

A novel and fully automatic spike sorting implementation with variable number of features

Fernando Chaure* (1,2,3,4), Hernan Rey* (1), Rodrigo Quian Quiroga* (1)

(1) Centre for Systems Neuroscience, University of Leicester, Leicester, LE1 7RH, United Kingdom.

(2) Instituto de Ingeniería Biomédica, UBA, Buenos Aires, Argentina.

(3) Estudios de Neurociencias y Sistemas Complejos (ENYS), CONICET - Hospital El Cruce - UNAJ, Florencio Varela, Argentina.

(4) Instituto de Biología Celular y Neurociencias "Prof. E. De Robertis", Facultad de Medicina, UBA, Buenos Aires, Argentina.

* These authors contributed equally

+ Corresponding author: rqqg1@le.ac.uk

Abstract

The most widely used spike sorting algorithms are semiautomatic in practice, requiring manual tuning of the automatic solution to achieve good performance. In this work, we propose a new fully automatic spike sorting algorithm that can capture multiple clusters of different size and densities. In addition, we introduce an improved feature selection method by using a variable number of wavelet coefficients based on the degree of non-gaussianity of their distributions. We evaluated the performance of the proposed algorithm with real and simulated data. With real data from single channel recordings, in about 95% of the cases the new algorithm replicated, in an unsupervised way, the solutions obtained by expert sorters, who manually optimized the solution of a previous semiautomatic algorithm. This was done while maintaining a low number of false positives. With simulated data from single channel and tetrode recordings, the new algorithm was able to correctly detect many more neurons compared to previous implementations and also compared to recently introduced algorithms, while significantly reducing the number of false positives. In addition, the proposed algorithm showed good performance when tested with real tetrode recordings.

New & Noteworthy

We propose a new fully automatic spike sorting algorithm including several steps that allow the selection of multiple clusters of different size and densities. Moreover, it defines the dimensionality of the feature space in an unsupervised way. We evaluated the performance of the algorithm with real and simulated data, from both single channel and tetrode recordings. The proposed algorithm was able to outperform manual sorting from experts, and other modern unsupervised algorithms.

INTRODUCTION

Extracellular recordings of single neuron activity are done by placing electrodes in brain tissue. The electrical potential changes measured at the electrode tip reflect the spiking activity of neurons close enough to the electrode, plus background activity elicited by neurons further away from the tip (black trace in the top panel of Figure 1A). In principle, the spikes fired by a neuron recorded by a given electrode have a distinct shape. This is mainly determined by the morphology of the neuron's dendritic tree, the distance and orientation relative to the recording site, the distribution of ion channels, and the properties of the extracellular medium (Gold et al. 2006). Spike sorting algorithms detect these spikes (Figure 1A, top) and, using features extracted from the waveforms, group them into clusters corresponding to the putative activity of different neurons (Figure 1A, bottom right) (Lewicki 1998; Quian Quiroga 2007). The importance of spike sorting is stressed by the fact that close-by neurons recorded from the same electrode can respond to completely different things and, therefore, it is crucial to know which spike corresponds to which neuron. This is the case, for example, in the human and the rat hippocampus, where nearby neurons fire to unrelated concepts in the first case (Rey et al. 2015a; De Falco et al. 2016) and to distant place fields in the latter (Redish et al. 2001).

The simultaneous recording of a large number of electrodes, thus accessing the activity of populations of neurons, is becoming an essential tool for understanding complex behaviors and network properties in the brain (Quian Quiroga and Panzeri 2009). Silicon probes have been developed in the last 10 years (Csicsvari 2003; Buzsáki 2004; Blanche et al. 2005) and large multielectrode arrays with up to thousands of electrodes are already being used for recording in retinal patches (Litke et al. 2004), cell cultures (Lambacher et al. 2011), or brain slices (Frey et al. 2009). Furthermore, a large number of channels is currently used to record from local circuits in behaving animals (Berényi et al. 2014). With such high channel count, the manual supervision of each single channel might turn into a very time-consuming task. In addition, it is well known that the subjectivity introduced by the human intervention creates an additional source of sorting errors (Harris et al. 2000; Wood et al. 2004). In parallel to these advances, there are also several cases where single electrode recordings are still routinely used, with automatic implementations being critical to avoid potential subjective biases of manual solutions. This is, for example, the case in invasive human recordings performed with depth electrodes implanted in patients with epilepsy (Rey et al. 2015a) or chronic implants for brain-machine interface (Homer et al. 2013). In this context, the unsupervised classification of single units, both from single channel recordings and large electrode arrays, has become the bottleneck to fully reach the potential of extracellular recordings (Rey et al. 2015b; Harris et al. 2016).

