods, the data is non-aggregate and isolated, so further analysis would need to be made.

7A cluster waveforms



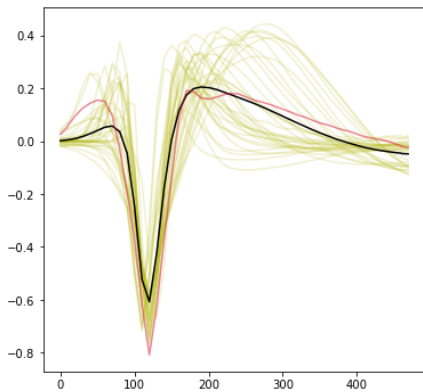
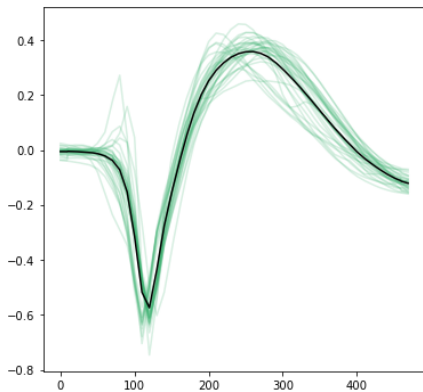
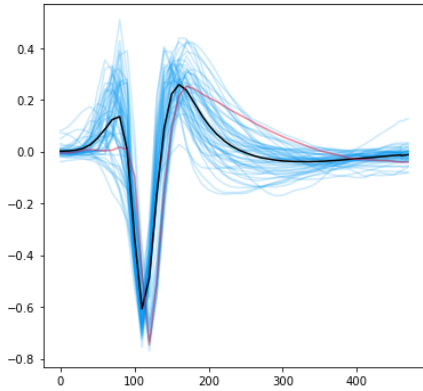
(a)

V4 cluster waveforms



(b)

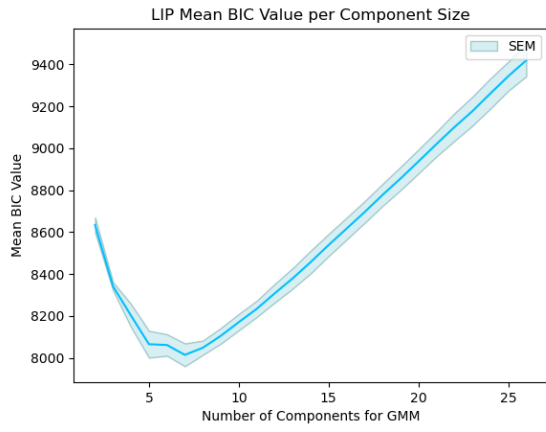
PFC cluster waveforms



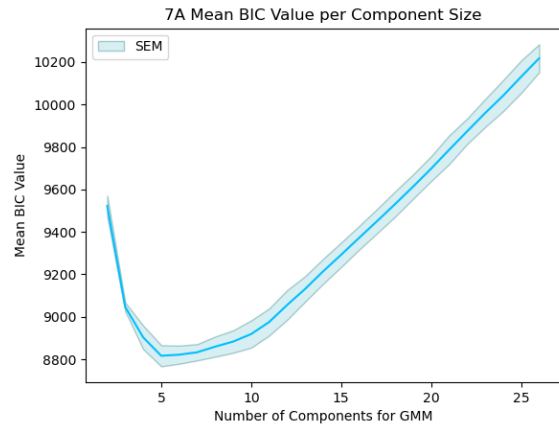
(c)

Figure 3.1: Average waveforms per preliminary clustering area. Outlier waveforms shown in red. Average waveform shown in black.

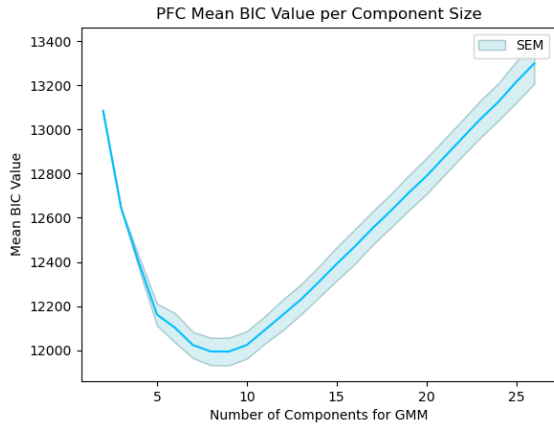
(a) 7A ; (b) PFC ; (c) V4



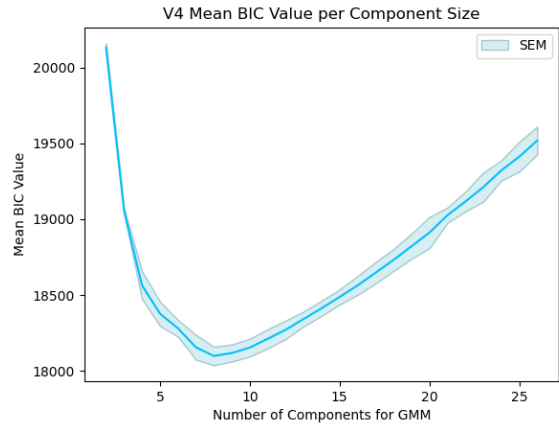
(a)



(b)

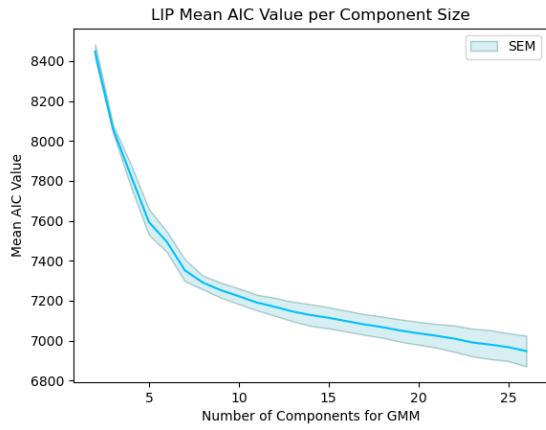


(c)

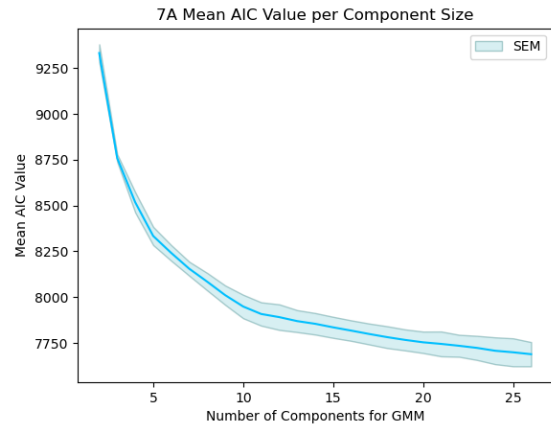


(d)

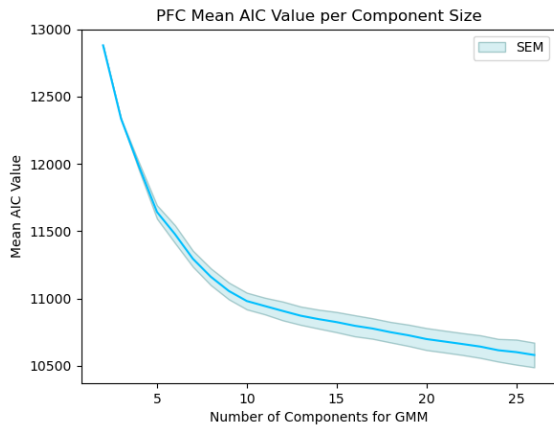
Figure 3.2: Mean BIC values and SEM values (shown as blue highlight) for each area
 (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



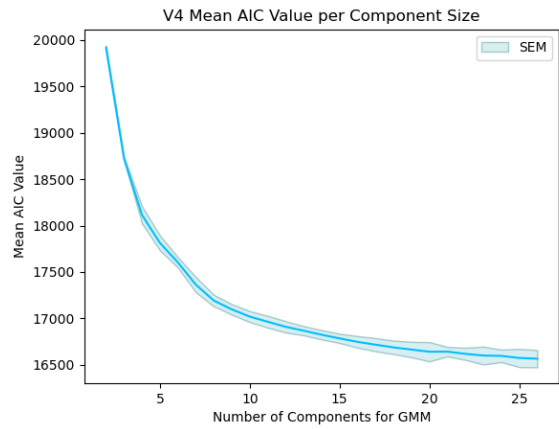
(a)



(b)

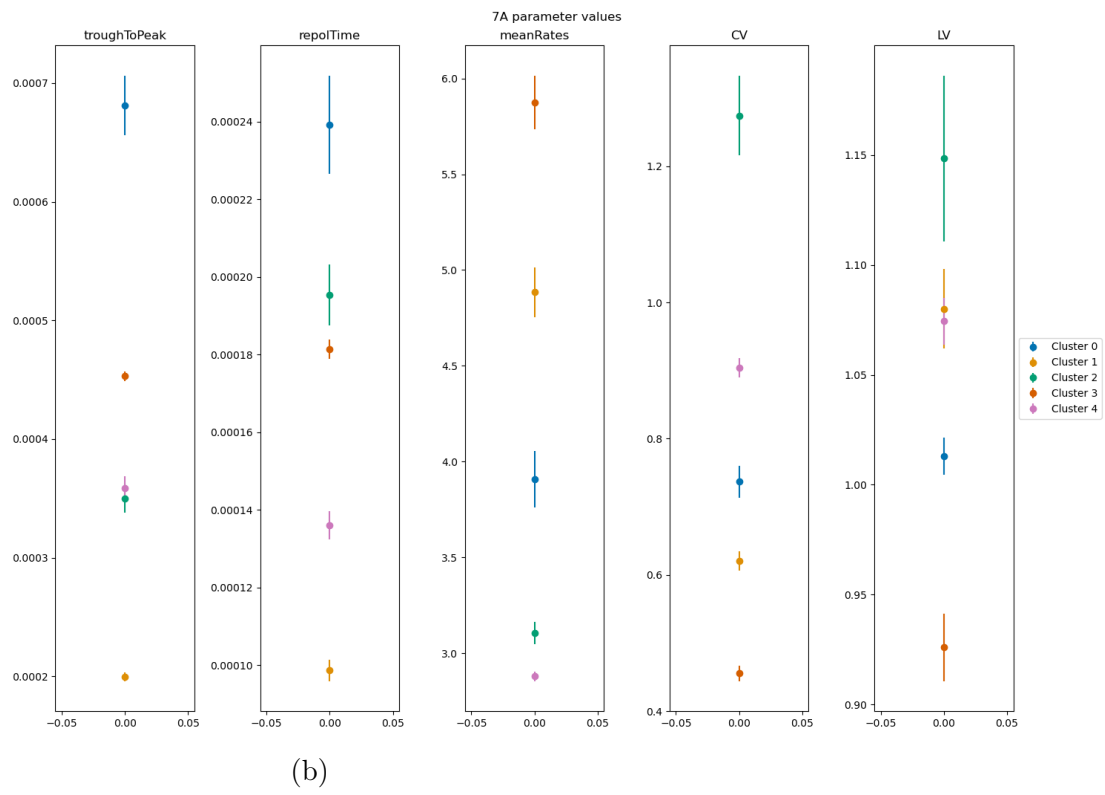
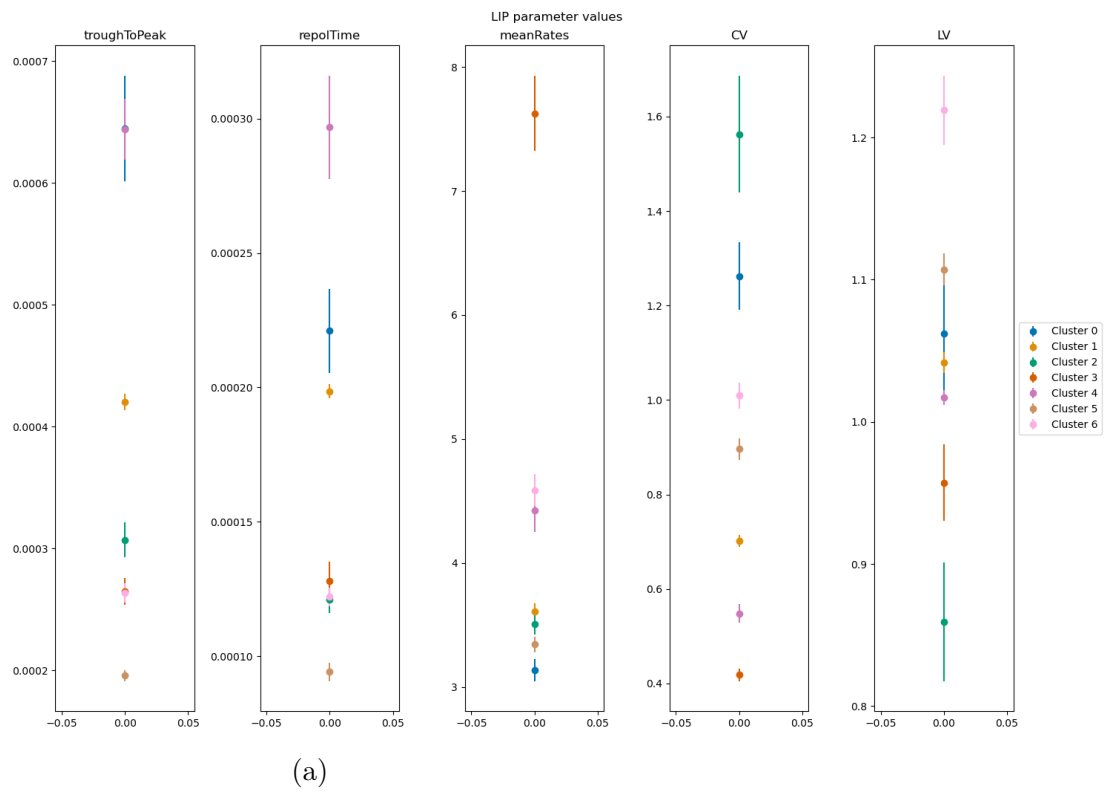


