Research article

# Functional clustering of neuronal signals with FMM mixture models

Cristina Rueda [a,b,*], Alejandro Rodríguez-Collado [a]

[a] *Department of Statistics and Operations Research, University of Valladolid, 47011 Valladolid, Spain*
[b] *Mathematics Research Institute of the University of Valladolid (IMUVA), 47011 Valladolid, Spain*

## ARTICLE INFO

## ABSTRACT

The identification of unlabeled neuronal electric signals is one of the most challenging open problems in neuroscience, widely known as Spike Sorting. Motivated to solve this problem, we propose a model-based approach within the mixture modeling framework for clustering oscillatory functional data called MixFMM. The core of the approach is the FMM (Frequency Modulated Möbius) waves, which are non-linear parametric time functions, flexible enough to describe different oscillatory patterns and simple enough to be estimated efficiently. In particular, specific model parameters describe the phase, amplitude and shape of the waveforms. A mixture model is defined using FMM waves as basic functions and gaussian errors, and an EM algorithm is proposed for estimating the parameters. Spike Sorting (SS) has received considerable attention in the literature, and different functional clustering approaches have been considered. We have conducted a fair comparative analysis of the MixFMM with three competitors. Two of them are traditional methods in functional clustering and widely used in Spike Sorting. The third is an approach that has proven superior to many others solving Spike Sorting problems. The datasets used for validation include benchmarking simulated and real cases. The internal and external validation indexes confirm a better performance of the MixFMM on real data sets against the three competitors and an outstanding performance in simulated data against traditional approaches.

## 1. Introduction

The analysis of the electric activity of the neurons is regarded as one of the most practical and effective approaches to studying the nervous system. The electric signals recorded by electrodes register rapid voltage rises, lasting a few milliseconds, called spikes. The voltage returns to the initial baseline level in each spike, describing a single oscillation. The study of spikes is crucial in neuroscience, as they serve as the information units between neurons, and their firing rate and shape determine the cell's morphological, functional, and genetic type. In particular, spikes fired by a particular neuron recorded under similar conditions, such as the electrode's distance and orientation, are assumed to have a specific shape. Spike Sorting (SS) is the collection of techniques to identify spikes corresponding to different neurons. The correct identification of spikes is crucial for studying the connectivity patterns between close-by neurons, [1], relating the firing of certain neurons to the memory process, [2], the treatment of epileptic patients, [3], or the development of high-accuracy brain-machine interfaces, [4], among many other questions. However, SS remains one of the most challenging open problems in neuroscience. The main reasons are the low signal-to-noise ratio, the waveform variability
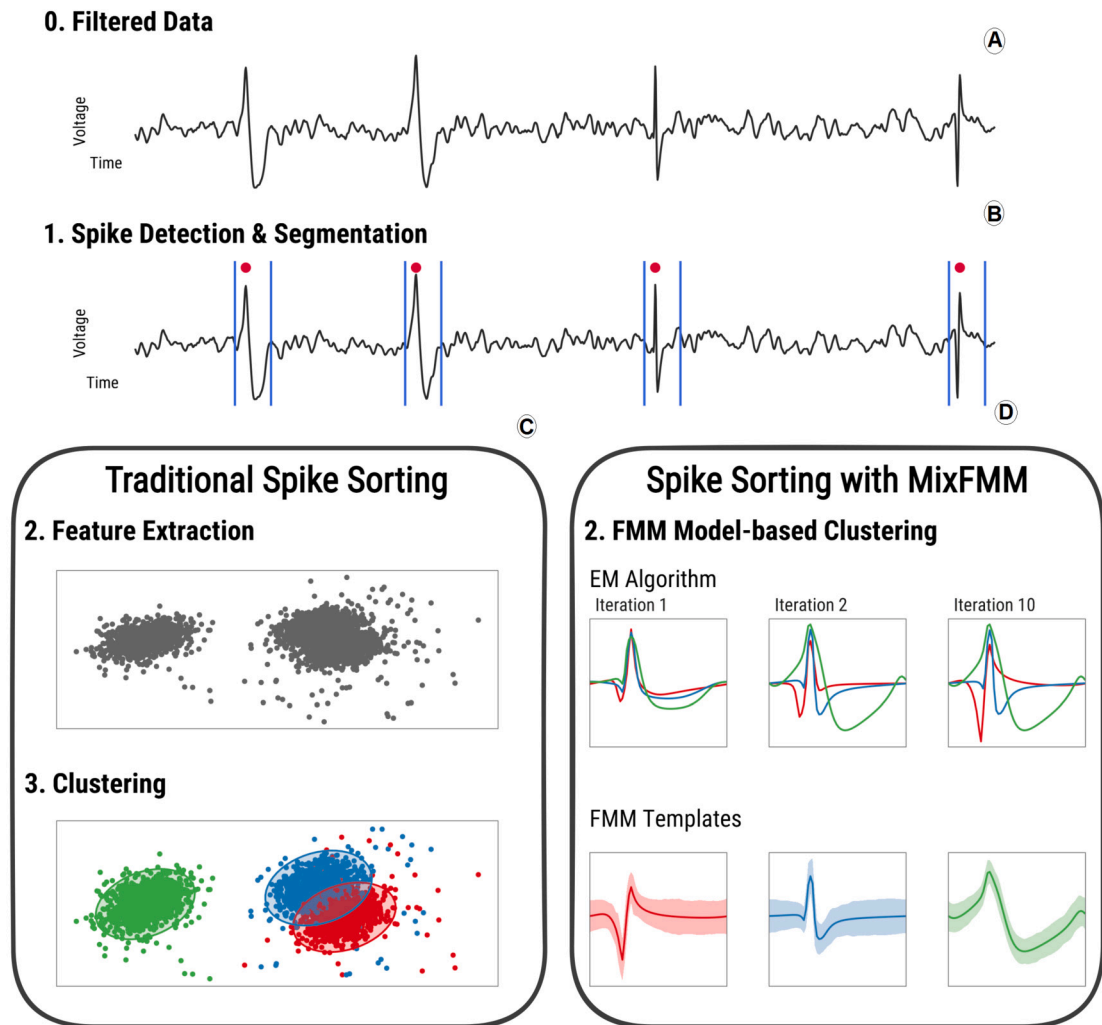
---

**Fig. 1.** Overview of SS: stages 0 and 1 (panels A and B), stages 2 and 3 in traditional SS procedures (panel C), and single model-based clustering stage in the MixFMM model (panel D).

of a particular neuron's spikes, the selection of adequate features to characterize the spike shape, or the overlapping of spikes fired nearby in time. Algorithms that simultaneously address all of these issues are an ever-increasing necessity, as manual spike curation has become unaffordable in light of the increased data availability in the next generation of recording techniques, [2].

The four traditional stages of SS are shown in Fig. 1, top and bottom left. In stage 0, voltage traces are preprocessed, often with a basic bandpass filtering. Stage 1 entails the automatic spike detection and the signal segmentation into various sub-signals, each containing a single spike. Finally, stages 2 and 3 are devoted to the extraction of discriminant features and clustering, respectively. The focus of this paper is on stages 2 and 3, which may be quite different from one approach to another; features are mainly extracted using principal component analysis (PCA), wavelet coefficients or direct measures of geometric features, while the most widely used clustering algorithms are k-means (KM) and Gaussian mixture models (GMM). However, a complete list of alternative techniques used in stages 2 and 3 in SS is large. Some of the most recent proposals are: [5,6,2,7–13], [14,15].

In the vast majority of these papers the authors claim that their method outperforms the others. However, a fair comparison of procedures is a difficult task. It greatly depends on the characteristics of the validation databases, such as the noise level and structure, whether the recordings are from simulated, real extracellular or real intracellular neurons, the waveforms, or whether overlapping spikes are considered, among other aspects. Furthermore, many different clustering approaches have been developed in the literature and there is no clear universal winner. The great attention the topic receives indicates that it is a crucial problem in neuroscience which is still under debate.

In this paper, an original parametric model-based functional clustering approach is proposed as an alternative to stages 2 and 3 in traditional SS (Fig. 1, bottom right).

Model-based functional clustering is also an active topic in the literature [16–21] to cite only a few. Different parametric models have been studied, mainly regression and additive models defined by basic time functions such as splines, Fourier waves, or wavelets.

The model choice depends on the available information, in terms of explicative variables, on the nature of the data, and on the sources of variation characterizing data, among other issues. For instance, Fourier methods might not be suitable for non-stationary signals. Specifically, two main sources of variation are considered in functional analysis, the *amplitude variation* that measures vertical variability, and the *phase variation* measuring horizontal (lateral) displacements: [22–24]. Depending on the problem at hand, either both sources of variation are the focus or one is just nuisance and must be removed. In addition, some methods do not work adequately without prior correction for phase differences.