Many spike sorting algorithms have been developed in the past years (see Rey et al. 2015b for a review). Some of these methods are based on Bayesian statistical frameworks, relying in some cases on a Gaussian model of the distribution of the spike waveforms (Harris et al. 2000; Rossant et al. 2016). Past studies have shown the non-Gaussian variability of the spike shapes and the non-stationarity of the extracellular recordings, due to, for example, small electrode drifts (Fee et al. 1996; Harris et al. 2016) or the presence of **bursting cells** (Henze et al. 2000), and motivated the development of **nonparametric approaches**. One such method is *Waveclus* (Quian Quiroga et al. 2004), a spike sorting algorithm that uses **wavelet decomposition to extract features of the spike waveforms and superparamagnetic**

clustering (SPC) to cluster the spikes in this feature space. However, *Wave_clus*, as well as the other most widely used algorithms, are semiautomatic in practice, requiring manual tuning of a first automatic solution to achieve good performance. To tackle this problem, we analyzed the actions typically performed by *Wave_clus* users to optimize spike sorting performance. Then, we implemented a set of heuristic modifications to the algorithm in order to reproduce the users' actions in a fully automatic way. Moreover, we present an automatic criterion for selecting the number of features —i.e. of wavelet coefficients— to be used for clustering, given that, intuitively, only a few features are required to separate relatively few neurons, and more features are necessary for recordings containing a larger (and in principle unknown) number of neurons.

We evaluated the performance of the proposed method using both real and simulated single channel and tetrode recordings. With the simulated data from single channels, the proposed algorithm significantly outperformed other automatic sorting algorithms, as well as the experts' manual solution with the former *Wave_clus* implementation. With the simulated tetrode recordings, the new algorithm also outperformed other recently introduced methods, as quantified by the number of hits (i.e. correctly identified clusters) and false positives. With the real single channel recordings, the algorithm retrieved, in a fully unsupervised way, about 95% of the clusters isolated by the sorting experts while keeping a low number of false positives. When assessed with real tetrode recordings, the new automatic algorithm outperformed the experts' manual solution with the former *Wave_clus* implementation.

MATERIALS AND METHODS

Former Wave_clus implementation

Our new method is based on the former Matlab implementation of *Wave_clus*, introduced in Quiroga et al. (2004). As with other spike sorting methods, this algorithm has four main steps: filtering, detection, feature extraction and clustering (Rey et al. 2015b). Zero-phase filtering is done by using a second-order bandpass elliptic filter in the range 300-3000 Hz. Spike detection is performed by setting a threshold as:

$$Thr = 5\sigma_n, \text{ with } \sigma_n = \text{median} \left\{ \frac{|x|}{0.6745} \right\},$$

where x is the bandpass filtered signal (Rey et al. 2015b). An example of a signal filtered for detection and thresholding is shown in the top panel of Figure 1A. For each detected spike, 64 samples are saved for further analysis, aligned to their maximum at data point 20. To avoid spike misalignments due to low sampling, spike maxima location is refined by using cubic splines interpolated waveforms with 320 samples. After realignment, the waveforms are downsampled back to 64 points.

Feature extraction is done using a 4-scale multiresolution decomposition with a Haar wavelet, resulting in 64 wavelet coefficients associated with each detected spike. To assess the ability of each coefficient to separate different clusters, the algorithm uses a Lilliefors test (a normality test based on the Kolmogorov–Smirnov test), retaining the 10 most significant ones. To minimize the effect of outliers in the test, only values within ± 3 standard deviations are considered for each coefficient. For tetrode recordings, the spikes detected at each

channel were concatenated. Then, if Nch channels are used with 64 samples per spike, there will be a total of $Nch \times 64$ wavelet coefficients, from which *Wave_clus* would pick the $Nch \times 10$ most significant ones.

Finally, a nonparametric clustering is performed in the feature space using **superparamagnetic clustering (SPC)**. SPC is an unsupervised approach in which the grouping of points into clusters depends on nearest-neighbor interactions (Blatt et al. 1996, 1997). SPC generates a family of solutions as a function of a parameter called the temperature, which is the key parameter to determine **how clusters are split** (Domany et al. 1999). In analogy with models in Statistical Mechanics, at **low temperatures**, all data points are **highly correlated** and are therefore **grouped into a single or relatively few clusters**. On the other hand, at **high temperatures** the **correlations are too weak** and clusters break up into **many groups with very few members in each group**. At **a certain temperature range between these two extremes**, **natural clusters appear** (i.e. the superparamagnetic regime) and only points corresponding to data from relatively high-density regions are grouped. With SPC, for each temperature a different data partition is generated, resulting in a “temperature plot” (Figure 1A, bottom left). *Wave_clus* uses **a range of temperatures from 0 to 0.25 in increments of 0.01**, where each temperature is labelled as T_i , with $i = \{0, \dots, 25\}$. At each temperature, clusters are sorted in decreasing order with respect to their size; thus, at temperature T_i the largest cluster is denoted as $C_1^{T_i}$, the second largest as $C_2^{T_i}$, etc.