(c)



(d)

Figure 3.3: Mean AIC values and SEM values (shown as blue highlight) for each area (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



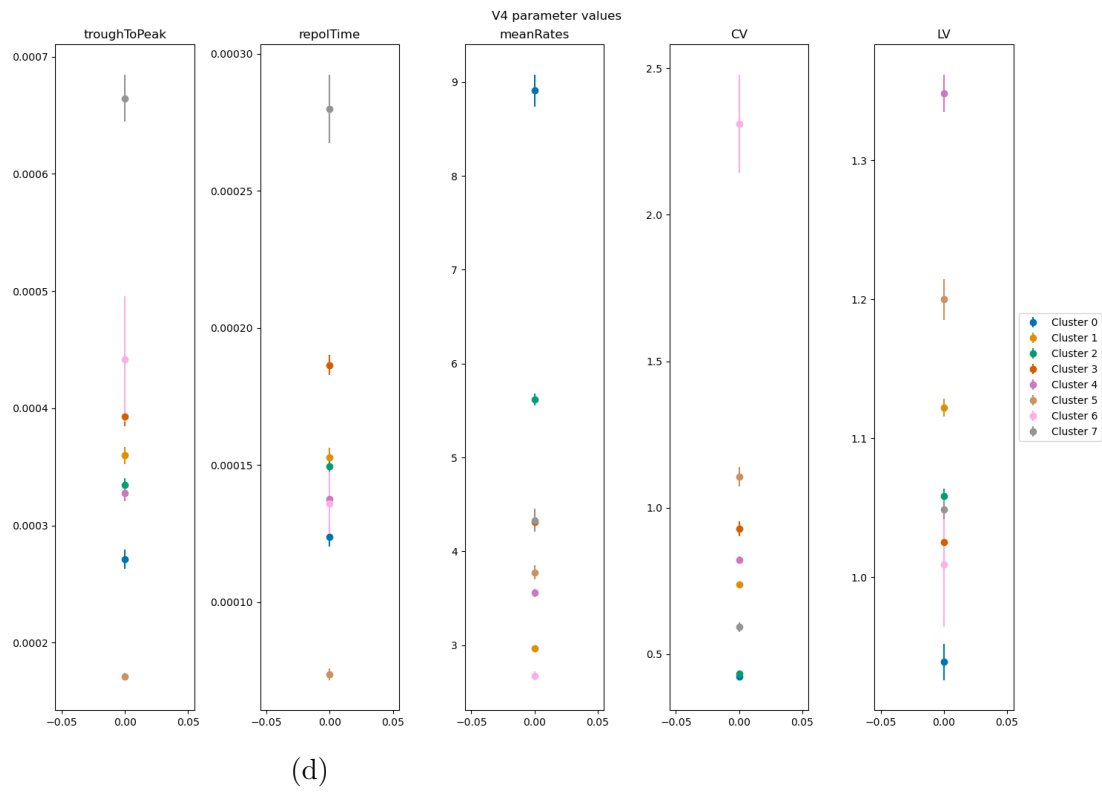
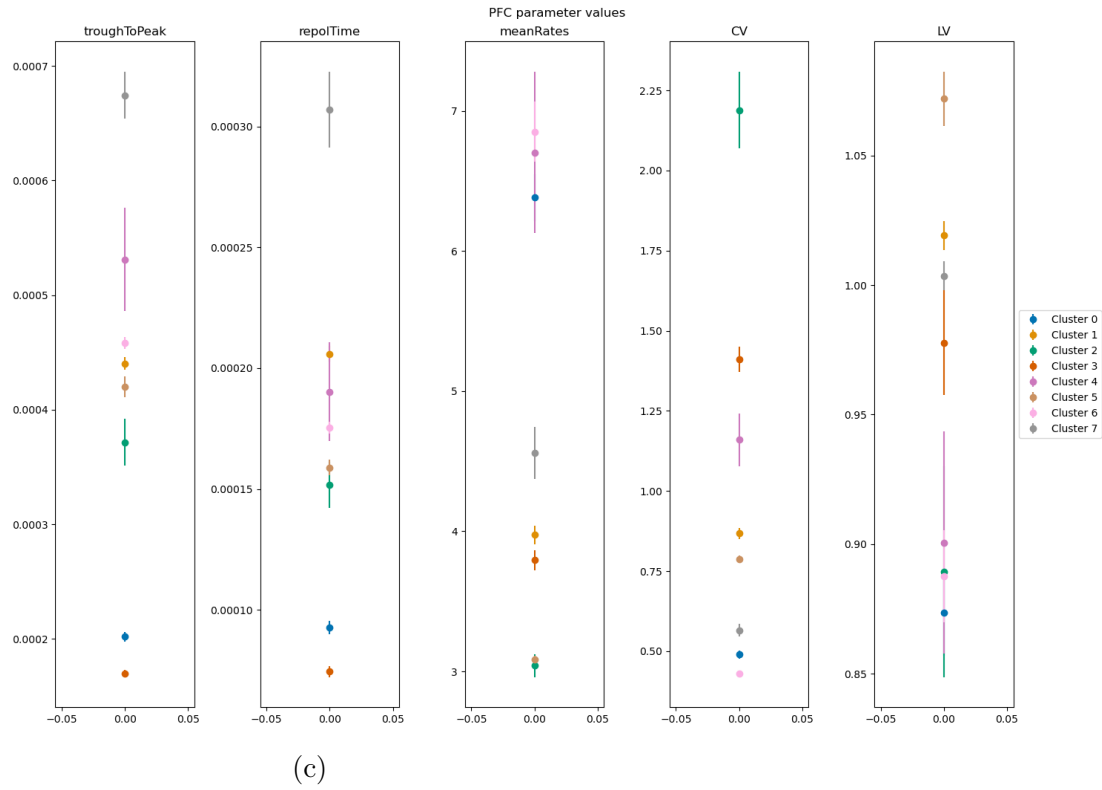
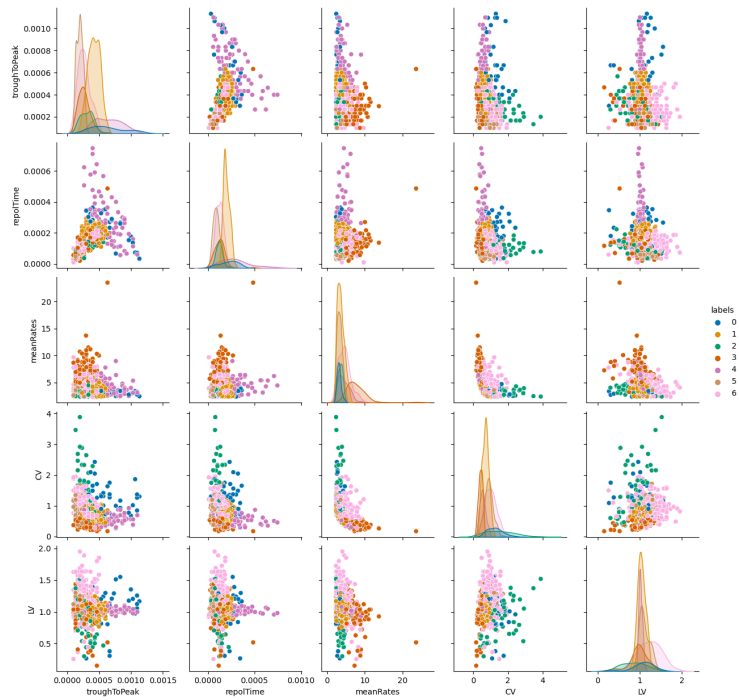
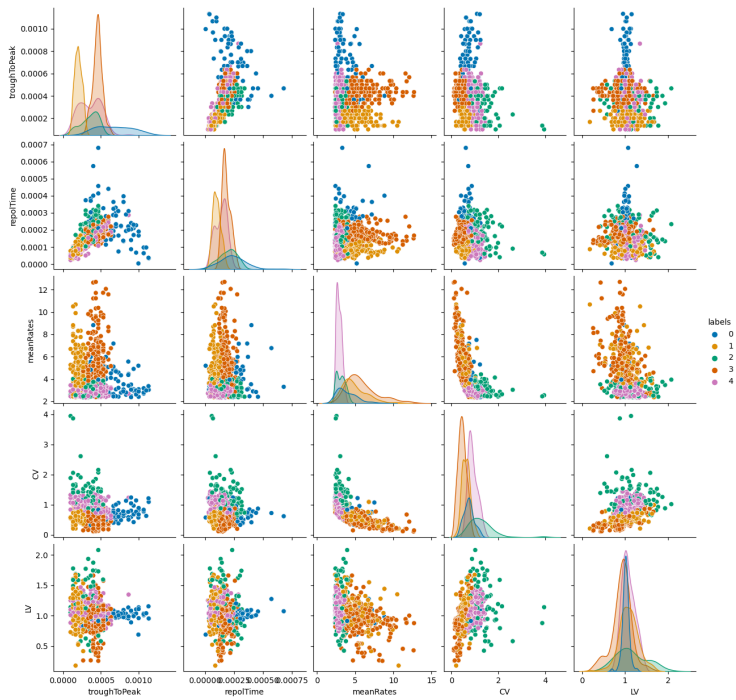


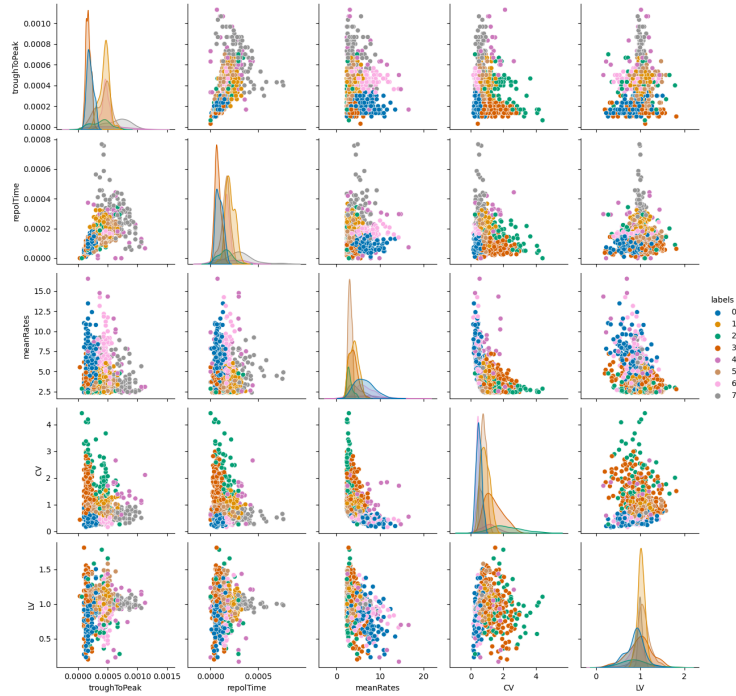
Figure 3.4: Five parameter values defining each cell class in areas (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



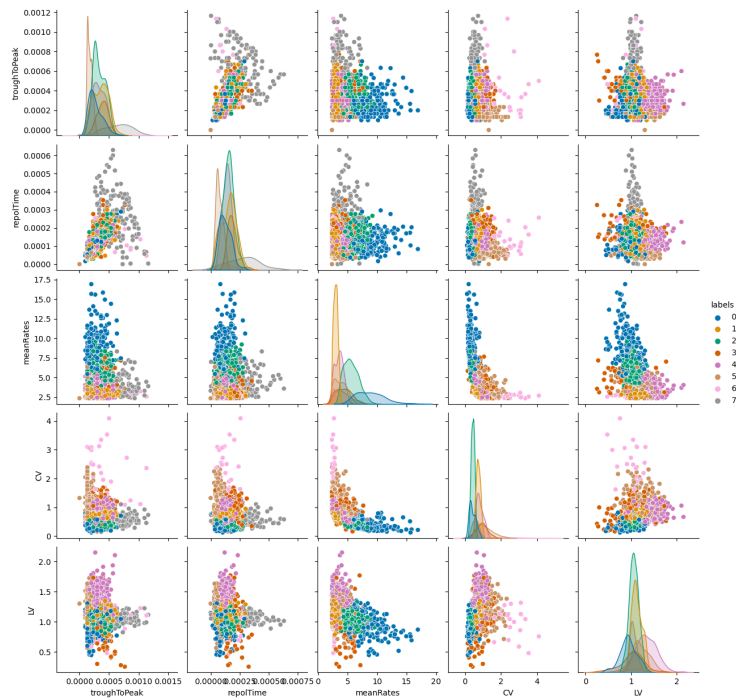
(a)



(b)



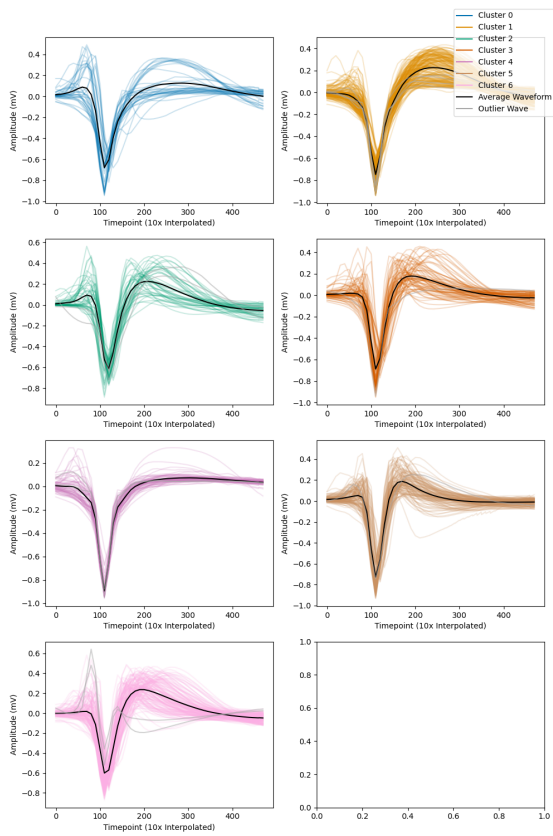
(c)