The model-based approach suitable for analyzing oscillatory signals that we present is the FMM mixture (MixFMM). It is defined by a mixture of gaussian distributions, where the means are sums of FMM waves or components. FMM waves are non-linear time functions that describe a single oscillation with four parameters $(A, \alpha, \omega, \beta)$, characterizing the amplitude $(A)$, phase $(\alpha)$, and shape $(\beta, \omega)$ of the waveform. In particular, the source of variation in amplitude and phase issues can be easily handled using this model. Moreover, the mean cluster waveforms or templates can be compared in a simple way. Models defined with a combination of $m$ FMM waves (FMM$_m$ models) have been successfully used to analyze oscillatory signals arising in different fields, such as ECG (electrocardiogram) and gene expression data, as well as neuronal spikes: [25–27]. In all of these works, the method attains highly accurate predictions and allows the extraction of interpretable features, among other assets. The number of components is often associated with the typical prominent peaks and troughs in the signal, these being five for ECG signals or three for neuronal spikes. The first component, also called the dominant component, usually identifies the prominent underlying biological process and, in some contexts, explains most of the data variability, as is the case of neuronal spikes. Finally, another remarkable property of FMM waves, very useful in functional analysis, is that they are truly dynamic models that can be formulated as an ordinary differential equations system. FMM$_m$ are defined as signal plus error models, where the error is assumed to be gaussian. The maximum likelihood estimators (MLE) of the parameters are derived using a backfitting algorithm, in which FMM$_1$ models are repeatedly fitted to the residue until a stop criterion is attained, [28].

We solve the MixFMM model's parameter estimation with an expectation-maximization (EM) algorithm designed ad-hoc that makes use of the FMM$_m$ estimation algorithm. In addition, we propose a likelihood-based method to select the number of clusters.

The new approach has been validated in 5 benchmarking SS datasets (four simulated and one real) and in an own-created dataset representing a variety of spike waveforms, noise levels, and recording conditions. The quality of clusters generated with the MixFMM is compared, using external and internal indexes, with three clustering approaches. Two of them are traditional methods in functional clustering and widely used in Spike Sorting, which uses PCA to extract the features plus KM or GMM. The third is a recent approach that uses shape, phase, and distribution features plus KM, and which has proven to be superior to other competitors for solving Spike Sorting problems in different datasets in [7]. A first advantage of the MixFMM against the PCA+KM or the PCA+GMM approaches is that the former exploit the underlying oscillatory structure of the signal with specific parameters describing location, amplitude and shape, while the latter do not. Hence, the former approach is more suitable when the spikes differ in shape and amplitude, which is illustrated in the results section. On the other hand, the FMM approach removes noise in a more reasonable way, because it does not remove prominent individual characteristics, while the PCA can do so. If a class that is small in size should contain signals with longer oscillations, then this could be a problem in a spike sorting context. Furthermore, the number of components may be modified, depending on the level of noise and the complexity of the spike morphology. The MixFMM is superior to its competitors in real data sets and achieves outstanding results against traditional approaches in simulated data sets. Furthermore, the proposed method estimates the number of clusters as well as other approaches.

The rest of the paper is organized as follows. Section 2 introduces the MixFMM model and the estimation algorithm. Section 3 describes the datasets analyzed, and Section 4 presents the main results. Finally, in Section 5, concluding remarks and future work are discussed.

## 2. Methods

In the following, let us assume that $\boldsymbol{x_1}, ..., \boldsymbol{x_n}$, $\boldsymbol{x_i} = (x_i(t_1), ..., x_i(t_p))'$, $\forall i \in \{1, ..., n\}$, with $t_0 \leq t_1 < t_2 < ... < t_p \leq T$, are observed signals, also called curves in the paper. Without loss of generality, we assume that $t \in [0, 2\pi)$. Otherwise, consider $t' = \frac{(t - t_0)2\pi}{T}$.

### 2.1. FMM background

An FMM wave is defined as follows,

$$W(t; A, \alpha, \beta, \omega) = A \cos\left(\beta + 2\arctan\left(\omega \tan\left(\frac{t - \alpha}{2}\right)\right)\right).$$

Specifically, $A \in \Re^+$ and $\alpha \in [0, 2\pi)$ are measures of the wave amplitude and phase location, respectively, while $\beta \in [0, 2\pi)$ and $\omega \in [0, 1]$ describe the shape. Fig. 2 shows the waveform patterns for different parameter configurations.

The multicomponent FMM model of order $m$, FMM$_m$, is a gaussian model in which the mean, $\mu(t; \theta)$, is defined as a sum of FMM waves as follows:

$$\mu(t; \theta) = M + \sum_{J=1}^{m} W(t; A_J, \alpha_J, \beta_J, \omega_J) \tag{1}$$

where $\theta = (M, A_1, \alpha_1, \beta_1, \omega_1, ..., A_m, \alpha_m, \beta_m, \omega_m)$ verifies:

1. $M \in \Re$; $(A_J, \alpha_J, \beta_J, \omega_J) \in \Theta_J = \Re^+ \times [0, 2\pi) \times [0, 2\pi) \times [0, 1]$; $J = 1, ..., m$
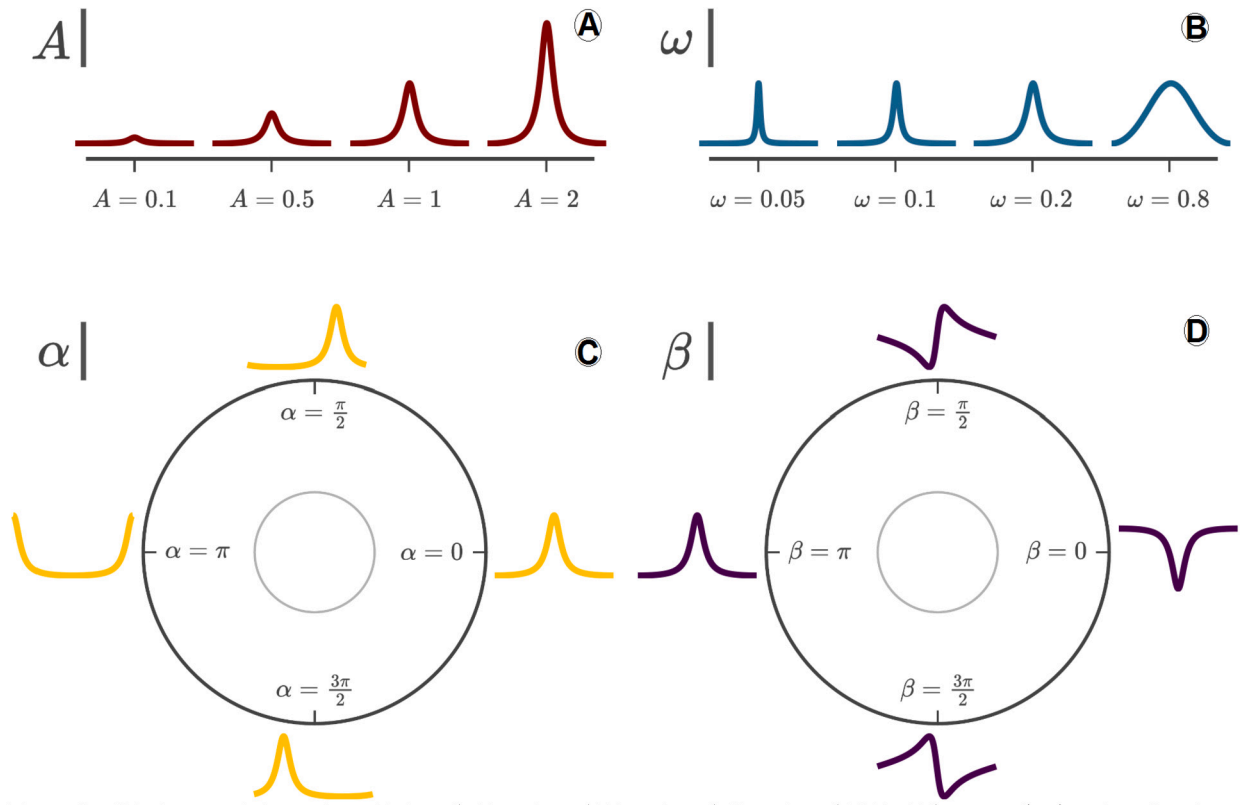
**Fig. 2.** $W(t; A, \alpha, \beta, \omega)$ for various (A (panel A), $\omega$ (panel B), $\alpha$ (panel C), $\beta$ (panel D))'. Unless stated otherwise, $A = 1, \alpha = 0, \beta = \pi, \omega = 0.2$.

2. $\alpha_1 \leq \alpha_2 \leq ... \leq \alpha_{m-1} \leq \alpha_m \leq \alpha_1$.
3. $A_1 = \max_{1 \leq J \leq m} A_J$.

The restrictions guarantee the identifiability of the model parameters and provide biophysiologically interpretable solutions. For instance, the $\alpha$ restrictions correspond to the assumption that the atrial depolarization is previous to the ventricle depolarization in the analysis of ECG signals, and that the cell membrane repolarization precedes its hyperpolarization in neuronal spikes: [25,28].

### 2.2. FMM mixture models (MixFMM)

The MixFMM$_m$ is a mixture model with density defined as follows:

$$f(\boldsymbol{x}|\boldsymbol{\Psi}) = \sum_{k=1}^{K} \gamma_k \, \mathcal{N}(\boldsymbol{x}; \mu(t; \boldsymbol{\theta_k}), \sigma_k^2 \boldsymbol{I_p}), \tag{2}$$

where $\boldsymbol{I_p}$ is the $p \times p$ identity matrix, $\boldsymbol{\Psi} = (\gamma_1, ..., \gamma_K, \boldsymbol{\theta_1}, ..., \boldsymbol{\theta_K}, \sigma_1, ..., \sigma_K)$ is the vector of the model's parameters, and $\gamma_1, ..., \gamma_K$ with $\gamma_k > 0$ and $\sum_{k=1}^{K} \gamma_k = 1$ are the mixture proportions. In addition, $K$ is assumed to be known and corresponds to the number of clusters in our application, while $\mu(t; \boldsymbol{\theta_1}), ..., \mu(t; \boldsymbol{\theta_K})$ are combinations of FMM waves defined as in equation (1).

The log-likelihood of (2) for a sample of size $n$ is given by:

$$\log L(\boldsymbol{\Psi}) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{K} \gamma_k \, \mathcal{N}(\boldsymbol{x_i}; \mu(t; \boldsymbol{\theta_k}), \sigma_k^2 \boldsymbol{I_p}) \right) \tag{3}$$

#### 2.2.1. MLE via EM algorithm

An EM algorithm is designed to find the MLE of the MixFMM model following the methodology in [29,30]. As equation (3) cannot be maximized in a closed form, the complete-data log-likelihood is maximized given the observed data, defined as follows:

$$\log L_c(\boldsymbol{\Psi}) = \sum_{i=1}^{n} \sum_{k=1}^{K} Z_{ik} \log(\gamma_k \, \mathcal{N}(\boldsymbol{x_i}; \mu(t; \boldsymbol{\theta_k}), \sigma_k^2 \boldsymbol{I_p}))$$

where $Z_{ik} \in \{1, ..., K\}, \forall i \in \{1, ..., n\}$ and $1 \leq k \leq K$, is an indicator binary variable such that $Z_{ik} = 1$ if $\boldsymbol{x_i}$ has been generated by the $k$th FMM model. As described below, the EM starts with an initial solution $\boldsymbol{\Psi}^{(0)}$ and alternates iteratively two steps, maximization

and expectation, until convergence is attained. The E-Step computes the expectation of the complete-data log-likelihood given the observed data and the current $\boldsymbol{\Psi}$; whereas, in the M-Step, $\boldsymbol{\Psi}$ is updated to maximize the expectation of the complete-data log-likelihood.

**E-Step**

The expectation of the complete-data log-likelihood, given the observed data and the current parameter vector $\boldsymbol{\Psi}^{(q)}$, is defined as:

$$Q(\boldsymbol{\Psi}; \boldsymbol{\Psi}^{(q)}) = E[\log L_c(\boldsymbol{\Psi})|\boldsymbol{x_1}, ..., \boldsymbol{x_n}, \boldsymbol{\Psi}^{(q)}] = \sum_{i=1}^{n} \sum_{k=1}^{K} \tau_{ik}^{(q)} \log(\gamma_k^{(q)} \mathcal{N}(\boldsymbol{x_i}; \mu(t; \boldsymbol{\theta_k}), \sigma_k^2 \boldsymbol{I_p})) \tag{4}$$

where $\tau_{ik}^{(q)}$ is the posterior probability of the curve $i$ being generated by the $k$th FMM model in the iteration $q$, defined as:

$$\tau_{ik}^{(q)} = P(Z_{ik} = 1|\boldsymbol{x_i}; \boldsymbol{\Psi}^{(q)}) = \frac{\gamma_k^{(q)} \mathcal{N}(\boldsymbol{x_i}; \mu(t; \boldsymbol{\theta_k}), \sigma_k^2 \boldsymbol{I_p})}{\sum_{h=1}^{K} \gamma_h^{(q)} \mathcal{N}(\boldsymbol{x_i}; \mu(t; \boldsymbol{\theta_h}), \sigma_h^2 \boldsymbol{I_p})}, \ \forall i \in \{1, ..., n\}, \ 1 \leq k \leq K \tag{5}$$

**M-Step**

The values of the parameter vector are updated maximizing equation (4):

$$\boldsymbol{\Psi}^{(q+1)} = \mathrm{argmax}_{\boldsymbol{\Psi}} Q(\boldsymbol{\Psi}; \boldsymbol{\Psi}^{(q)}) \tag{6}$$

The estimators of the FMM parameters that solve equation (6) are updated as follows:

$$\boldsymbol{\theta_k^{(q+1)}} = \mathrm{argmin}_{\theta \in \Theta} \sum_{j=1}^{p} (\overline{x}_k^{(q)}(t_j) - \mu(t_j; \boldsymbol{\theta}))^2, \quad 1 \leq k \leq K \tag{7}$$

where $\overline{x}_k^{(q)}(t_j) = \frac{\sum_{i=1}^{n} \tau_{ik}^{(q)} x_i(t_j)}{\sum_{i=1}^{n} \tau_{ik}^{(q)}}, \forall j \in \{1, ..., p\}$, is the mean waveform for cluster $k$ in the iteration $q$. Equation (7) is solved using the standard FMM backfitting algorithm [28].

Then, the estimators of the standard deviation terms are obtained as follows:

$$\sigma_k^{(q+1)} = \sqrt{\frac{1}{p \sum_{i=1}^{n} \tau_{ik}^{(q)}} \sum_{i=1}^{n} \tau_{ik}^{(q)} ||\boldsymbol{x_i}(t) - \mu(t; \boldsymbol{\theta_k^{(q+1)}})||^2}, \quad 1 \leq k \leq K \tag{8}$$

Finally, the cluster proportions are updated as follows:

$$\gamma_k^{(q+1)} = \frac{\sum_{i=1}^{n} \tau_{ik}^{(q)}}{n}, \quad 1 \leq k \leq K$$

The homoscedasticity restriction $\sigma = \sigma_1 = ... = \sigma_K$, often assumed in SS problems, reduces equation (8) to a single one.

Our proposal for the initial values is: $\gamma_k^{(0)} = 1/k, 1 \leq k \leq K$, and $\theta_1^{(0)}, ..., \theta_K^{(0)}$ are the FMM parameters obtained by fitting $FMM_m$ models to the mean waveforms of a random cluster assignation.

Success in terms of convergence to the MLE is not initially guaranteed, although the solution converges to a local minimum. We propose several random initializations to avoid local minimums. We have included numeric studies in the paper that show how the proposal reaches a good solution in practice.

The predicted label for each curve can be obtained by assigning the curve to the cluster with the highest estimated posterior probability. This is denoted as $I_i$, and defined as follows: $I_i = \mathrm{argmax}_{1 \leq k \leq K} \tau_{ik}^{(q)}, \forall i \in \{1, ..., n\}$.

Details concerning the steps of the MixFMM algorithm are given in a flow-chart below.

*2.3. Validity indexes*

Indexes for cluster validation can be classified as external or internal. The former uses externally known information, mainly the class labels; whereas the latter, suitable when the true class labels are unknown, only uses the internal information of the clustering process to evaluate its performance. Accuracy is considered here as the most direct measure when ground truth spikes are available and is the most commonly used criterion across almost all SS applications. As the internal index, we consider the Davies-Bouldin index ($DB$), which is often used in SS applications. The definitions of these two indexes are included below to facilitate the interpretation of the results.

Let $\boldsymbol{D}$ be the $r \times r$ confusion matrix obtained from crossing predicted and real classes, in such a way that the sum of the diagonal frequencies is maximum. The accuracy is defined as:

$$\mathrm{Accuracy} = \frac{\sum_{k=1}^{r} d_{kk}}{n}$$

Now, let $n_k$, $\boldsymbol{c^k} = (c^k(t_1), ..., c^k(t_n))'$, where $c^k(t_j) = \frac{1}{n_k} \sum_{i/I_i=k} x_i(t_j), j \in \{1, ..., p\}$, and $\delta^k = \frac{1}{n_k} \sum_{i/I_i=k} ||\boldsymbol{x_i} - \boldsymbol{c^k}||^2$ denote, respectively, the number of observations, the barycenter and the dispersion of a cluster $k$. The $DB$ is the mean of each cluster's maximum ratio between the sum of the cluster dispersions and the distances between their barycenters:

$$DB = \frac{1}{K} \sum_{k=1}^{K} \max_{k \neq k'} \left( \frac{\delta^k + \delta^{k'}}{||c^k - c^{k'}||} \right).$$

The better the cluster partitions, the higher the expected values of *Accuracy* and the lower the expected values of *DB*.