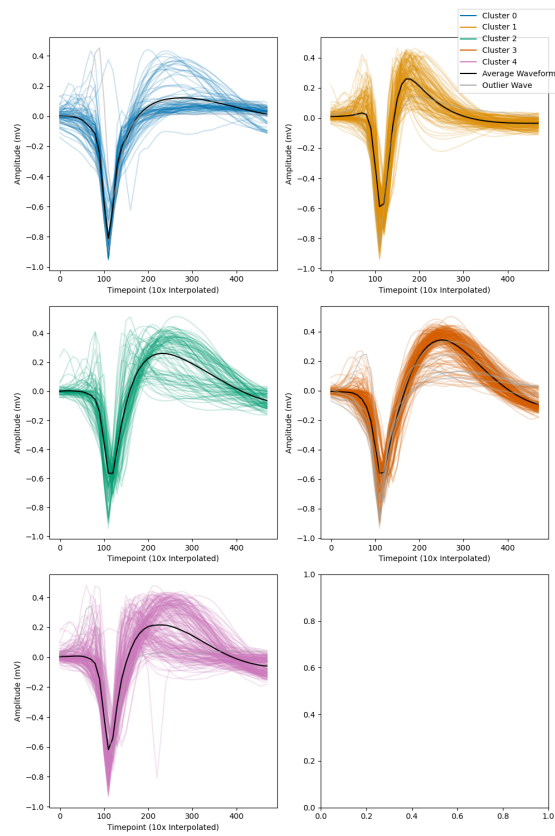
(d)

Figure 3.5: Visualization of clusters in the 2D space of each pair of features. Marginal distribution of each feature on the diagonal.

(a) LIP ; (b) 7A ; (c) PFC ; (d) V4



(a)



(b)

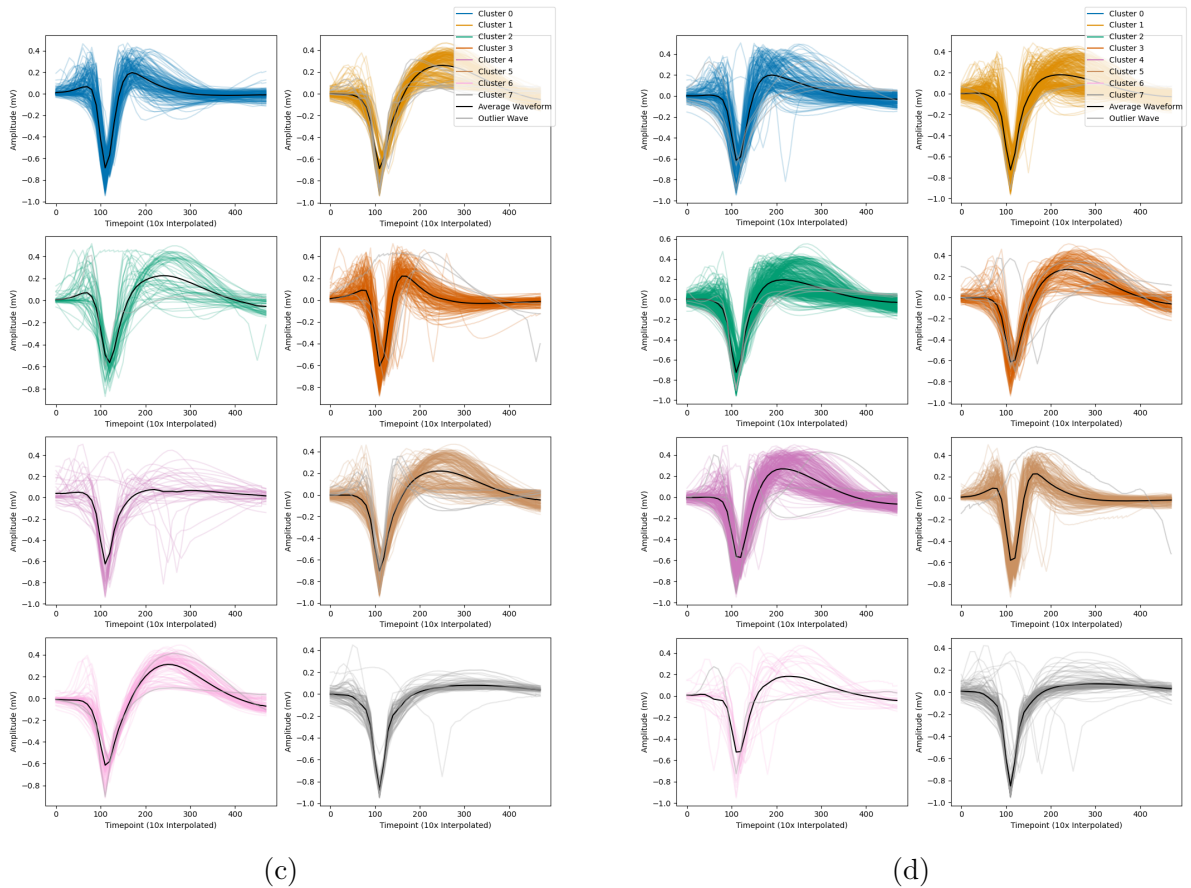
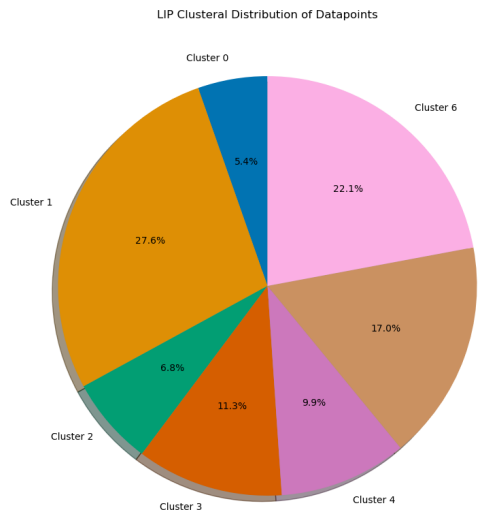
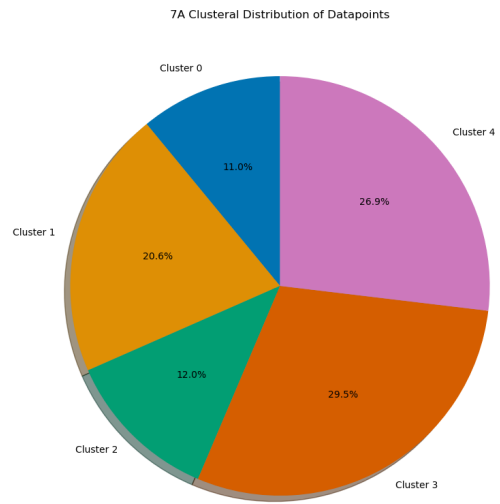


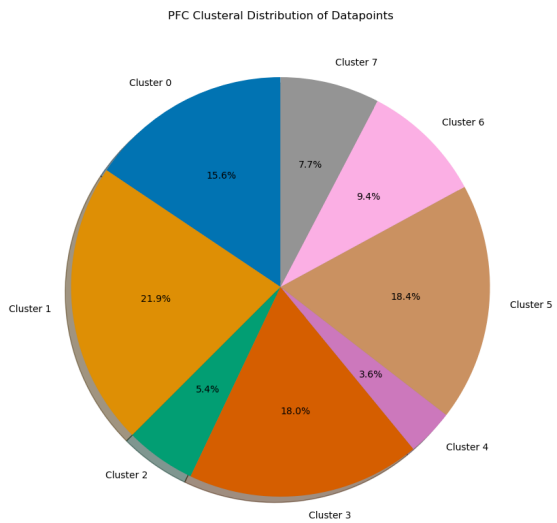
Figure 3.6: Average waveforms per preliminary clustering area. Outlier waveforms shown in light grey. Average waveform shown in black.
 (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



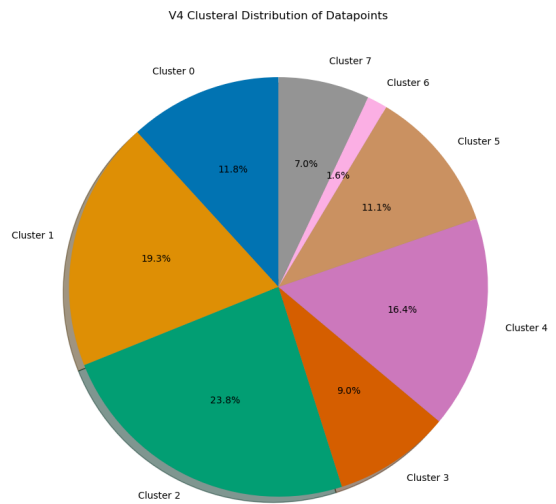
(a)



(b)

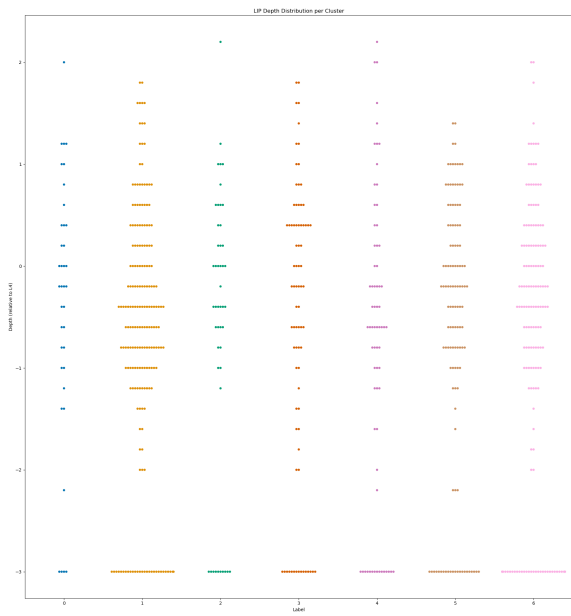


(c)

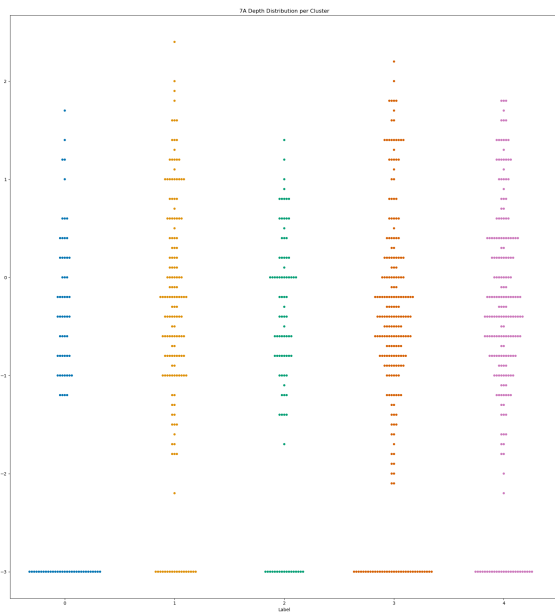


(d)

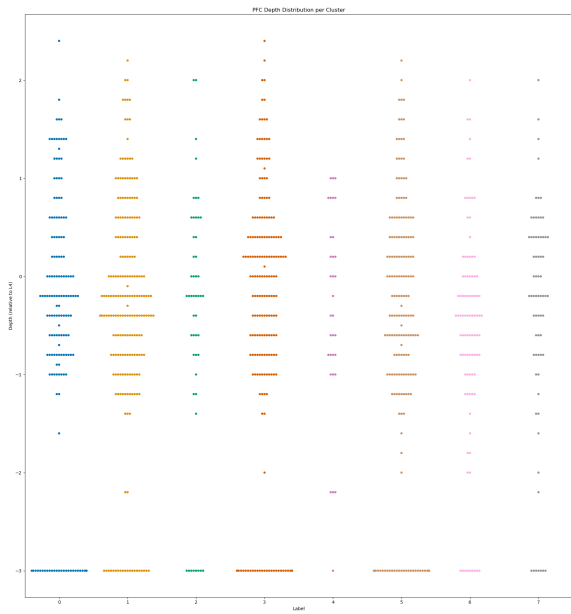
Figure 3.7: Percentage of data points in each cluster for areas
(a) LIP ; (b) 7A ; (c) PFC ; (d) V4



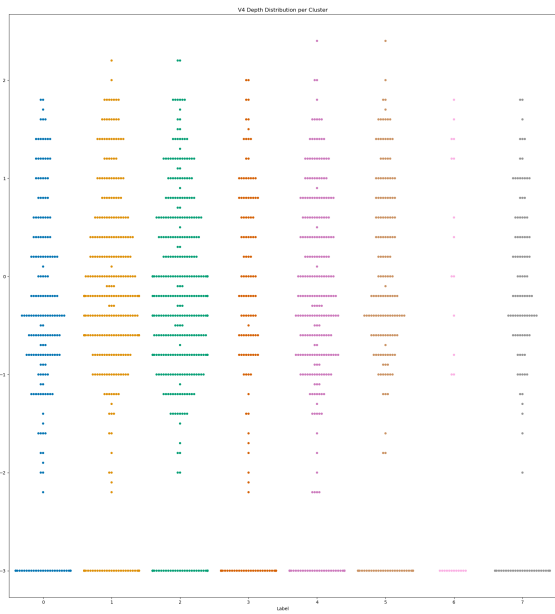
(a)



(b)

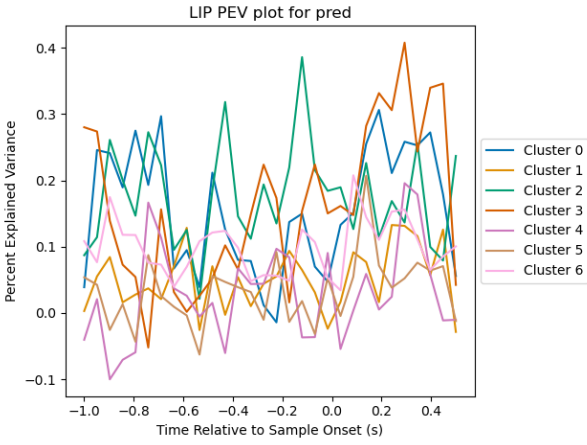


(c)

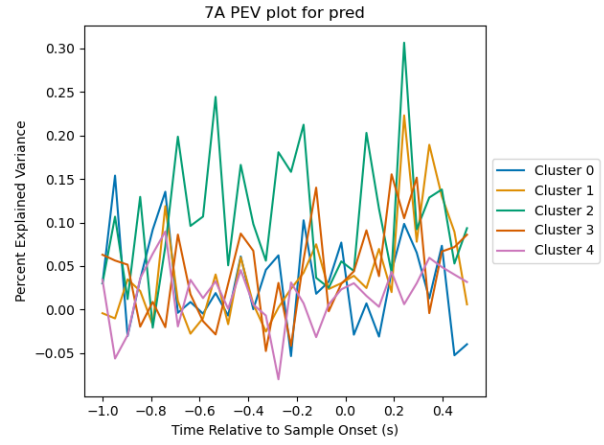


(d)

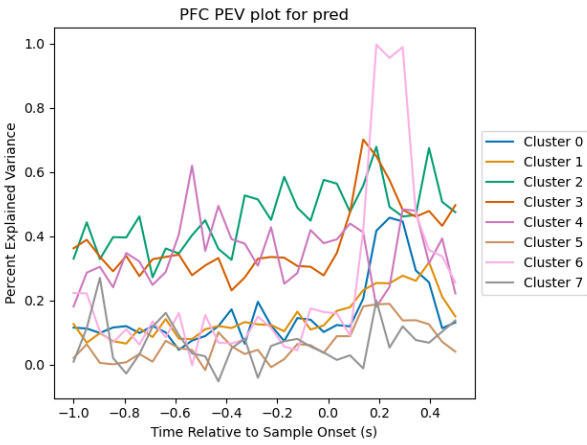
Figure 3.8: Depth distribution for clusters in areas
 (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



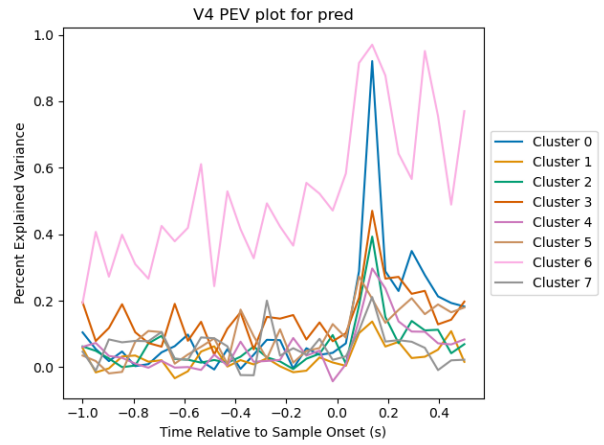
(a)



(b)

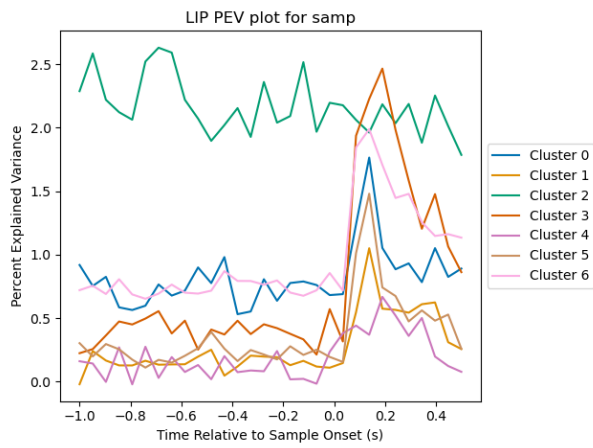


(c)



(d)

Figure 3.9: Percent explained variance of trial predictability on cluster firing rate
 (a) LIP ; (b) 7A ; (c) PFC ; (d) V4

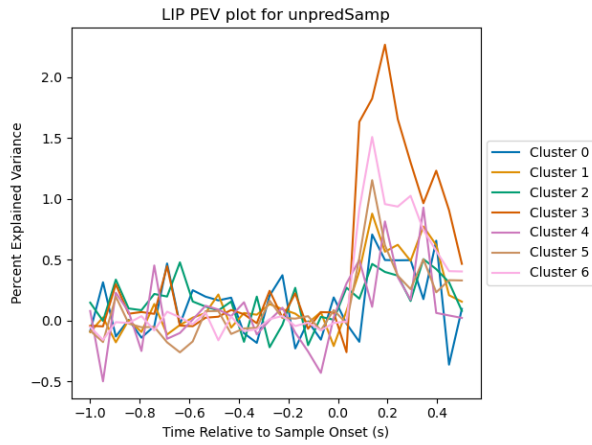


(a)

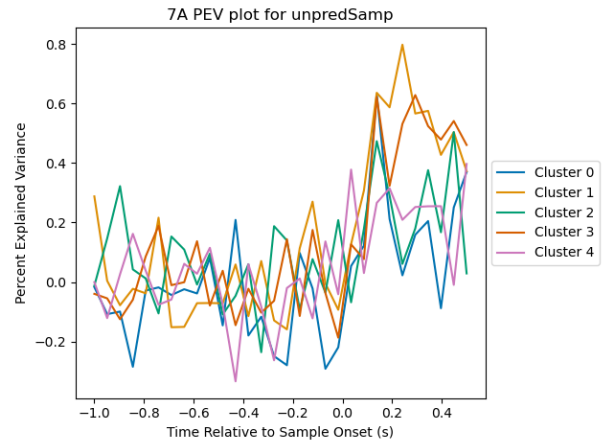


(b)

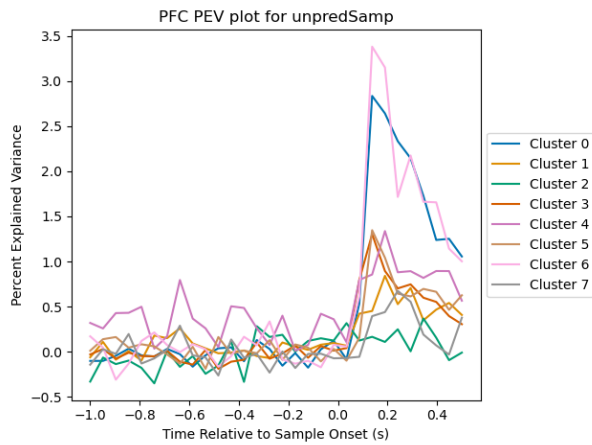
Figure 3.10: (a) PEV of each cluster on sample type in LIP ; (b) Raster plot of neuronal activity over different sample trials in LIP



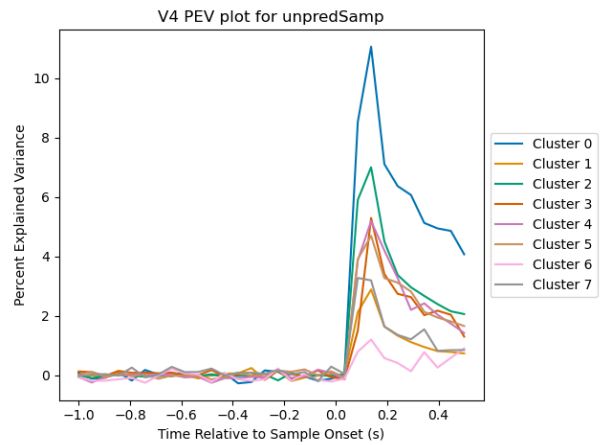
(a)



(b)

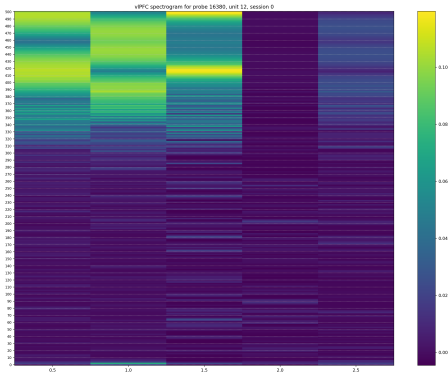


(c)

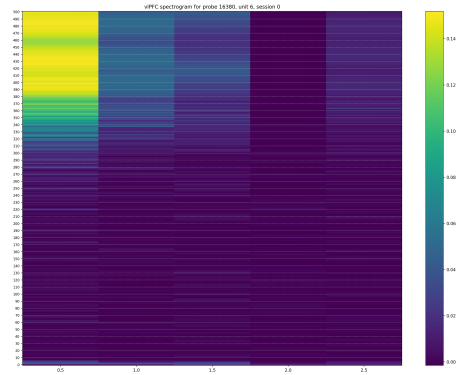


(d)

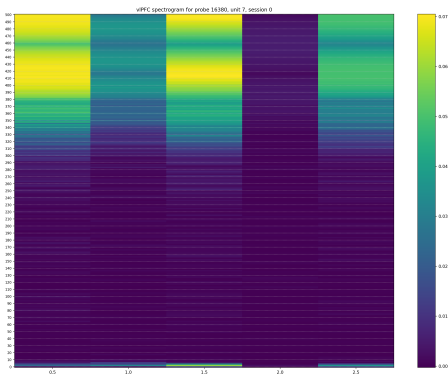
Figure 3.11: Percent explained variance of trial sample type on cluster firing rate, with predictable trial blocks removed
 (a) LIP ; (b) 7A ; (c) PFC ; (d) V4



(a)



(b)



(c)

Figure 3.12: Spike-field coupling synchrony of select PFC units with superficial LFP signals
 (a) Cluster 0 unit (Group N1) ; (b) Cluster 5 unit (Group L1) ; (c) Cluster 7 unit (Group B1)

Chapter 4

Discussion

4.1 Distinct cell classes

Using data in awake macaque monkeys performing a working memory task, five distinct cell classes were found in area 7A, 7 in area LIP, and 8 in areas PFC and V4. These cluster sizes are similar to previous studies, though more varied due to separating the data by area. While the size of the number of cell classes differed between areas, there were certain classes that seemed to be correlates of each other across areas. Identified in this project were group B1 (broadest waveforms, low firing rates, and semi-regular spiking), group N1 (narrowest waveforms, high firing rates, and regular spiking), group C1 (narrower/variable waveforms, low firing rates, non-bursty, but most irregular spiking), and group L1 (variable waveforms, lower firing rates, regular spiking, high burstiness). While these were identified comparing trends in mean feature values across areas, a more robust, computational method is needed to see if these cell classes across areas are actually the same electrophysiological phenotype. This would also be beneficial for identifying the defining features of each class more thoroughly. A potential strategy would be to try validating the results with other unsupervised learning methods, such as the simpler K-means clustering. While GMM is good for handling complex, non-linear data, the results can be harder to interpret and parse.

Using a simpler algorithm to see the comparison in results could offer insight into the clusters given by the GMM.

4.2 Functional and distributional differences between cell classes

Different clusters were shown to be functionally distinct from another, through analysis of the differing firing rate for each cluster due to the difference in sample type, or predictability of the trial. Group N1 and clusters with similar attributes were able to explain variance in predictability in V4 and PFC, as well as sample type in LIP, PFC and V4. This implies that group N1 does a lot of the information coding during mentally demanding working memory tasks such as the DMTS task, in both aspects of differing sample and in prediction. Preliminary spike-field coupling data also showed promise for differing synchrony periods between different cell classes. More investigation into different measures of function would help distinguish the roles of the classes even further. While N1 seems to be important in information coding, it would be interesting to see what role the other classes play, if any at all. It's possible that there's not a lot of functional information in certain areas, such as 7A, which showed little to no effect of both sample type and predictability on any class' firing rate.

Clusters were also shown to be variable in both distribution across areas and depths. In some areas, certain groups would be much more prominent in ratio than others, such as N1 in 7A, or C1 in V4. Cluster depth distribution revealed patterns in certain clusters that were more prevalent in deep or superficial layers. The correlation between depth preference and any functional or informational differences between clusters should be investigated.

4.3 Correlates between classification methods

While it's hard to tell if all clusters have an identifiable morphological or molecular correlate, the characteristics of groups B1 and N1, which are probably the closest to the over-arching broad-spiking and narrow-spiking cell classes identified in previous studies, appear to be similar to regular-spiking pyramidal cells and fast-spiking GABAergic chandelier cells [13]. However, there are certain discrepancies that suggest a bigger picture – for example, pyramidal cells are known to make up around 70-85% of all neurons in the brain [14], but the percentages of datapoints the B1 cell classes made up out of all datapoints were all less than 15%. This might be due to the distribution of probes, and where they recorded from, but it seems suspect. As for the other cell classes, which are less easily separable, it may be to match it up with a previous cell class identified through morphological or molecular methods before more stringently pinning down each class' identifying features.