### 2.4. Alternative procedures

We compare the performance of the novel MixFMM with three procedures.

**PCA+KM**: PCA for feature selection (dimension 4) and KM for cluster generation.

**PCA+GMM**: PCA for feature selection (dimension 4) and GMM for cluster generation.

**SPDF+KM**: a procedure proposed by [7] for feature selection (dimension 24) and KM for cluster generation. The acronym SPDF is for Shape, Phase, Distribution and Feature.

PCA+KM and PCA+GMM are likely the most widely used clustering algorithms in functional clustering, with better results in SS than wavelets, among other approaches, [5,6,13]. The number of principal components selected ranges from two to four in the different SS applications. Here, only the results using four components are shown, as this choice gives better results in these datasets.

SPDF+KM is a recent approach that considers a set of 24 phase, shape, and distribution features and KM to derive clusters. We have chosen SPDF+KM for three main reasons; first because we would like to validate the discriminative power of the set of 24 features in our datasets; second, because it has achieved very good results compared with other methods in [7], and finally, because we have checked that it is also superior to those approaches presented in [14] for the simulated data sets common to the latter and the present paper, eliminating overlapping spikes, which is the setting in [14]. This latter comment also indicates that it is a very difficult competitor to beat and that we are presenting a fair evaluation of our method.

For each of the clustering procedures and datasets, the best model from multiple initializations has been selected.

### 2.5. Selection of the number of clusters

The optimal number of clusters has been estimated using the slope BIC heuristic for PCA+KM, PCA+GMM and SPDF+KM.

For the MixFMM models, we propose a log-likelihood based approach. Specifically, the number of clusters is determined using a slope heuristic approach, similar to those proposed in [31–33]. Specifically, the log-likelihood $\log L(\hat{\Psi})$ values for different numbers of clusters starting from one are calculated and plotted consecutively. An exponential trend is expected at the beginning. The minimum value for which the exponential is reduced to a linear trend is considered the optimal number of clusters.

### 2.6. Programming languages

Data from the Allen Cell Types Database was obtained using Python and the Allen SDK [34]. The rest of the experimentation was developed with R, including spike detection, spike segmentation, and functional clustering. The MixFMM approach has been implemented using the functions from the FMM package [35]. The PCA and KM implementations used are from the package stats [36], whereas the GMM are from mclust [37]. Moreover, the package clusterCrit [38] was used to calculate the internal indexes, while ggplot2 [39] and plotly [40] were used to create the visualizations of the manuscript. The mixFMM method takes around 1 sec to fit an FMM model, although we are working on reducing this time further. Except for this, which takes around K*I (K = number of clusters, I = number of EM iterations) seconds in total for a fixed K, the time spent by the algorithm is similar to that taken by GM or KM.

### 2.7. Flowchart

Fig. 3 presents a diagram showing the sequence of steps of the MixFMM algorithm, which includes the optimal selection of the number of classes. It requests three input parameters: the number of FMM components, $m$; the maximum number of clusters, $Ng$; and the number of iterations of the EM algorithm, $I$. The input data is a voltage time series, while the output is the optimal number of classes and the assignment of each spike to one of those classes.

## 3. Datasets for evaluation

In this work, six different datasets have been analyzed for which the real class labels (ground truth) are known, encompassing a wide variety of real and simulated spike waveforms. The simulated datasets correspond to extracellular recordings provided by [41], and have been used by many other authors as benchmarking datasets since then, such as [8,12,4,13] and [14], among others. Specifically, we consider four sets of data E11, E21, D11, and D21 with the same noise level equal to 0.1 and three different waveforms each. The data have been preprocessed in the same way as in other papers. Specifically, spikes are detected with a threshold, $k = \text{Med}(|z|/0.6745)$ where $z$ is the filtered signal, and the whole voltage traces have been fragmented into signals of 64 samples, each containing a single spike with its maximum aligned with the point 20. The noise artifacts erroneously detected as spikes, less than 4% of the signals, have initially been discarded for the rest of the study. The first real dataset corresponds to four temporal lobe neuron types from Epileptic patients (EPT), which is openly available and provided by [42]. Finally, we also present the analysis of our own created dataset, labeled VGA, with signals from three GABAergic neuron types (Pvalb, Htr3a, and
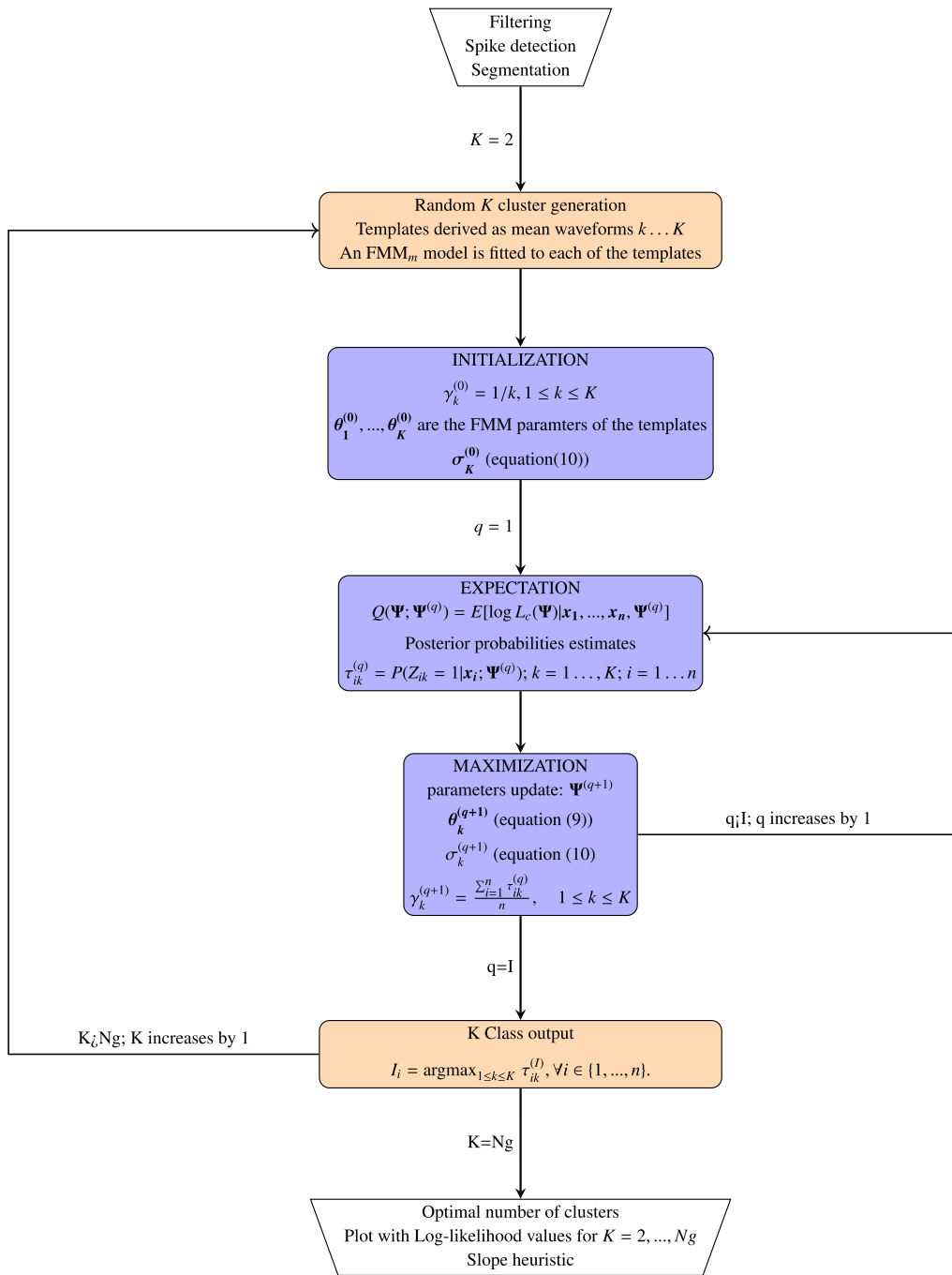
Filtering
Spike detection
Segmentation

$K = 2$

Random $K$ cluster generation
Templates derived as mean waveforms $k \ldots K$
An $FMM_m$ model is fitted to each of the templates

INITIALIZATION

$\gamma_k^{(0)} = 1/k, 1 \le k \le K$

$\boldsymbol{\theta_1^{(0)}}, ..., \boldsymbol{\theta_K^{(0)}}$ are the FMM paramters of the templates

$\sigma_K^{(0)}$ (equation(10))

$q = 1$

EXPECTATION
$Q(\boldsymbol{\Psi}; \boldsymbol{\Psi}^{(q)}) = E[\log L_c(\boldsymbol{\Psi})|\boldsymbol{x_1}, ..., \boldsymbol{x_n}, \boldsymbol{\Psi}^{(q)}]$

Posterior probabilities estimates

$\tau_{ik}^{(q)} = P(Z_{ik} = 1|\boldsymbol{x_i}; \boldsymbol{\Psi}^{(q)}); k = 1 \ldots, K; i = 1 \ldots n$

MAXIMIZATION
parameters update: $\boldsymbol{\Psi}^{(q+1)}$
$\boldsymbol{\theta_k^{(q+1)}}$ (equation (9))
$\sigma_k^{(q+1)}$ (equation (10)
$\gamma_k^{(q+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(q)}}{n}, \quad 1 \le k \le K$

q¡I; q increases by 1

q=I

K¿Ng; K increases by 1

K Class output

$I_i = \text{argmax}_{1 \le k \le K} \tau_{ik}^{(I)}, \forall i \in \{1, ..., n\}.$

K=Ng

Optimal number of clusters
Plot with Log-likelihood values for $K = 2, ..., Ng$
Slope heuristic

**Fig. 3.** MixFMM algorithm. $m$ is the number of FMM components, $Ng$ the maximum number of classes and $I$ the EM iterations.

Vip positive) of the mouse visual cortex from the [43], a repository of intracellular electrophysiological neuronal recordings. In particular, we have selected those generated by the short square stimulus with the lowest stimulus amplitude that elicited a single spike for each cell in the database. To facilitate the comparison, the signals have been down-sampled to have 64 samples, as in other datasets, aligning with the maximum at data point 20. Even though the mean spike waveforms are highly similar, as can be seen in Fig. 5, the distinction of GABAergic neuron spikes is crucial as they are radically different in terms of the physiological function [44].

Table 1 gives information on these datasets. Furthermore, 2D descriptions are shown in Fig. 4 using the first two principal components; while the mean spike waveform by class (solid line), along with the global mean (dashed line), are displayed in Fig. 5. The plots in Fig. 4 show the ability of the first two principal components for class discrimination. The percentage of variability explained by the two principal components is low in general, except for the VGA case. The first component is fundamentally associated with

**Table 1**

Datasets description: Label, number of spikes and template similarity level.

| Dataset Label | N° classes | N° Spikes | Templates Similarity |
|---|---|---|---|
| **E11** | 3 | $n_1 = 1123, n_2 = 1061, n_3 = 1212$ | Low |
| **E21** | 3 | $n_1 = 1125, n_2 = 1092, n_3 = 1172$ | Low |
| **D11** | 3 | $n_1 = 1136, n_2 = 1126, n_3 = 1103$ | Middle |
| **D21** | 3 | $n_1 = 1164, n_2 = 1119, n_3 = 1101$ | High |
| **EPT** | 4 | $n_1 = 1050, n_2 = 6638, n_3 = 1449, n_4 = 58$ | Middle |
| **VGA** | 3 | $n_1 = 163, n_2 = 221, n_3 = 123$ | High |



**Fig. 4.** First two principal components projections of the spikes in the datasets colored by class. Components' explained variance is shown in the corresponding axis. Colors are assigned with the data set class labels as follows: Class1 (red), Class2 (blue), Class3 (green), and Class4 (purple). Panel labels correspond to dataset labels.

the differences in amplitude around the prominent peak, which is more evident in the EPT case. Fig. 5 shows differences between the original class templates, which are highest for the E11 or E21 cases, moderate for D11 and EP, and lowest for the D21 or VGA sets. The results concerning what was observed in this figure are discussed below.

## 4. Numerical studies

In this Section, the MixFMM approach, PCA+KM, PCA+GMM and SPDF+KM have been compared in the six datasets described above; specifically, the $MixFMM_1$ and $MixFMM_3$ homoscedastic models have been considered. A preliminary analysis with het-
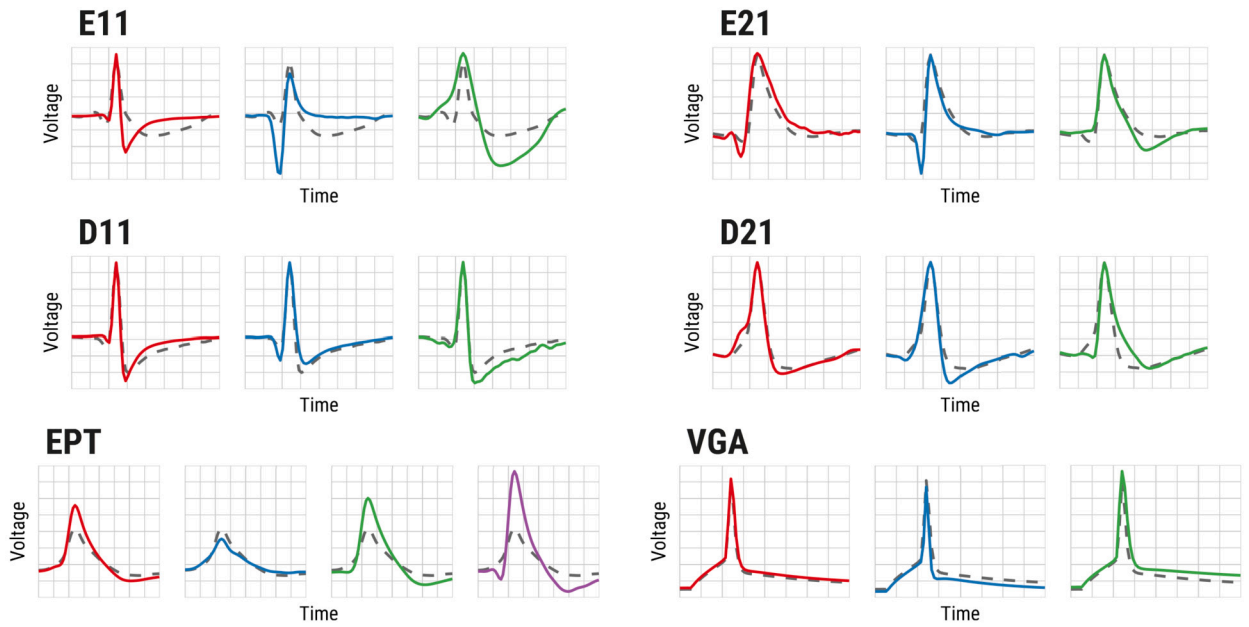
**Fig. 5.** Templates by class (solid lines) and global mean (dashed) across datasets. Panel labels corresponds to dataset labels. Colors are assigned with the data set class labels as follows: Class1 (red), Class2 (blue), Class3 (green), and Class4 (purple).

**Table 2**
Estimated number of clusters, Accuracy and $DB$ values, across datasets and methods.

|  | Number of clusters | | | | | |
|---|---|---|---|---|---|---|
|  | E11 | E21 | D11 | D21 | EPT | VGA |
| **PCA+KM** | **3** | **3** | **3** | **3** | **3** | **3** |
| **PCA+GMM** | 4 | 4 | 4 | 3 | 4 | 3 |
| **SPDF+KM** | **3** | **3** | **3** | **3** | 2 | **3** |
| **MixFMM$_1$** | **3** | **3** | **3** | **3** | 2 | **3** |
| **MixFMM$_3$** | **3** | **3** | **3** | **3** | 2 | **3** |
|  | Accuracy | | | | | |
|  | E11 | E21 | D11 | D21 | EPT | VGA |
| **PCA+KM** | 0.992 | 0.875 | 0.665 | 0.720 | 0.589 | 0.617 |
| **PCA+GMM** | 0.960 | 0.851 | 0.905 | 0.640 | 0.680 | 0.710 |
| **SPDF+KM** | 0.992 | **0.981** | **0.981** | **0.967** | 0.861 | 0.720 |
| **MixFMM$_1$** | 0.992 | 0.936 | 0.895 | 0.705 | **0.874** | **0.755** |
| **MixFMM$_3$** | **0.994** | 0.966 | 0.941 | 0.743 | **0.874** | 0.746 |
|  | $DB$ | | | | | |
|  | E11 | E21 | D11 | D21 | EPT | VGA |
| **PCA+KM** | 0.727 | 1.455 | 2.531 | 2.144 | 1.601 | 1.213 |
| **PCA+GMM** | 2.411 | 2.421 | 3.149 | 3.291 | 2.355 | 1.580 |
| **SPDF+KM** | 0.788 | **1.072** | **1.341** | **1.700** | 1.311 | 1.503 |
| **MixFMM$_1$** | 0.728 | 1.357 | 1.785 | 2.091 | 0.744 | **1.144** |
| **MixFMM$_3$** | **0.726** | 1.341 | 1.717 | 2.070 | **0.742** | **1.144** |

eroscedastic models resulted in quite similar variance estimators across the datasets, as well as less accurate and more unstable solutions than those of the homoscedastic models.