References

- [1] S. R. Cajal, N. Swanson, and L. Swanson, *Histology of the nervous system, vol. 1*, 1995.
- [2] H. Zeng and J. R. Sanes, “Neuronal cell-type classification: Challenges, opportunities and the path forward,” *Nature Reviews Neuroscience*, vol. 18, no. 9, pp. 530–546, 2017.
- [3] S. Ardid, M. Vinck, D. Kaping, S. Marquez, S. Everling, and T. Womelsdorf, “Mapping of functionally characterized cell classes onto canonical circuit operations in primate prefrontal cortex,” *Journal of Neuroscience*, vol. 35, no. 7, pp. 2975–2991, 2015, ISSN: 0270-6474. [Online]. Available: <https://www.jneurosci.org/content/35/7/2975>.
- [4] C. Trainito, C. von Nicolai, E. K. Miller, and M. Siegel, “Extracellular spike waveform dissociates four functionally distinct cell classes in primate cortex,” *Current Biology*, vol. 29, no. 18, pp. 2973–2982, 2019.
- [5] K. Banaie Boroujeni, P. Tiesinga, and T. Womelsdorf, “Interneuron-specific gamma synchronization indexes cue uncertainty and prediction errors in lateral prefrontal and anterior cingulate cortex,” *eLife*, vol. 10, S. Haegens and M. J. Frank, Eds., e69111, Jun. 2021, ISSN: 2050-084X. [Online]. Available: <https://doi.org/10.7554/eLife.69111>.

- [6] C. Constantinidis and P. S. Goldman-Rakic, “Correlated discharges among putative pyramidal neurons and interneurons in the primate prefrontal cortex,” *Journal of neurophysiology*, vol. 88, no. 6, pp. 3487–3497, 2002.
- [7] E. K. Lee, H. Balasubramanian, A. Tsolias, S. U. Anakwe, M. Medalla, K. V. Shenoy, and C. Chandrasekaran, “Non-linear dimensionality reduction on extracellular waveforms reveals cell type diversity in premotor cortex,” *Elife*, vol. 10, e67490, 2021.
- [8] M. Dasilva, C. Brandt, S. Gotthardt, M. A. Gieselmann, C. Distler, and A. Thiele, “Cell class-specific modulation of attentional signals by acetylcholine in macaque frontal eye field,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 40, pp. 20 180–20 189, 2019.
- [9] K. Banaie Boroujeni, M. Oemisch, S. A. Hassani, and T. Womelsdorf, “Fast spiking interneuron activity in primate striatum tracks learning of attention cues,” *Proceedings of the National Academy of Sciences*, vol. 117, no. 30, pp. 18 049–18 058, 2020.
- [10] J. Gu, “Comparative analysis based on clustering algorithms,” in *Journal of Physics: Conference Series*, IOP Publishing, vol. 1994, 2021, p. 012 024.
- [11] J. Jones, “Clustering: Out of the black box,” *Medium*, 2021. [Online]. Available: <https://towardsdatascience.com/clustering-out-of-the-black-box-5e8285220717>.
- [12] S. Brincat, *Spynal*, version 0.1.2, Sep. 2023. DOI: [10.5281/zenodo.8346152](https://doi.org/10.5281/zenodo.8346152). [Online]. Available: <https://github.com/sbrincat/spynal>.
- [13] L. S. Krimer, A. V. Zaitsev, G. Czanner, S. Kroner, G. González-Burgos, N. V. Povysheva, S. Iyengar, G. Barrionuevo, and D. A. Lewis, “Cluster analysis-based physiological classification and morphological properties of inhibitory neurons in layers 2–3 of monkey dorsolateral prefrontal cortex,” *Journal of neurophysiology*, vol. 94, no. 5, pp. 3009–3022, 2005.

- [14] Y. Wang, M. Ye, X. Kuang, Y. Li, and S. Hu, “A simplified morphological classification scheme for pyramidal cells in six layers of primary somatosensory cortex of juvenile rats,” *IBRO reports*, vol. 5, pp. 74–90, 2018.
- [15] E. Huang, *Clustering*, Feb. 2024. [Online]. Available: <https://github.com/ehuas/clustering>.

Appendix A

Code listing

Below is all code written for this project, separated by file/ function. Access to the repository containing this code (<https://github.com/ehuas/clustering>) can be given upon request [15]. Minor adjustments for clean up and documentation may have been made since the writing of this thesis.

A.1 Loading Data

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Wed Nov 30 13:25:19 2022
5
6 @author: huange
7 """
8 from spynal.matIO import loadmat
9 import numpy as np
10 import pandas as pd
11 from analysis import *
12 from preprocessing import preProcessing, featExtract, coupling
13 import os
14
15 def select_area(unit_info, data, area_name):
16     '''
17     Isolates data recorded from a select area.
18
19     Input: unit_info: dataframe containing unit recording
20           information, such as area recorded
```

```

20         spike_times (n_units, n_timepts): ndarray of spike
           timestamps
21         spike_waves (n_units, n_timepts): ndarray of spike
           waveforms
22         area_name: area of interest
23
24         Output: spike times, waves of select area
25
26     '''
27     areas = unit_info['area'].to_numpy()
28     area_idx = np.where(areas == area_name)[0]
29
30     data = data[:, area_idx]
31
32     return data, area_idx
33
34 def shape_data(spike_times, spike_waves = None):
35
36     for i in range(spike_times.shape[0]):
37         for j in range(spike_times.shape[1]):
38             spike = spike_times[i,j]
39             if spike_waves is not None:
40                 wave = spike_waves[i,j]
41                 if len(wave.shape) < 2:
42                     wave = np.expand_dims(wave, axis = 1)
43             if type(spike) != float:
44                 if np.size(spike) != 0:
45                     trunc_spike = np.where((-1 < spike) & (spike < 2))
46                     [0]
47
48                     spike_times[i,j] = np.atleast_1d(spike[trunc_spike
49 ])
50
51                     if spike_waves is not None:
52                         trunc_wave = wave[:, trunc_spike]
53                         spike_waves[i,j] = np.atleast_2d(trunc_wave)
54
55                     else:
56                         spike_times[i, j] = []
57
58                     if spike_waves is not None:
59                         spike_waves[i,j] = np.empty((48,0))
60
61             else:
62                 if -1 < spike < 2:
63                     spike_times[i,j] = [spike]
64
65                     if spike_waves is not None:
66                         spike_waves[i,j] = np.expand_dims(wave, 1)
67
68                     else:
69                         spike_times[i,j] = []
70
71                     if spike_waves is not None:
72                         spike_waves[i,j] = np.empty((48, 0))
73
74

```

```

69     return spike_times, spike_waves
70
71 def load_data(path):
72     spike_times, spike_times_schema, unit_info, trial_info,
73     session_info, spike_waves, spike_waves_schema = \
74     loadmat(path,
75             variables=['spikeTimes', 'spikeTimesSchema', 'unitInfo', '
76                        trialInfo', 'sessionInfo', 'spikeWaves', 'spikeWavesSchema
77                        '],
78             typemap={'unitInfo': 'DataFrame', 'trialInfo': 'DataFrame'})\
79
80     shape_spikes, shape_waves = shape_data(spike_times, spike_waves =
81     spike_waves)
82
83     return shape_spikes, spike_times_schema, unit_info, trial_info,
84     session_info, shape_waves, spike_waves_schema
85
86 def load_osc_data(path):
87     lfp, lfp_schema, electrode_info, spike_times, unit_info = \
88     loadmat(path,
89             variables=['lfp', 'lfpSchema', 'electrodeInfo', 'spikeTimes',
90                        'unitInfo'],
91             typemap={'electrodeInfo': 'DataFrame', 'unitInfo': 'DataFrame'
92                       })\
93
94     spike_times, _ = shape_data(spike_times)
95
96     return lfp, lfp_schema, electrode_info, spike_times, unit_info
97
98 def concat_sessions(paths, area):
99     comb = pd.DataFrame(columns=['meanRates', 'troughToPeak', '
100                                repolTime', 'CV', 'LV'])
101     PEV_samp_concat = np.empty((0, 30))
102     PEV_pred_concat = np.empty((0, 30))
103     align_waves_concat = np.empty((470, 0))
104     depths_concat = np.empty((0, ))
105     pred_concat = np.empty((0, ))
106     samp_concat = np.empty((0, ))
107     unit_count = 0
108
109     PEV_unpredSamp_concat = np.empty((0, 30))
110
111     for path in paths:
112         spike_times, _, unit_info, trial_info, session_info,
113         spike_waves, spike_waves_schema = load_data(path)
114         area_spike_times, area_idx = select_area(unit_info,
115         spike_times, area)
116         area_spike_waves = spike_waves[:, area_idx]
117
118         validTrials, validNeurons, meanRates, ISIs, meanAlignWaves,
119         smpRate, rates, _, _, _, predInfo, sampInfo, depths, _, _
120         = \
121         preProcessing(area_spike_times,

```



```

110         trial_info,
111         session_info,
112         area_spike_waves,
113         spike_waves_schema,
114         unit_info,
115         area,
116         unit_count) \
117
118     if validTrials.size < 2 or len(validNeurons) < 2:
119         pass
120     else:
121         unit_count += len(validNeurons)
122         features = featExtract(meanRates, ISIs, meanAlignWaves,
123                               smpRate, rates)
124         comb = pd.concat([comb, features], ignore_index=True)
125         #pev_samp = pev_func(rates, sampInfo)
126         pev_pred = pev_func(rates, predInfo)
127
128         #PEV_samp_concat = np.concatenate((PEV_samp_concat, np.
129                                           squeeze(pev_samp, axis=0)), axis = 0)
130         PEV_pred_concat = np.concatenate((PEV_pred_concat, np.
131                                           squeeze(pev_pred, axis=0)), axis = 0)
132
133         # pev_unpredSamp = pev_func(rates[trial_trials, :, :],
134                                     unpredSampInfo)
135         # PEV_unpredSamp_concat = np.concatenate((
136             PEV_unpredSamp_concat, np.squeeze(pev_unpredSamp, axis
137                                               =0)), axis = 0)
138
139         align_waves_concat = np.concatenate((align_waves_concat,
140                                               meanAlignWaves), axis = 1)
141
142         depths_concat = np.concatenate((depths_concat, depths),
143                                       axis = 0)
144
145         pred_concat = np.concatenate((pred_concat, predInfo.
146                                     to_numpy()), axis = 0)
147         samp_concat = np.concatenate((samp_concat, sampInfo.
148                                     to_numpy()), axis = 0)
149
150     return comb, PEV_pred_concat, align_waves_concat, depths_concat,
151           pred_concat, samp_concat
152
153 def osc_concat(paths, areas):
154     session = 0
155     for path in paths:
156         lfp, lfp_schema, electrode_info, spike_times, unit_info =
157             load_osc_data(path)
158
159     for area in areas:
160         if unit_info['area'].isin([area]).any():

```

```

150         area_lfp, area_idx = select_area(electrode_info, lfp,
151                                         area)
152         lfp_trunc = area_lfp[1000:4001, :, :]
153         area_spikes = spike_times[:, area_idx]
154
155         smp_rate = lfp_schema['smpRate']
156
157         if '/mnt/common/datasets/wmPredict/mat/mainTask' in
158             path:
159             depth_var = 'betaGammaDepth'
160         else:
161             depth_var = 'laminarDepth'
162
163         coupling(lfp_trunc, area_idx, depth_var,
164                 electrode_info, unit_info, area_spikes, area,
165                 session, smp_rate)
166
167     session += 1
168
169 def main():
170     directories = ['/mnt/common/datasets/wmPredict/mat/mainTask', '/
171                  /mnt/common/scott/laminarPharm/mat']
172     paths = []
173     for directory in directories:
174         for filename in os.listdir(directory):
175             if filename == 'laminarPharm_databases.mat' or filename
176                 == 'spikesOnly' or filename == 'wmPredict_databases.
177                 mat':
178                 pass
179             else:
180                 f = os.path.join(directory, filename)
181                 paths.append(f)
182
183     areas = ['v1PFC', 'd1PFC', '7A', 'V4', 'LIP']
184
185     #comb, PEV_pred, waves, depths, pred, samp = concat_sessions(
186         paths, area)
187
188     osc_concat(paths, areas)
189
190     # unpredSamp_df = pd.DataFrame(PEV_samp)
191     # unpredSamp_df.to_csv('/home/ehua/clustering/090623_data/{}_
192         _PEV_unpredSamp.csv'.format(area))
193
194     # df = pd.DataFrame(comb)
195     # df.to_csv('/home/ehua/clustering/090623_data/{}_df.csv'.
196         format(area))
197
198     # waves_df = pd.DataFrame(waves)
199     # waves_df.to_csv('/home/ehua/clustering/090623_data/{}_waves
200         .csv'.format(area))
201
202     # samp_df = pd.DataFrame(PEV_samp)

```

```

192     # samp_df.to_csv('/home/ehua/clustering/090623_data/{}
        _PEV_samp.csv'.format(area))
193
194     # pred_df = pd.DataFrame(PEV_pred)
195     # pred_df.to_csv('/home/ehua/clustering/090623_data/{}
        _PEV_pred.csv'.format(area))
196
197     # depths_df = pd.DataFrame(depths)
198     # depths_df.to_csv('/home/ehua/clustering/090623_data/{}
        _depths_jitter.csv'.format(area))
199
200     # pred_df = pd.DataFrame(pred)
201     # pred_df.to_csv('/home/ehua/clustering/090623_data/{}_pred.
        csv'.format(area))
202
203     # samp_df = pd.DataFrame(samp)
204     # samp_df.to_csv('/home/ehua/clustering/090623_data/{}_samp.
        csv'.format(area))
205
206
207 if __name__ == "__main__":
208     main()

```

A.2 Pre-processing

```

1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Tue Sep 27 22:52:30 2022
5
6 @author: huange
7 """
8 from spynal import spikes, utils, sync, spectra
9 import numpy as np
10 import pandas as pd
11 from utils import *
12 import copy
13 import math
14 import matplotlib.pyplot as plt
15
16 def trialSelection(trial_info, session_info):
17     '''
18     Selects valid trials from data.
19
20     Input: trialInfo (n_trials, n_variables) DataFrame for single
           session
21     Output: (valid_trials,) vector of indices of trials to keep
22
23     Valid trials are defined as:

```