Table 2 summarizes the main results across different approaches. The procedure to estimate the number of clusters using the MixFMM gives better results than PCA+GMM and is comparable to other approaches. In terms of both, accuracy and $DB$ values, the MixFMM approach attains the best result in three out of the six datasets, achieving significant improvements over the PCA-based
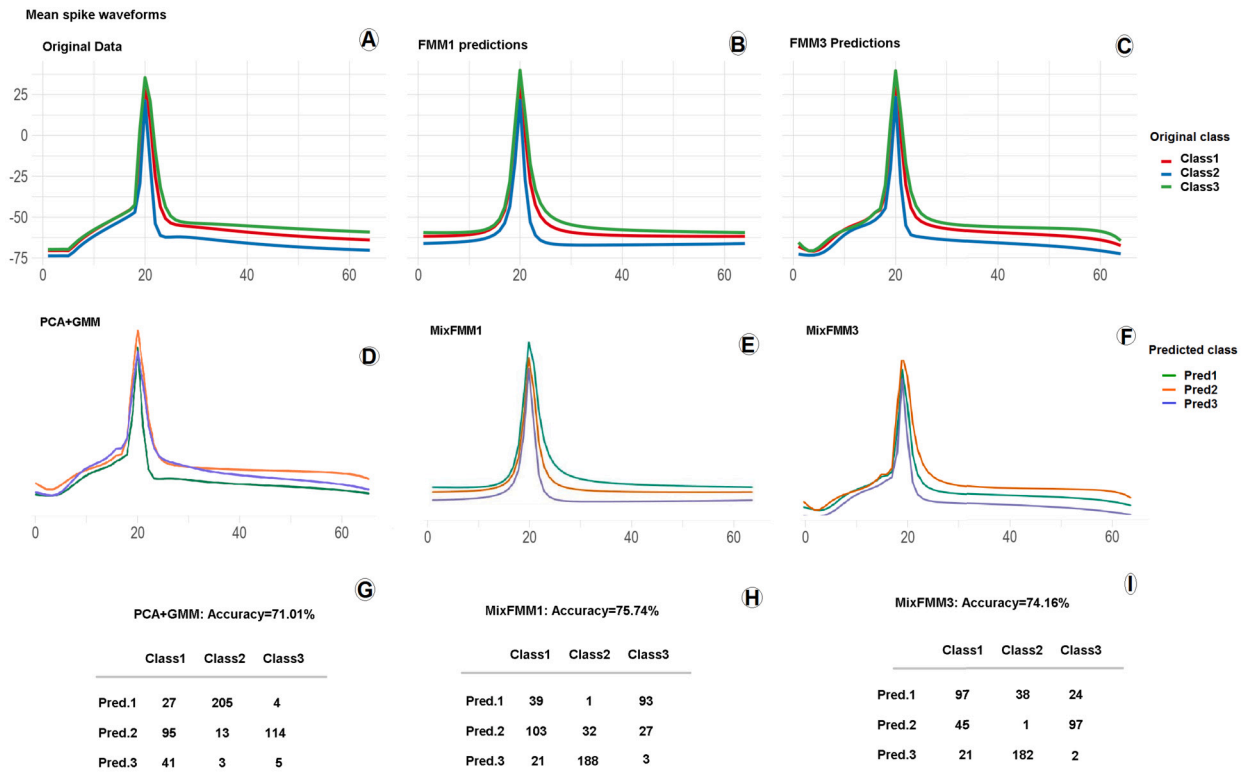
**Fig. 6.** VGA dataset analysis. Templates for the original classes (panel A), and their FMM1 and FMM3 predictions (panels B and C). Templates for the predicted classes with PCA+GMM (panel D), MixFMM1 (panel E) column) and MixFMM3 (panel F). Original class colors are assigned with the data set class labels as follows: Class1 (red), Class2 (blue), Class3 (green), and Class4 (purple). Predicted class colors are assigned as follows: Pred1 (jade green), Pred2 (ocher), Pred3 (slate blue). Confusion matrices derived from using PCA+GMM (panel G), MixFMM1 (panel H), and MixFMM3 (panel I).

methods in some cases. In particular, MixFMM is the best clustering approach in the real datasets, while the SPDF+KM approach gives better results in most of the simulated datasets. Regarding the two FMM approaches, MixFMM$_3$ performs better than MixFMM$_1$ in most cases. However, MixFMM$_1$ could be a suitable alternative in real cases with high noise level data. The index values for the clusters obtained using the MixFMM approaches agree with what is shown Fig. 5. However, other methods do not give the same expected results. For example, the PCA+KM does not work as well in discriminating the spikes of the sets D11 and EPT as for those of D21. The results for the VGA dataset are presented in more detail. Fig. 6 illustrates the difference between templates, across the original and predicted classes for the PCA+GMM and MixFMM approaches. The templates for predicted KM classes are similar to those derived with the GMM classification, so they are not shown. The first row in Fig. 6 shows the templates for the original classes (first column) and their FMM1 (second column) and FMM3 (third column) predictions. The three templates of the MixFMM predicted classes resemble those of the original classes. However, the PCA+GMM templates for classes 1 and 3 are not so different from each other. In fact, different initializations generate predicted classes in which the original classes 1 and 3 are mixed differently. Furthermore, the FMM parameters $A$ and $\omega$ discriminate between the three templates. Specifically, using sub-indexes with the original class numbers, we have that, $A_1 = 47$, $\omega_1 = 0.07$; $A_2 = 45$, $\omega_2 = 0.05$; and $A_3 = 50$, $\omega_3 = 0.08$. The class 2 template is described with the smallest $\omega$ and $A$ values. In fact, most class 2 spikes are grouped together in the first iteration of the MixFMM algorithm. Class 1 and 3 templates differ on the $A$ values and, to a lesser extent, on the $\omega$ values. We need more iterations to differentiate the other two classes. The predicted classes after 5 iterations have templates that mimic those of the original class. On the other hand, predicted class templates obtained with other approaches do not differ by a clear characteristic and therefore fail to represent the original classes 1 and 3.

Note that our results for PCA+GM are better than those given in [13] in the common datasets; this is probably due to the fact that we have used R's mclust instead of Matlab's gm.

Finally, in order to validate the robustness of the different approaches, the distribution of the accuracy has been studied across a hundred repetitions for the methods, with ten initializations in each dataset. Fig. 7 shows the associated boxplots, which illustrate that the MixFMM models achieve an accurate solution independently of the initialization in all the scenarios analyzed.

## 5. Discussion

We present a novel model-based functional data clustering method and prove its potential to solve the SS problem. In particular, we show that the MixFMM approach achieves better results than PCA-based methodologies often used in SS applications, with
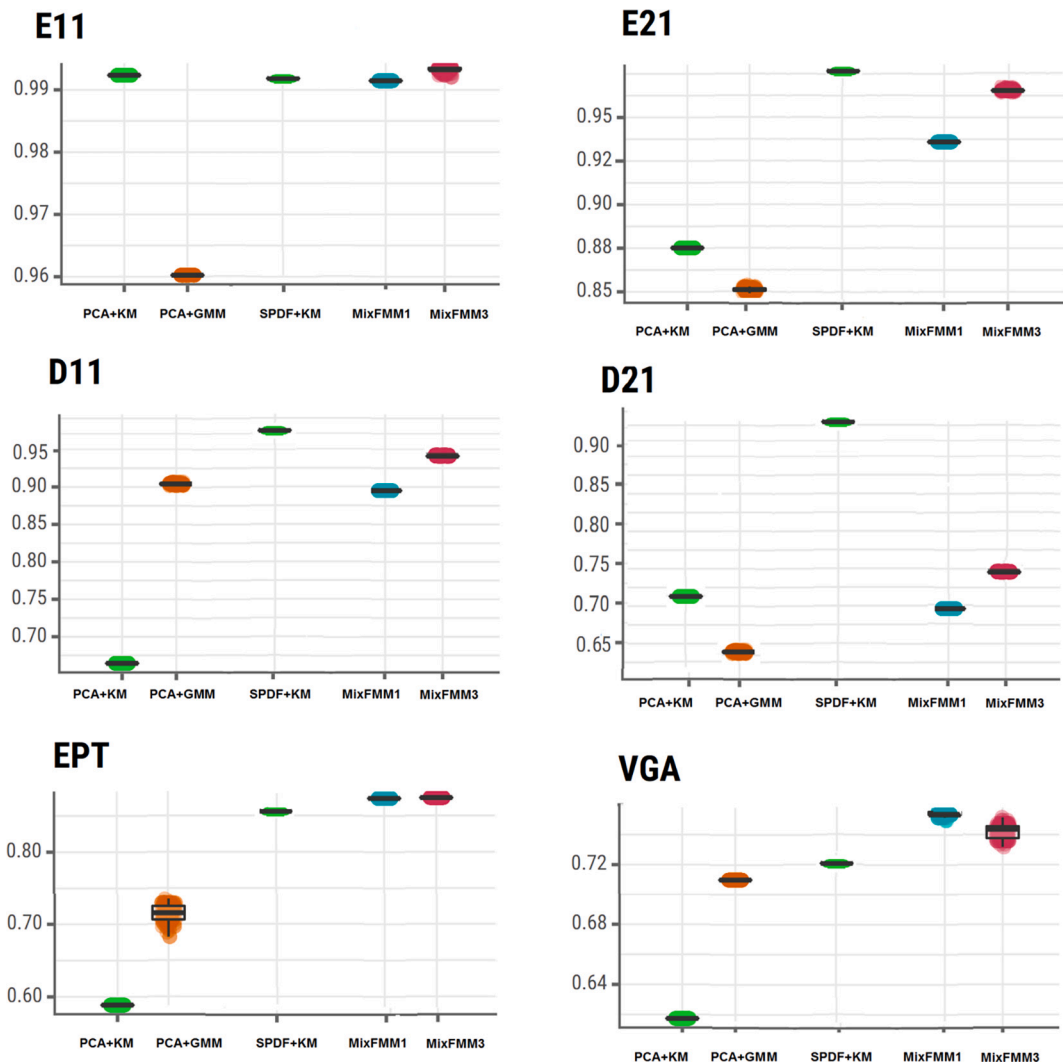
**Fig. 7.** Boxplots of the accuracy across a hundred repetitions for each dataset and the five clustering methods: KM (green), GMM (orange), MiXFMM1 (blue) and MixFMM3 (red). Panel labels correspond to dataset labels.

significant improvements in some scenarios, including very noisy cases. Compared with SPDF+KM, our method gives better results in the two real data sets and 1 out of 4 simulated datasets.