```

24         (a) correct and
25         (b) <= 5 trials before drug injection onset
26
27     '''
28     ### isolates only trials that are <= 5 trials before drug onset
29     if 'drugStartTrial' in session_info:
30         drugStartTrial = session_info['drugStartTrial']
31         beforeDrugTrials = trial_info.loc[trial_info['trial'] <=
32             drugStartTrial-5]
33         trials_df = beforeDrugTrials.loc[beforeDrugTrials['correct']]
34         trials_keep = np.where(beforeDrugTrials['correct'])[0]
35     else:
36         trials_keep = np.where(trial_info['correct'])[0]
37         trials_df = trial_info.loc[trial_info['correct']]
38
39     sampInfo = copy.deepcopy(trials_df['sample'])
40
41     ### identifies the non predictable and predictable trials
42     block_trials = np.where(trials_df['blockType'] == 'block')[0]
43     trial_trials = np.where(trials_df['blockType'] == 'trial')[0]
44
45     ### sets block types to 1, trial types to 0
46     trials_df.loc[trials_df['blockType'] == 'block'] = 1
47     trials_df.loc[trials_df['blockType'] == 'trial'] = 0
48     predInfo = copy.deepcopy(trials_df['blockType'])
49     unpredSampInfo = sampInfo[trial_trials]
50
51     return trials_keep, predInfo, sampInfo, block_trials,
52         trial_trials, unpredSampInfo
53
54 def trialsKeep(trials_keep, spike_times, spike_waves):
55     '''
56     Filters out non-valid trials in given time, wave data
57     '''
58     times_trials = spike_times[trials_keep, :]
59     waves_trials = spike_waves[trials_keep, :]
60
61     return times_trials, waves_trials
62
63 def neuronSelection(times_trials):
64     '''
65     Selects valid neurons from data.
66
67     Input: spikeTimes (n_trials, n_units) object array for single
68         session
69     Output: (n_units,) bool vector indicating which units to keep
70
71     Valid neurons are defined as:
72     (a) having an overall mean spike rate > 1 spike/s (weeding
73         out silent cells) and
74     (b) having < 0.1% ISIs within 1 ms (weeding out poorly
75         isolated single neurons)
76
77     '''

```

```

72     '''
73     rates, timepts = spikes.rate(times_trials, method='bin', lims =
74         [-1, 0.5])
75
76     #takes mean over all trials & timepts
77     meanRates = np.mean(np.mean(rates, axis = 2), axis = 0)
78
79     #find all indices where mean rate > 1
80     neurons_keep = np.where(meanRates > 1)[0]
81
82     valid_spikes = times_trials[:, neurons_keep]
83
84     allISIs = spikes.isi(valid_spikes)
85     ms_ISI = np.multiply(allISIs, 1000)
86     concat = utils.concatenate_object_array(ms_ISI, 0)
87     if concat.size > 1:
88         #every ISI per unit
89         flatAll = np.squeeze(utils.concatenate_object_array(ms_ISI,
90             0))
91     else:
92         flatAll = concat[0]
93
94     for idx, neuron in enumerate(neurons_keep):
95         #all ISIs less than 1 ms for neuron
96         shorts = np.where(flatAll[idx] <= 1)[0]
97
98         num_shorts = np.size(shorts)
99         total_isi = np.size(flatAll[idx])
100        if num_shorts/total_isi >= 0.1:
101            #don't keep the neuron
102            neurons_keep = np.delete(neurons_keep, idx)
103
104    return neurons_keep
105
106 def neuronsKeep(neurons_keep, times_trials, waves_trials):
107     '''
108     Filters out non-valid units in given time, wave data
109     '''
110     times_data = times_trials[:, neurons_keep]
111     waves_data = waves_trials[:, neurons_keep]
112
113     return times_data, waves_data
114
115 def depth(neurons_keep, unit_info, jitter = True):
116     '''
117     Saves depth information for two dataests. Depth data is
118     discretized -- jitter can be added to depths for ease of
119     plotting.
120
121     *Scaled Andre's data to match Alex's data
122     '''
123     try:

```

```

121     depths = unit_info['laminarDepth'][neurons_keep]
122 except:
123     ### Scale Andre's data
124     depths = (unit_info['betaGammaDepth'][neurons_keep])/1000
125
126 if jitter:
127     ### Create random vector of jitter
128     x = pd.Series(np.random.uniform(low = -0.02, high = 0.02,
129                                     size = len(depths)))
129
130     ### Add jitter to original depths. Fill nan values with -3.
131     depths = x.add(depths.reset_index(drop=True), fill_value =
132                    -3)
132
133 return depths
134
135 def rateData(time_data):
136     '''
137     Returns:
138     (a) meanRates: mean spike rate / unit
139     (b) rates: full np array of spike rates (units x timepts x
140         trials)
141     '''
142     rates, _ = spikes.rate(time_data, method='bin', lims = [-1, 0.5])
143     meanRates = np.mean(np.mean(rates, axis = 2), axis = 0)
144     meanRates = anscombe(np.expand_dims(meanRates, axis=0))
145
146     return meanRates, rates
147
148 def isiData(time_data):
149     '''
150     Returns: ISI data from given spike data. On the scale of
151         milliseconds (scaled by 1000).
152     '''
153     allISIs = spikes.isi(time_data)
154     ms_ISI = np.multiply(allISIs, 1000)
155
156     return ms_ISI
157
158 def waveAlign(waves_data, spike_waves_schema, trial_subset_indices =
159 None):
160     '''
161     Gets the mean-aligned waveform from data.
162
163     Input: spikeWaves (n_trials, n_units) object array for single
164         session
165         spikeWavesSchema:
166         trials_keep: all valid trials
167         neurons_keep: all valid units
168         trial_subset_indices: optional subset of (valid)
169             trials of interest

```

```

167
168     Output: meanAlignWaves (n_trials, n_units): aligned mean
169           waveforms for each unit
170           smpRate: 10x interpolated sampling rate
171     '''
172
173     if trial_subset_indices: #if we pass in some subset
174         waves_data = waves_data[trial_subset_indices, :]
175
176     n_trials, n_units = np.shape(waves_data)
177     timepts = spike_waves_schema['elemIndex'][0]
178
179     num_timepts = (np.size(timepts)-1)*10
180     meanAlignWaves = np.zeros((num_timepts, n_units))
181
182     for neuron in range(n_units):
183         spikesAll = utils.concatenate_object_array(waves_data[:,
184             neuron], axis = 0, elem_axis = 1)
185         n_timepts, n_spikes = np.shape(spikesAll) #get # of time pts
186         x = np.arange(1, n_timepts+1)
187         xinterp = np.arange(1, n_timepts, 0.1) #keep length, divide
188             step by 10
189         waves_interp = utils.interp1(x, spikesAll, xinterp, axis = 0)
190         meanWave = np.mean(waves_interp, axis=1) #get mean waveform
191             over all spikes
192         meanTroughIdx = np.argmin(meanWave) #get mean trough idx
193
194         for spike_idx in range(n_spikes):
195             spike = waves_interp[:, spike_idx]
196             spikeTroughIdx = np.argmin(spike)
197             diff = spikeTroughIdx - meanTroughIdx
198             newSpike = np.full(np.shape(spike), np.nan)
199             if diff > 0:
200                 #if the spike's trough is shifted ahead of mean trough
201                 newSpike[:-diff] = spike[diff:]
202                 #move it back
203             elif diff < 0:
204                 newSpike[abs(diff):] = spike[:diff]
205             else:
206                 newSpike = spike
207             waves_interp[:, spike_idx] = newSpike
208
209         # new waves_interp with aligned spikes
210         meanAlignWave = np.nanmean(waves_interp, axis=1) #take mean
211             of all spikes
212         meanAlignWaves[:, neuron] = meanAlignWave
213
214     smpRate = spike_waves_schema['smpRate']*10
215
216     return meanAlignWaves, smpRate
217
218 def LV(ISIs):

```

```

215     '''
216     Returns:
217     '''
218
219     n_trials, n_units = np.shape(ISIs)
220     allLV = np.zeros((1, n_units))
221
222     for neuron in range(n_units):
223         neuronLV = np.zeros((n_trials,))
224         for trial in range(n_trials):
225             if len(ISIs[trial, neuron]) <= 1: #if there are no ISIs
226                 to compare against each other
227                 neuronLV[trial] = float('NaN')
228             else:
229                 LV = spikes.isi_stats(ISIs[trial, neuron], stat='LV')
230                 neuronLV[trial] = LV
231             meanLV = np.nanmean(neuronLV, axis=0)
232             allLV[:, neuron] = meanLV
233     return allLV
234
235 def waveform_check(repolTime):
236     '''
237     Returns: passed_neurons (n_units, ): all units with non-inverted
238     waveforms.
239     '''
240
241     passed_neurons = []
242     row, num_neurons = np.shape(repolTime)
243     for i in range(num_neurons):
244         if not math.isnan(repolTime[:, i]):
245             passed_neurons.append(i)
246
247     newRepolTime = repolTime[:, passed_neurons]
248     return passed_neurons
249
250 def spike_i(spikes, predInfo, sampInfo, area, idx):
251     '''
252     Saves the time data for each unit (data is 1 x trials)
253     '''
254     trials, units = np.shape(spikes)
255     for i in range(units):
256         np.save('/home/ehua/clustering/090623_data/spikes/{}_spikes_
257             {}.npy'.format(area, i+idx), spikes[:, i])
258
259         predInfo.to_csv('/home/ehua/clustering/090623_data/info/{}
260             _predInfo_{}.csv'.format(area, i+idx))
261         sampInfo.to_csv('/home/ehua/clustering/090623_data/info/{}
262             _sampInfo_{}.csv'.format(area, i+idx))
263
264 def filterSingleElectrodes(electrode_info, depths, lfp, area_idx,
265     lfp_probe_idx):
266     '''
267     Returns: indices of non-singular electrodes

```



```

262     '''
263     idx_keep = np.where(electrode_info['elecType'][area_idx][
264         lfp_probe_idx] != 'single')[0]
265     depths = depths[idx_keep]
266     lfp = lfp[:, idx_keep, :]
267
268     return depths, lfp
269 def preProcessing(spike_times, trial_info, session_info, spike_waves,
270 spike_waves_schema, unit_info, area, unit_count):
271     ### trial selection
272     trials_keep, predInfo, sampInfo, block_trials, trial_trials,
273     unpredSampInfo = trialSelection(trial_info, session_info)
274     times_trials, waves_trials = trialsKeep(trials_keep, spike_times,
275     spike_waves)
276
277     ### neuron selection
278     neurons_keep = neuronSelection(times_trials)
279     time_data, waves_data = neuronsKeep(neurons_keep, times_trials,
280     waves_trials)
281
282     meanAlignWaves, smpRate = waveAlign(waves_data,
283     spike_waves_schema)
284     repolTime = spikes.waveform_stats(meanAlignWaves, stat='
285     repolarization', smp_rate=smpRate)
286     passed_neurons = waveform_check(repolTime)
287
288     time_data = time_data[:, passed_neurons]
289
290     # if time_data.shape[1] >= 2:
291     #     spike_i(time_data, predInfo, sampInfo, area, unit_count)
292
293     waves_data = waves_data[:, passed_neurons]
294     meanAlignWaves = meanAlignWaves[:, passed_neurons]
295
296     depths = depth(passed_neurons, unit_info)
297
298     meanRates, rates = rateData(time_data)
299     meanNeuronRate = np.mean(rates, axis=0)
300     blockRates = np.mean(rates[block_trials, :, :], axis = 0)
301     trialRates = np.mean(rates[trial_trials, :, :], axis = 0)
302
303     ISIs = isiData(time_data)
304
305     return trials_keep, passed_neurons, meanRates, ISIs,
306     meanAlignWaves, smpRate, rates, meanNeuronRate, blockRates,
307     trialRates, predInfo, sampInfo, depths, unpredSampInfo,
308     trial_trials
309
310 def coupling(area_lfp, area_idx, depth_var, electrode_info, unit_info
311 , spike_times, area, session, smp_rate):
312     probeIDs = electrode_info['probeID'][area_idx].unique()
313

```

```

304     for probeID in probeIDs:
305         lfp_probe_idx = np.where(electrode_info['probeID'][area_idx]
306                                 == probeID)[0]
307         spk_probe_idx = np.where(unit_info['probeID'][area_idx] ==
308                                 probeID)[0]
309
310         depths = electrode_info[depth_var][area_idx][lfp_probe_idx].
311                 to_numpy()
312
313         lfp = area_lfp[:, lfp_probe_idx, :]
314         spk = spike_times[:, spk_probe_idx]
315
316         if depth_var == 'betaGammaDepth':
317             depths, lfp = filterSingleElectrodes(electrode_info,
318                                                 depths, lfp, area_idx, lfp_probe_idx)
319
320             ### get idx of depth - superficial is negative, deep is
321             positive. labels of 0 (layer 4) ignored
322             sup_idx = np.where(depths < 0)[0]
323             deep_idx = np.where(depths > 0)[0]
324
325             lfp_sup = np.squeeze(np.mean(lfp[:, sup_idx, :], axis = 1))
326             lfp_deep = np.squeeze(np.mean(lfp[:, deep_idx, :], axis = 1))
327
328             spike_trains = spikes.times_to_bool(spk, lims=(-1,2))[0]
329
330             _, n_units, _ = spike_trains.shape
331
332             for unit in range(n_units):
333                 unit_spikes = np.transpose(np.squeeze(spike_trains[:,
334                                                         unit, :]))
335                 osc_sup, freqs_sup, timepts_sup, n_sup, phi_sup = \
336                     sync.spike_field_coupling(np.transpose(unit_spikes),
337                                               np.transpose(lfp_sup),
338                                               time_axis = 1,
339                                               smp_rate = smp_rate,
340                                               return_phase = True) \
341
342                 osc_deep, freqs_deep, timepts_deep, n_deep, phi_deep = \
343                     sync.spike_field_coupling(np.transpose(unit_spikes),
344                                               np.transpose(lfp_deep),
345                                               time_axis = 1,
346                                               smp_rate = smp_rate,
347                                               return_phase = True) \
348
349                 np.save('/home/ehua/clustering/090623_data/osc/{
350                         _osc_sup_{}__{}_}'.format(area, session, probeID, unit
351                     ), np.squeeze(osc_sup))
352                 np.save('/home/ehua/clustering/090623_data/osc/{
353                         _osc_deep_{}__{}_}'.format(area, session, probeID,
354                     unit), np.squeeze(osc_deep))
355
356

```