The approach uses combinations of non-linear parametric functions, the so-called FMM waves, to describe the signals. It has important advantages over the alternatives for two powerful reasons: firstly, the accuracy with which these waves describe any type of oscillatory morphology, much better than linear functions; and secondly, that the meaning that the parameters accounting for phase, amplitude and shape variations have. Furthermore, the interpretation of the FMM parameters allows open problems in spike sorting to be addressed. On the one hand, overlapping spikes can be automatically detected by two significant FMM components with close $\alpha$ values and high $A$ values. In addition, positive or negative spikes are recognize with values of $\beta$ far from $0$ and $\pi$, respectively. On the other hand, the questions of high channel-count analysis and spike shape variation caused by electrode drift can be addressed with the aid of recent results in [45,46]. These papers present the adequate multicomponent and multidimensional FMM model when sequential events are observed from different locations. It can thus be used to model multi-channel spike data. In addition, the said papers extend the role of the parameters beyond their ability to describe morphology. Specifically, it is shown that the $\omega$ and $\alpha$ parameters, common to the signals recorded from different locations, describe the three-dimensional process. Meanwhile, the $A$ and $\beta$ parameters describe the relative location from which the data has been recorded. Hence, only these last parameters must be analyzed to identify possible electrode drift by comparing the values of these parameters for typical spikes with those of the same individual recorded in other experiments. In [45], multi-lead ECG signals are analyzed. That paper shows how the multidimensional and multicomponent FMM model can identify spikes common to different channels, even if the amplitude is too small in some of them or the waveform is corrupted by noise. The development of multivariate models and strategies to address the questions mentioned so far, as well as to cluster using multichannel data and real dataset problems, is beyond the scope of this paper, but is part of our

future work. We will address the comparison of the FMM classification with others derived using alternative distances, such as [47], and, reproducible methods such as [48]. An extensive review of cluster multi-channel spike recordings can be found in [49].

Regarding the analysis of real cases, the results of the MixFMM clustering entails interesting contributions in neuroscience. In particular, neuronal signals from the temporal lobe, such as EPT, are related to specific visual stimuli [2]. The function in the nervous system of the different kinds of GABAergic neurons is radically different, so the distinction of their signals is crucial. In VGA, signals from three types of neurons are differentiated: Pvalb-positive neurons, key in learning and memory process, as detailed in [50]; Vip-positive neurons, related to circadian rhythms [51]; and Htr3a-positive neurons, related to bladder dysfunction [52]. The reader interested in how the FMM can shed more light on the modeling of neuronal action potential is referred to [27] and [53]. Due to the flexibility of the FMM model for discriminating different waveforms, the MixFMM algorithm would also be suitable for classifying spikes arising in other cells and experimental conditions, as well as in other contexts, such as those from calcium signals; but always after an ad hoc preprocessing to filter and cut the spike segments. Specifically, the FMM approach may be useful to classify, in a supervised or unsupervised manner, epileptic and control APs, as the former show an increase in waveform duration, a decrease in the waveform amplitude, and an increase in the signal noise. The two first characteristics are just what the $\omega$ and $A$ parameters measure. Meanwhile, a noise level measure is given with the proportion of the variation that can be predicted from the FMM model.

The motivation of the present paper has been the classification of neuronal signals, but the methodology has a wide range of applications. In fact, the MixFMM is a general approach in functional clustering when oscillatory or quasi-oscillatory signals are the target. In a field in which the literature is ever-growing, and many excessively specific and complex proposals are arising, the need for flexible and robust proposals is vital. The MixFMM brings together all of these characteristics with a simple but rich parametric formulation that allows differences in phase and/or amplitude to be taken into account. Oscillatory signals appear in astronomy, economy, spectrometry, environmental science, medical imaging, or electrocardiography, among other fields. In addition, from a theoretical perspective, many extensions of the basic mixture framework could be explored, such as a regularized estimation algorithm or the inclusion of mixed-effects [30].

Finally, two limitations of this study are, on the one hand, that the computational cost is relatively high. However, it could be significantly reduced by implementing the parameter search in a more computationally efficient programming language, such as C, and using GPU-based computing. On the other hand, validation indexes have not shown a clear winner between MixFMM and SPDF+KM. The latter is superior to the former in some simulated cases and the former is better in the real datasets. This may be due to the waveforms and the noise structure chosen to generate the simulated datasets, a question that deserves further research.

## 6. List of acronyms

FMM: Frequency Modulated Möbius.
MixFMM: Mixture Frequency Modulated Möbius.
SS: Spike Sorting.
PCA:Principal Component Analysis.
KM: K-Means.
GMM: Gaussian Mixture Models.
ECG: Electrocardiogram.
MLE: Maximum Likelihood Estimator.
EM: Expectation-maximization.
$DB$: Davies-Bouldin index
SPDF: Shape, Phase, and Distribution Feature.
EPT: Epileptic.
VGA: Dataset label.

## CRediT authorship contribution statement

C.R. and A.R.-C: Conceived and designed the experiments, performed the experiments, analyzed and interpreted the data and wrote the paper.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Data availability

The data used for this research is publicly available, and you can find repository instructions in [41] and [43].

## Acknowledgments and funding sources

## References

[1] G. Buzsáki, Large-scale recording of neuronal ensembles, Nat. Neurosci. 7 (2004) 446–451.

[2] H.G. Rey, C. Pedreira, R.Q. Quiroga, Past, present and future of spike sorting techniques, Brain Res. Bull. 119 (2015) 106–117.

[3] N.K. Sharma, C. Pedreira, M. Centeno, U.J. Chaudhary, et al., A novel scheme for the validation of an automated classification method for epileptic spikes by comparison with multiple observers, Clin. Neurophysiol. 128 (2017) 1246–1254, https://www.sciencedirect.com/science/article/pii/S1388245717301621.

[4] M. Moghaddasi, M. Aliyari Shoorehdeli, Z. Fatahi, A. Haghparast, Unsupervised automatic online spike sorting using reward-based online clustering, Biomed. Signal Process. Control 56 (2020) 101701, https://www.sciencedirect.com/science/article/pii/S1746809419302824.

[5] R.Q. Quiroga, Z. Nadasdy, Y. Ben-Shaul, Unsupervised spike detection and sorting with wavelets and superparamagnetic clustering, Neural Comput. 16 (2004) 1661–1687.

[6] C. Ekanadham, D. Tranchina, E.P. Simoncelli, A unified framework and method for automatic neural spike identification, J. Neurosci. Methods 222 (2014) 47–55, https://www.sciencedirect.com/science/article/pii/S0165027013003415.

[7] C.R. Caro-Martín, J.M. Delgado-García, A. Gruart, R. Sánchez-Campusano, Spike sorting based on shape, phase, and distribution features, and k-tops clustering with validity and error indices, Sci. Rep. 8 (2018) 1–28.

[8] B. Souza, V. Lopes dos Santos, J. Bacelo, A. Tort, Spike sorting with Gaussian mixture models, Sci. Rep. 9 (2019) 3627.

[9] M. Rácz, C. Liber, E. Németh, R. Fiáth, et al., Spike detection and sorting with deep learning, J. Neural Eng. 17 (2019) 016038.

[10] D. Carlson, L. Carin, Continuing progress of spike sorting in the era of big data, Curr. Opin. Neurobiol. 55 (2019) 90–96, https://www.sciencedirect.com/science/article/pii/S0959438818301375.

[11] J. Sukiban, N. Voges, T.A. Dembek, R. Pauli, et al., Evaluation of spike sorting algorithms: application to human subthalamic nucleus recordings and simulations, Neuroscience 414 (2019) 168–185, https://www.sciencedirect.com/science/article/pii/S0306452219304750.