```

347     np.save('/home/ehua/clustering/090623_data/osc/{}_
        _phi_sup_{}_{}_{}'.format(area, session, probeID, unit
        ), phi_sup)
348     np.save('/home/ehua/clustering/090623_data/osc/{}_
        _phi_deep_{}_{}_{}'.format(area, session, probeID,
        unit), phi_deep)
349
350     sup_path = '/home/ehua/clustering/090623_data/figures/{}_
        _spec_sup_{}_{}_{}'.format(area, session, probeID,
        unit)
351     deep_path = '/home/ehua/clustering/090623_data/figures/{}_
        _spec_sup_{}_{}_{}'.format(area, session, probeID,
        unit)
352
353     sup_img, sup_ax = spectra.plot_spectrogram(timepts_sup,
        freqs_sup, np.squeeze(osc_sup), sup_path, area,
        session, probeID, unit)
354     deep_img, deep_ax = spectra.plot_spectrogram(
        timepts_deep, freqs_deep, np.squeeze(osc_deep),
        deep_path, area, session, probeID, unit)
355
356     # np.save('/home/ehua/clustering/090623_data/osc/{}_
        _phi_sup_session_{}_unit_{}.csv'.format(area, session, i),
        phi_sup)
357     # np.save('/home/ehua/clustering/090623_data/osc/{}_
        _phi_deep_session_{}_unit_{}.csv'.format(area, session, i)
        , phi_deep)
358
359 def featExtract(meanRates, ISIs, meanAlignWaves, smpRate, rates):
360     '''
361     Extracts features of interest from data.
362
363     Input: meanRates (n_units): mean spike rates for each unit
364            ISIs (n_trials, n_units): ISIs for each trial and unit
365            meanAlignWaves (n_timepts, n_units): mean aligned wave
366            for each unit
367            smpRate: 10x interpolated sampling rate
368     Output: featuresDF: dataframe containing features of interest
369            meanRates, troughToPeak, repolTime, CV,
370            LV
371
372     '''
373     troughToPeak = spikes.waveform_stats(meanAlignWaves, stat='width'
        , smp_rate=smpRate) #axis is 0?
374     repolTime = spikes.waveform_stats(meanAlignWaves, stat='
        repolarization', smp_rate=smpRate)
375
376     mean_timepts = np.mean(rates, axis=2)
377     CV = spikes.rate_stats(mean_timepts, stat='CV', axis=0) #deal
        with timepts
378     allLV = LV(ISIs)

```

```

379     features = {'meanRates': np.squeeze(meanRates).tolist(), '
                'troughToPeak': np.squeeze(troughToPeak).tolist(), 'repolTime':
                    np.squeeze(repolTime.tolist()), 'CV': np.squeeze(CV).tolist()
                , 'LV': np.squeeze(allLV).tolist()}
380     if np.shape(troughToPeak) == (1, 1):
381         featuresDF = pd.DataFrame(data=features, index = [0])
382     else:
383         featuresDF = pd.DataFrame(data=features)
384
385     return featuresDF
386
387 def main():
388     #validTrials, validNeurons, meanRates, ISIs, meanAlignWaves,
        smpRate, rates = preProcessing(spike_times, trial_info,
        session_info, spike_waves, spike_waves_schema)
389     #featuresDF = featExtract(meanRates, ISIs, meanAlignWaves,
        smpRate, rates)
390     pass
391
392 if __name__ == "__main__":
393     main()

```

A.3 Clustering

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Thu Nov  3 16:14:04 2022
5
6  @author: huange
7  """
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 from sklearn.mixture import GaussianMixture
12 from statistics import mode
13 from sklearn.preprocessing import StandardScaler
14 import pandas as pd
15 from plotting import *
16 from analysis import *
17
18
19 def GMM(features, num_reps, area):
20     '''
21     Performs GMM clustering on data.
22
23     Input: features (n_features, n_datapts): matrix of feature
           values for each datapoint

```

```

24         num_reps: number of times clustering is performed for
                a certain component value
25
26     Output: gmm_min: model fitted using the best number of
                components
27             min_labels: cluster assignments for each data point
28             average_min_comp = number of components used for
                clustering
29
30     '''
31     components = np.arange(2, 27) # 2-9 clusters
32     bics = np.zeros((num_reps, len(components)))
33     aics = np.zeros((num_reps, len(components)))
34     models = np.empty((0, len(components)))
35     labels = []
36     min_comps = []
37
38     for rep in range(num_reps): # what is num_reps
39         min_bic = np.inf
40         rep_models = []
41         rep_labels = []
42         for comp in components: # for each cluster #
43             gmm = GaussianMixture(n_components=comp, random_state=rep
44             )
45             gmm.fit(features)
46             label = gmm.predict(features)
47
48             bic = gmm.bic(features)
49             aic = gmm.aic(features)
50
51             bics[rep, comp-2] = bic
52             aics[rep, comp-2] = aic
53             rep_models.append(gmm)
54             rep_labels.append(label)
55
56             if bic < min_bic:
57                 min_bic = bic
58                 min_comp = comp
59         min_comps.append(min_comp)
60         models = np.append(models, np.array(rep_models).reshape((1,
61         25)), axis=0)
62         labels.append(rep_labels)
63
64     plt.figure(0)
65     bics_mean = np.mean(bics, axis=0)
66     bics_stds = np.std(bics, axis = 0)
67     plt.plot(components, bics_mean, 'k', color='#CC4F1B')
68     plt.fill_between(components, bics_mean-bics_stds, bics_mean+
69     bics_stds,
70                     alpha=0.5, edgecolor='#CC4F1B', facecolor='#
71                     FF9848')
72     plt.title(area + " bics")
73     plt.savefig('/home/ehua/clustering/090623_data/figures/{}_bics.

```

```

    png'.format(area))
70
71 plt.figure(1)
72 aics_mean = np.mean(aics, axis=0)
73 aics_stds = np.std(aics, axis = 0)
74 plt.plot(components, aics_mean, 'k', color='#CC4F1B')
75 plt.fill_between(components, aics_mean-aics_stds, aics_mean+
    aics_stds,
76                 alpha=0.5, edgecolor='#CC4F1B', facecolor='#
    FF9848')
77 plt.title(area + " aics")
78 plt.savefig('/home/ehua/clustering/090623_data/figures/{}_aics.
    png'.format(area))
79
80 average_min_comp = mode(min_comps)
81 print(min_comps)
82 min_model_idx = np.argmin(bics[:, average_min_comp-2])
83 min_model = models[min_model_idx, average_min_comp-2]
84 min_labels = labels[min_model_idx][average_min_comp-2]
85
86 return min_model, min_labels, average_min_comp\
87
88 def main():
89     area = 'dlPFC'
90     feat_df = pd.read_csv('/home/ehua/clustering/090623_data/{}_df.
        csv'.format(area), index_col = 0)
91     waves_df = pd.read_csv('/home/ehua/clustering/090623_data/{}_
        _waves.csv'.format(area), index_col = 0)
92     depths = pd.read_csv('/home/ehua/clustering/090623_data/{}_depths
        .csv'.format(area), index_col = 0)
93     waves = waves_df.to_numpy()
94     waves_ptp = waves.ptp(axis = 0)
95     waves_norm = np.divide(waves, waves_ptp)
96
97     all_params = ['troughToPeak', 'repolTime', 'meanRates', 'CV', 'LV
        ']
98
99     cluster_stats = feat_df[all_params].to_numpy()
100     scaler = StandardScaler()
101     cluster_stats_norm = scaler.fit_transform(cluster_stats)
102
103     _, min_labels, _ = GMM(cluster_stats_norm, 500, area)
104
105     cluster_stats_df = pd.DataFrame(cluster_stats_norm)
106     cluster_stats_df['labels'] = min_labels
107     labels_df = feat_df.copy(deep=True)
108     labels_df['labels'] = min_labels
109     labels_df['depths'] = depths
110
111     labels_df.to_csv('/home/ehua/clustering/090623_data/clusters/{}_
        _labels_df.csv'.format(area))
112     min_labels_df = pd.DataFrame(min_labels)
113     min_labels_df.to_csv('/home/ehua/clustering/090623_data/clusters

```

```

        /{}_labels.csv'.format(area))
114
115 if __name__ == "__main__":
116     main()

```

A.4 Analysis

```

1     #!/usr/bin/env python3
2     # -*- coding: utf-8 -*-
3     """
4     Created on Tue Feb 21 13:06:39 2023
5
6     @author: ehua
7     """
8
9     from spynal import spikes, info, randstats
10    import pandas as pd
11    import numpy as np
12    import matplotlib.pyplot as plt
13
14    def cluster_count(labels, comp_num):
15        counts = np.zeros(comp_num)
16        for label in labels['0']:
17            counts[label] += 1
18        return counts
19
20    def pev_func(data, labels):
21        pev = info.neural_info(data, labels)
22        return pev
23
24    def anova(pev, labels):
25        p = np.squeeze(randstats.one_way_test(pev, labels))
26        #time_vec = np.linspace(-1.5, 2.5, 60)
27        #plt.figure()
28        #plt.plot(time_vec, p)
29        return p
30
31    def ttest(data, labels):
32        plt.figure()
33        comp_num = max(labels)
34        fig, axs = plt.subplots(comp_num, 1)
35        fig.tight_layout()
36        colors = ['xkcd:azure', 'mediumseagreen', 'tab:olive', 'xkcd:lavender']
37        time_vec = np.linspace(-1, 0.5, 30)
38
39        for i in range(comp_num):
40            ax = axs[i]
41            cluster_units = labels == i

```

```

42     p = randstats.one_sample_test(data[cluster_units, :])
43     ax.plot(time_vec, np.squeeze(p), color = colors[i])
44
45
46 def main():
47     allPEVDf = pd.read_csv('/home/ehua/clustering/allPEV_samp_V4.csv'
48                             , index_col = 0)
49     allPEV = allPEVDf.to_numpy()
50
51     labels_df = pd.read_csv('/home/ehua/clustering/V4_labels.csv',
52                             index_col = 0)
53     labels = labels_df['labels']
54
55     p = anova(allPEV, labels)
56     print(p)
57
58     ttest(allPEV, labels)
59
60 if __name__ == "__main__":
61     main()

```

A.5 Plotting

```

1     #!/usr/bin/env python3
2     # -*- coding: utf-8 -*-
3     """
4     Created on Thu Feb 9 15:47:31 2023
5
6     @author: ehua
7     """
8     import matplotlib.pyplot as plt
9     import seaborn as sns
10    import numpy as np
11    from matplotlib.ticker import NullFormatter
12    from scipy import stats
13    import pandas as pd
14    from sklearn.manifold import TSNE
15    from spynal import spikes
16    import math
17    from copy import deepcopy
18    from analysis import cluster_count
19
20    def outlier_id(labels_df, comp_num):
21        rows, _ = labels_df.shape
22        max_std = 2.5
23
24        outlier_col = np.empty((rows, ))
25        for i in range(comp_num):

```



```

26     cluster_units = labels_df.loc[labels_df['labels'] == i]
27     cluster_idx = cluster_units.index
28
29     troughToPeak = cluster_units['troughToPeak'].to_numpy()
30     repolTime = cluster_units['repolTime'].to_numpy()
31     ttp_std = np.std(troughToPeak)
32     ttp_mean = np.mean(troughToPeak)
33     ttp_std_filter = np.where(np.logical_or(troughToPeak <
34         ttp_mean - (max_std*ttp_std), troughToPeak > ttp_mean + (
35             max_std*ttp_std)), 1, 0)
36     rpt_std = np.std(repolTime)
37     rpt_mean = np.mean(repolTime)
38     rpt_std_filter = np.where(np.logical_or(repolTime < rpt_mean
39         - (max_std*rpt_std), repolTime > rpt_mean + (max_std*
40             rpt_std)), 1, 0)
41
42     outlier_all = np.logical_or(ttp_std_filter, rpt_std_filter)
43
44     outlier_col[cluster_idx] = outlier_all
45
46 labels_df['outliers'] = outlier_col
47 return labels_df
48
49 def pairplot(labels_df, outliers_df, comp_num, area, outlier=False):
50 all_params = ['troughToPeak', 'repolTime', 'meanRates', 'CV', 'LV
51     ']
52 if outlier:
53     for i in range(comp_num):
54         cluster_units = outliers_df.loc[outliers_df['labels'] ==
55             i]
56         g = sns.pairplot(cluster_units, hue = "outliers", kind='
57             scatter',
58                 diag_kind='kde', palette = 'muted
59                 ', x_vars = all_params,
60                 y_vars = all_params)
61         g.fig.suptitle(area + " outlier pairplot for comp " + i,
62             y = 1.03, fontsize = 20)
63 else:
64     g = sns.pairplot(labels_df, hue = "labels", kind='scatter',
65         diag_kind='kde', palette = 'muted',
66         x_vars = all_params, y_vars =
67             all_params)
68     g.fig.suptitle(area + " cluster pairplot", y = 1.03, fontsize
69         = 20)
70     plt.savefig('/home/ehua/clustering/090623_data/figures/{
71         }_pairplot.png'.format(area))
72
73
74 def plot_avg_wave(allAlignWaves, df, comp_num, area):
75     fig, axs = plt.subplots(comp_num, 1)
76     fig.tight_layout()
77     fig.set_figheight(15)
78     fig.set_figwidth(10)