[12] P.K. Wang, S.H. Pun, C.H. Chen, E.A. McCullagh, et al., Low-latency single channel real-time neural spike sorting system based on template matching, PLoS ONE 14 (2019) 1–30, https://doi.org/10.1371/journal.pone.0225138.

[13] R. Veerabhadrappa, M. Ul Hassan, J. Zhang, A. Bhatti, Compatibility evaluation of clustering algorithms for contemporary extracellular neural spike sorting, Front. Syst. Neurosci. 14 (2020) 34, https://www.frontiersin.org/article/10.3389/fnsys.2020.00034.

[14] Y. Zhang, et al., A robust spike sorting method based on the joint optimization of linear discrimination analysis and density peaks, Sci. Rep. 12 (2022) 1–16.

[15] M. Radmanesh, A.A. Rezaei, M. Jalili, A. Hashemi, M.M. Goudarzi, Online spike sorting via deep contractive autoencoder, Neural Netw. 155 (2022) 39–49, https://www.sciencedirect.com/science/article/pii/S089360802200301X.

[16] J. Jacques, C. Preda, Functional data clustering: a survey, Adv. Data Anal. Classif. 8 (2014) 24, https://hal.inria.fr/hal-00771030.

[17] W. Lee, M.F. Miranda, P. Rausch, V. Baladandayuthapani, et al., Bayesian semiparametric functional mixed models for serially correlated functional data, with application to glaucoma data, J. Am. Stat. Assoc. 114 (2019) 495–513.

[18] Y. Lim, H.-S. Oh, Y.K. Cheung, Multiscale clustering for functional data, J. Classif. 36 (2019) 368–391.

[19] Q. Zhong, H. Lin, Y. Li, Cluster non-Gaussian functional data, Biometrics 77 (2020) 852–865, https://onlinelibrary.wiley.com/doi/abs/10.1111/biom.13349, https://onlinelibrary.wiley.com/doi/pdf/10.1111/biom.13349.

[20] T.L.J. Ng, T.B. Murphy, Model-based clustering for random hypergraphs, Adv. Data Anal. Classif. 16 (2022) 691–723.

[21] K. You, C. Suh, Parameter estimation and model-based clustering with spherical normal distribution on the unit hypersphere, Comput. Stat. Data Anal. 171 (2022) 107457, https://www.sciencedirect.com/science/article/pii/S0167947322000378.

[22] J.S. Marron, J. Ramsay, L.M. Sangalli, A. Srivastava, Functional data analysis of amplitude and phase variation, Stat. Sci. 30 (2015) 468–484, https://doi.org/10.1214/15-STS524.

[23] J. Park, J. Ahn, Clustering multivariate functional data with phase variation, Biometrics 73 (2017) 324–333, https://onlinelibrary.wiley.com/doi/abs/10.1111/biom.12546, https://onlinelibrary.wiley.com/doi/pdf/10.1111/biom.12546.

[24] G. Claeskens, E. Devijver, I. Gijbels, Nonlinear mixed effects modeling and warping for functional data using B-splines, Electron. J. Stat. 15 (2021) 5245–5282, https://doi.org/10.1214/21-EJS1917.

[25] C. Rueda, Y. Larriba, S.D. Peddada, Frequency modulated Möbius model accurately predicts rhythmic signals in biological and physical sciences, Sci. Rep. 9 (2019) 1–10.

[26] C. Rueda, Y. Larriba, A. Lamela, The hidden waves in the ecg uncovered revealing a sound automated interpretation method, Sci. Rep. 11 (2021) 1–11.

[27] A. Rodríguez-Collado, C. Rueda, Electrophysiological and transcriptomic features reveal a circular taxonomy of cortical neurons, Front. Human Neurosci. 15 (2021) 684950, https://www.frontiersin.org/article/10.3389/fnhum.2021.684950.

[28] C. Rueda, A. Rodríguez-Collado, Y. Larriba, A novel wave decomposition for oscillatory signals, IEEE Trans. Signal Process. 69 (2021) 960–972.

[29] A.P. Dempster, N. Laird, D. Rubin, Maximum likelihood from incomplete data via the em algorithm, J. R. Stat. Soc., Ser. B, Methodol. 39 (1977) 1–22.

[30] F. Chamroukhi, H.D. Nguyen, Model-based clustering and classification of functional data, WIREs Data Min. Knowl. Discov. 9 (2019) e1298, https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/widm.1298, https://wires.onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1298.

[31] L. Birgé, P. Massart, Minimal penalties for Gaussian model selection, Probab. Theory Relat. Fields 138 (2007) 33–73.

[32] C. Bouveyron, E. Côme, J. Jacques, The discriminative functional mixture model for a comparative analysis of bike sharing systems, Ann. Appl. Stat. 9 (2015) 1726–1760, http://www.jstor.org/stable/43826443.

[33] M.L.L. García, R. García-Ródenas, A.G. Gómez, K-means algorithms for functional data, Neurocomputing 151 (2015) 231–245.

[34] Allen Brain Institute, Allen SDK, https://allensdk.readthedocs.io/en/latest/index.html, 2021, Python package version v2.13.0.

[35] I. Fernández, A. Rodríguez-Collado, Y. Larriba, A. Lamela, C. Canedo, C. Rueda, FMM: an R package for modeling rhythmic patterns in oscillatory systems, R Journal (2022) 361–379.

[36] R Core Team, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, 2020, https://www.R-project.org/.

[37] L. Scrucca, M. Fop, T.B. Murphy, A.E. Raftery, mclust 5: clustering, classification and density estimation using Gaussian finite mixture models, R J. 8 (2016) 289–317, https://doi.org/10.32614/RJ-2016-021.

[38] B. Desgraupes, clusterCrit: Clustering Indices, University of Paris, Ouest, 2018, https://CRAN.R-project.org/package=clusterCrit, R package version 1.2.8.

[39] H. Wickham, ggplot2: Elegant Graphics for Data Analysis, second edn., Springer - Verlag, New York, 2016, https://ggplot2.tidyverse.org.

[40] P.T. Inc, Collaborative Data Science, Plotly Technologies Inc., Montreal, QC, 2015, https://plot.ly.

[41] R.Q.Quiroga, Simulated Datasets, 2009, http://www2.le.ac.uk/centres/csn/research-2/spike-sorting. Available in https://leicester.figshare.com/articles/dataset/Simulated_dataset/11897595.

[42] R.Q. Quiroga, Dataset: Human Single-Cell Recording, 2019, https://leicester.figshare.com/articles/Dataset-Human-single-cell-recording/11302427/1. Available in https://leicester.figshare.com/articles/Dataset-Human-single-cell-recording/11302427/1.

[43] Allen Brain Institute, Allen Cell Types Database, 2021, https://celltypes.brain-map.org/data. Available in https://celltypes.brain-map.org/data.

[44] N.W. Gouwens, S.A. Sorensen, F. Baftizadeh, A. Budzillo, et al., Integrated morphoelectric and transcriptomic classification of cortical gabaergic cells, Cell 183 (2020) 935–953.

[45] C. Rueda, A. Rodríguez-Collado, I. Fernández, C. Canedo, et al., A unique cardiac electrocardiographic 3D model. Toward interpretable AI diagnosis, https://www.sciencedirect.com/science/article/pii/S2589004222018892, 2022.

[46] C. Rueda, I. Fernández, C. Canedo, Y. Larriba, Fda with Möbius waves: estimation of signals and their derivatives with applications, Preprint, 2023.

[47] P. Larionov, T. Juergens, T. Schanze, Correlation-based spike sorting of multivariate data, Curr. Dir. Biomed. Eng. 5 (2019) 113–116.

[48] P. Yger, G.L. Spampinato, E. Esposito, B. Lefebvre, et al., A spike sorting toolbox for up to thousands of electrodes validated with ground truth recordings in vitro and in vivo, eLife 7 (2018) e34518, https://doi.org/10.7554/eLife.34518.

[49] R.B. Bod, et al., From end to end: gaining, sorting, and employing high-density neural single unit recordings, Front. Neuroinform. 16 (2022).

[50] F. Donato, S. Rompani, P. Caroni, Parvalbumin-expressing basket-cell network plasticity induced by experience regulates adult learning, Nature 504 (2013) 272–276.

[51] C. Mazuski, S. Chen, E. Herzog, Different roles for vip neurons in the neonatal and adult suprachiasmatic nucleus, J. Biol. Rhythms 35 (2020) 465–475.

[52] E. Ritter, M. Southard-Smith, Dynamic expression of serotonin receptor 5-ht3a in developing sensory innervation of the lower urinary tract, Front. Neurosci. 10 (2017) 592, https://www.frontiersin.org/article/10.3389/fnins.2016.00592.

[53] A. Rodríguez-Collado, C. Rueda, A simple parametric representation of the Hodgkin-Huxley model, PLoS ONE 16 (2021) 1–19, https://doi.org/10.1371/journal.pone.0254152.