```

```

65
66     clusters = []
67
68     colors = sns.color_palette("muted")
69     for i in range(comp_num):
70         ax = axs[i]
71         clusters.append(i)
72         cluster_units = df.loc[df['labels'] == i]
73         cluster_units_idx = cluster_units.index
74         outlier_idx = cluster_units.loc[cluster_units["outliers"] ==
75             1].index
76         cluster_units_idx = list(set(cluster_units_idx) - set(
77             outlier_idx))
78         cluster_waves = allAlignWaves[:, cluster_units_idx]
79         outlier_waves = allAlignWaves[:, outlier_idx]
80
81         ax.plot(cluster_waves, color = colors[i], alpha = 0.2)
82         mean_wave = np.mean(cluster_waves, axis = 1)
83         ax.plot(mean_wave, color = 'k')
84         if outlier_waves.size != 0:
85             ax.plot(outlier_waves, color="crimson", alpha = 0.5)
86     fig.suptitle(area + " cluster waveforms", y = 1.03, fontsize =
87         20)
88     fig.legend(clusters)
89
90     plt.savefig('/home/ehua/clustering/090623_data/figures/{
91         _avg_waves.png'.format(area))
92
93
94
95
96 def elbow_plot(components, data_mean, data_std):
97     plt.figure(0)
98     plt.plot(components, data_mean, 'k', color='#CC4F1B')
99     plt.fill_between(components, data_mean-data_std, data_mean+
100         data_std,
101         alpha=0.5, edgecolor='#CC4F1B', facecolor='#
102             FF9848')
103
104 def tsne_plot(ax, perplexity, df_tsne, area):
105     ax.set_title(area + " for Perp=%d" % perplexity)
106     sns.scatterplot(data=df_tsne, x='comp1', y='comp2', marker='o',
107         hue=df_tsne.label.astype('category').cat.codes, ax = ax)
108     ax.xaxis.set_major_formatter(NullFormatter())
109     ax.yaxis.set_major_formatter(NullFormatter())
110     ax.axis("tight")
111
112 def area_dist(df, comp_num, area):
113     colors = sns.color_palette("muted")
114     plt.figure(figsize = (10, 10))
115     datapts, _ = df.shape
116     labels = []
117     percs = []
118     for i in range(comp_num):
119         cluster_pts = df.loc[df['0'] == i].shape[0]

```

```

111     perc = cluster_pts/datapts
112     percs.append(perc)
113     labels.append(str(i))
114
115     plt.pie(percs, labels=labels, autopct='%1.1f%%',
116            shadow=True, startangle=90, colors=colors)
117     plt.title(area + " datapoints per cluster")
118     plt.savefig('/home/ehua/clustering/090623_data/figures/{}_dist.
119                png'.format(area))
120
121 def param_values(df, all_params, comp_num, area):
122     fig, axs = plt.subplots(1, len(all_params))
123     fig.tight_layout()
124     fig.set_figheight(10)
125     fig.set_figwidth(15)
126
127     colors = sns.color_palette("muted")
128     clusters = []
129
130     for i in range(len(all_params)):
131         ax = axs[i]
132         param = all_params[i]
133         param_df = df[[param, 'labels']]
134         for j in range(comp_num):
135             clusters.append(j)
136             cluster_units = param_df.loc[param_df['labels'] == j]
137             mean_param = np.mean(cluster_units[param])
138             sem_param = stats.sem(cluster_units[param])
139             ax.errorbar(x=0, y=mean_param, yerr=sem_param, fmt='o',
140                       color = colors[j])
141             ax.set_title(param)
142
143     plt.legend(clusters)
144     plt.title(area + " parameter values")
145     plt.savefig('/home/ehua/clustering/090623_data/figures/{}_
146                _param_values.png'.format(area))
147
148 def psth(df, comp_num, blockRates, trialRates):
149     time_vec = np.linspace(-1, 1.95, 60)
150     fig, axs = plt.subplots(comp_num, 1)
151     fig.tight_layout()
152     fig.set_figheight(15)
153     fig.set_figwidth(10)
154
155     for i in range(comp_num):
156         ax = axs[i]
157         cluster_units = df.loc[df['labels'] == i]
158         cluster_units_idx = cluster_units.index
159
160         ax.plot(time_vec, np.mean(blockRates[cluster_units_idx, :],
161                                   axis = 0), color="crimson", linewidth = 2)
162         ax.plot(time_vec, np.mean(trialRates[cluster_units_idx, :],
163                                   axis = 0), color="cyan", linewidth = 2)

```

```

159     fig.legend(["block", "trial"])
160
161
162 def pev_plot(data, labels, comp_num, area, label_type):
163     plt.figure()
164     #fig, axs = plt.subplots(comp_num, 1)
165     #fig.tight_layout()
166     colors = sns.color_palette("muted")
167     clusters = []
168     time_vec = np.linspace(-1, 0.5, 30)
169
170     for i in range(comp_num):
171         #ax = axs[i]
172         clusters.append(i)
173         cluster_units = labels == i
174
175         data_mean = np.mean(data[cluster_units.to_numpy().flatten(),
176                               :], axis = 0)
177         plt.plot(time_vec, data_mean, color=colors[i])
178
179         sems = stats.sem(data)
180         plt.fill_between(time_vec, data_mean-sems, data_mean+sems,
181                          alpha=0.5, edgecolor='#CC4F1B', facecolor='#
182                               FF9848')
183
184     plt.title(area + ' PEV plot for ' + label_type)
185     plt.savefig('/home/ehua/clustering/090623_data/figures/{}_PEV_{}.
186                 png'.format(area, label_type))
187     plt.legend(clusters)
188
189 def feat_reduction(df, min_labels, area):
190     '''
191     Reduces N-d data to a 2-d feature space using TSNE method.
192
193     Input: df (n_features, n_datapts): dataframe of features to
194            be reduced
195            min_labels: cluster assignments for each datapoint
196            area: cortical area of data
197
198     Output: scatterplot of data in 2-d space.
199
200     '''
201     num_pts, num_vars = df.shape
202     perplexities = [10, 30, 40, 50, 60, 70, 80, 100]
203     (fig, subplots) = plt.subplots(2, 4, figsize=(16, 8))
204     axes = subplots.flatten()
205
206     for i, perplexity in enumerate(perplexities):
207         ax = axes[i]
208
209         tsne = TSNE(
210             n_components=2,
211             init="random",

```

```

208         perplexity=perplexity,
209         learning_rate="auto",
210         n_iter = 5000
211     )
212     df_embedded = tsne.fit_transform(df)
213
214     df_tsne = pd.DataFrame(df_embedded, columns=['comp1', 'comp2',
215         ])
216     df_tsne["label"] = min_labels
217
218     tsne_plot(ax, perplexity, df_tsne, area)
219
220 def raster(area, labels, cluster, cond_data, cond_type, unit_num,
221     cluster_units):
222     '''
223     Generates raster plots of spike data per cluster for different
224     PEV labels.
225
226     Input: spike_data (n_units, n_timepts): ndarray of spike
227         timestamps
228         labels (n_units, ): column of cluster labels
229         comp_num: number of clusters
230         cond_data (n_units, ): labels of conditions
231         cond_type: condition type (samp or pred)
232
233     Output: scatterplot of data in 2-d space.
234
235     '''
236     fig, axs = plt.subplots(math.ceil(unit_num/10), 10)
237     fig.suptitle('Raster Plot of Spikes for ' + cond_type + ' in ' +
238         area)
239     cond_name = str(cond_type + 'Info')
240     ax_count = 0
241
242     for i in cluster_units.index:
243
244         spikes_i = np.load('/home/ehua/clustering/090623_data/spikes
245             /{}_spikes_{}.npz'.format(area, i), allow_pickle=True)
246         info_i = pd.read_csv('/home/ehua/clustering/090623_data/info
247             /{}_{}_{}.csv'.format(area, cond_name, i), index_col = 0)
248
249         r, c = divmod(ax_count, 10)
250         ax = axs[r, c]
251
252         spikes_df = pd.DataFrame(spikes_i)
253         info_i.index = spikes_df.index
254         df = pd.concat([spikes_df, info_i], axis = 1)
255
256         if cond_type == 'samp':
257             df = df.sort_values('sample')
258             df = df.reset_index(drop = True)

```

```

254     spikes.plot_raster(df[0], ax = ax)
255
256     last1Trial = df['sample'].where(df['sample']==1.0).
        last_valid_index()
257     last2Trial = df['sample'].where(df['sample']==2.0).
        last_valid_index()
258     last3Trial = df['sample'].where(df['sample']==3.0).
        last_valid_index()
259     ax.axhspan(0, last1Trial, facecolor='b', alpha=0.3)
260     ax.axhspan(last1Trial, last2Trial, facecolor='m', alpha
        =0.3)
261     ax.axhspan(last2Trial, last3Trial, facecolor='y', alpha
        =0.3)
262     #plt.axhline(y = last1Trial, color = 'r', linestyle =
        '-')
263     #plt.axhline(y = last2Trial, color = 'r', linestyle =
        '-')
264
265     else:
266         df = df.sort_values('blockType')
267         df = df.reset_index(drop = True)
268
269         spikes.plot_raster(df[0], ax=ax)
270
271         last1Trial = df['blockType'].where(df['blockType']==1.0).
            last_valid_index()
272         last2Trial = df['blockType'].where(df['blockType']==2.0).
            last_valid_index()
273         plt.axhspan(0, last1Trial, facecolor='b', alpha=0.3)
274         plt.axhspan(last1Trial, last2Trial, facecolor='m', alpha
            =0.3)
275
276         ax_count += 1
277
278     def depth_analysis(depths, labels, counts, area):
279         depths = depths.rename(columns={"0": "depth"})
280         labels = labels.rename(columns={"0": "label"})
281         df = pd.concat([depths, labels], axis=1)
282
283         fig, ax = plt.subplots(figsize=(19,19))
284         df.label = df.label.astype("category")
285         sns.swarmplot(data=df, x="label", y="depth", palette='muted', ax=
            ax)
286
287         # weights = np.ones(len(labels))
288         # for i in range(len(labels)):
289         #     weights[i] = 1/counts[labels['label'][i]]
290
291         # sns.histplot(data=df, x="depth", hue="label", multiple="stack",
            weights = weights, palette='muted')
292
293         plt.title(area + 'Depth Distribution per Cluster')
294         plt.savefig('/home/ehua/clustering/090623_data/figures/{}

```

```

        _cluster_depth_nojitter.png'.format(area))
295
296 def main():
297     areas = ['7A', 'V4', 'LIP', 'PFC']
298     pev_types = ['samp', 'pred']
299
300
301     for area in ['LIP']:
302         #feat_df = pd.read_csv('/home/ehua/clustering/090623_data/{
303             _df.csv'.format(area), index_col = 0)
304         waves_df = pd.read_csv('/home/ehua/clustering/090623_data/{
305             _waves.csv'.format(area), index_col = 0)
306         waves = waves_df.to_numpy()
307         waves_ptp = waves.ptp(axis = 0)
308         waves_norm = np.divide(waves, waves_ptp)
309
310         #spikes_df = pd.read_csv('/home/ehua/clustering/090623_data
311             /{}_spikes.csv'.format(area), index_col = 0)
312         #spikes = spikes_df.to_numpy()
313
314         # jitter_df = pd.read_csv('/home/ehua/clustering/090623_data
315             /{}_depths_jitter.csv'.format(area), index_col = 0)
316         # jitter = jitter_df.to_numpy()
317         # depths_df = pd.read_csv('/home/ehua/clustering/090623_data
318             /{}_depths.csv'.format(area), index_col = 0)
319         # depths = depths_df.to_numpy()
320
321         all_params = ['troughToPeak', 'repolTime', 'meanRates', 'CV',
322             'LV']
323         labels_df = pd.read_csv('/home/ehua/clustering/090623_data/
324             clusters/{}_labels_df.csv'.format(area), index_col = 0)
325         labels = pd.read_csv('/home/ehua/clustering/090623_data/
326             clusters/{}_labels.csv'.format(area), index_col = 0)
327         comp_num = max(labels['0']+1)
328         counts = cluster_count(labels, comp_num)
329
330         #depth_analysis(depths_df, labels, counts, area)
331
332         #feat_reduction(feat_df, labels, area)
333
334         # outliers_df = outlier_id(labels_df, comp_num)
335         # plot_avg_wave(waves_norm, outliers_df, comp_num, area)
336         # pairplot(labels_df, outliers_df, comp_num, area)
337
338         # area_dist(labels, comp_num, area)
339         # param_values(labels_df, all_params, comp_num, area)
340
341         #psth(labels_df, comp_num, allBlockRates, allTrialRates)
342
343     for pev_type in pev_types:
344         pev_df = pd.read_csv('/home/ehua/clustering/090623_data
345             /{}_PEV_{}.csv'.format(area, pev_type), index_col = 0)

```

```

338         pev_data = pev_df.to_numpy()
339
340         pev_plot(pev_data, labels, comp_num, area, pev_type)
341
342         for cond_type in pev_types:
343             cluster = 4
344             cond_data = pd.read_csv('/home/ehua/clustering/090623
                 _data/{}_PEV{}.csv'.format(area, cond_type),
                 index_col = 0)
345
346             cluster_units = labels.loc[labels['0'] == cluster]
347             unit_num = cluster_units.size
348             raster(area, labels, cluster, cond_data, cond_type,
                 unit_num, cluster_units)
349
350 if __name__ == "__main__":
351     main()

```

A.6 Utilities

```

1     #!/usr/bin/env python3
2     # -*- coding: utf-8 -*-
3     """
4     Created on Tue Feb 14 15:22:24 2023
5
6     @author: ehua
7     """
8     import numpy as np
9
10    def anscombe(x):
11        return 2.0*np.sqrt(x + 3.0/8.0)

```