In *Wave_clus*, clusters are automatically identified according to a thresholding procedure based on their size. The idea is that **as the temperature is increased, new clusters will appear**. *Wave_clus* selects the highest temperature where at least one of the sorted clusters $C_i^{T_n}$ increases its size **a minimum number of spikes (parameter N_inc)** (Quiñero et al. 2004). At this temperature, the cluster that increased its size by at least N_inc spikes, as well as all the clusters with a larger size are selected. The parameter N_inc is introduced to **avoid overclustering**—i.e. choosing a very high temperature where data is grouped into many clusters with a few members each. Once the clusters have been identified, a **template matching** procedure is used to assign the remaining unclassified waveforms. For each cluster, the centroid (mean waveform) and a measure of its total variance, σ_T , are computed, where $\sigma_T = \sqrt{\sum_{i=1}^{64} var(x_i)}$, with $var(x_i)$ denoting the variance at the i -th sample across the waveforms of a given cluster. Then, each spike is assigned to the cluster with the smallest **Euclidian distance** to its centroid, as long as this distance is smaller than $3\sigma_T$ (waveforms with a larger distance to the centroid are considered to be noise).

New Wave_clus implementation

To improve the former implementation of *Wave_clus*, we propose the use of heuristic modifications that are inspired by the actions that are usually taken by experts to optimize the automatic solution given by the previous implementation. The main limitation of the previous automatic implementation is that **a single temperature** is chosen for clustering, and in many cases, **there are clusters appearing at different temperatures**. The rationale of the new implementation is to choose all putative clusters from the different temperatures in the temperature plot (peak selection step), and then get rid of false positives by avoiding double detections (inclusion criterion) —as the same clusters may appear at different

temperatures— and also avoiding overclustering at large temperatures (regime border detection). A final improvement is the introduction of an automatic criterion to select a variable number of features (i.e. wavelet coefficients) for clustering.

Finally, we also included some improvements in the implementation of the algorithm to reduce the computing time. Specifically, we used a MEX implementation of the Matlab function *filtfilt* for the detection process, and eliminated time-consuming “for loops” in the interpolation procedure for spike detection and the wavelet decomposition of the detected waveforms (using optimized matrix manipulations). The codes and documentation of the new sorting implementation are available on GitHub (https://github.com/csn-le/wave_clus).

1) Peak Selection in the temperature plot

We first select every cluster $C_i^{T_n}$ where its size $|C_i^{T_n}|$ is increased by at least N_{inc} spikes as the temperature is increased from T_{n-1} to T_n , i.e.,

$$|C_i^{T_n}| - |C_i^{T_{n-1}}| \geq N_{inc}, \quad \text{with } |C_i^{T_{-1}}| = 0.$$

To consider other relevant members of a partition at temperature T_n in which cluster $C_i^{T_n}$ has been selected by the peak criterion, all clusters with a larger size (i.e. all clusters $C_j^{T_n}$ with $j < i$) are also selected.

2) Regime border detection

The end of the superparamagnetic regime is normally associated with an abrupt decrease of the principal cluster as we transition from $C_1^{T_{B-1}}$ to $C_1^{T_B}$, with the appearance of several small clusters (overclustering). However, in some cases, before reaching the end of the regime there is an abrupt decrease in the principal cluster in conjunction with the appearance of a large new and relevant cluster, typically being the second largest (C_2). For this reason, we first define the largest increment at a certain temperature as $LI_{T_i} = \max(|C_j^{T_i}| - |C_j^{T_{i-1}}|)$, for $j > 1$. Then, we find the temperature T_B as the minimum T_i fulfilling the condition:

$$T_B = \min(T_i) \text{ for which } \frac{|C_1^{T_i}| + LI_{T_i}}{|C_1^{T_{i-1}}|} < Thr_border,$$

with Thr_border being a threshold parameter (we used a value of 0.4 in our implementation, although we have seen that the results are not changed for $0.25 < Thr_border < 0.45$). This way, when the ratio is small, the principal cluster is largely split beyond what could have been related with a new cluster (associated with the increment LI).

3) Inclusion criterion

First, we define the overlapping coefficient

$$O_{i,j}^{T_n,T_m} = \frac{|C_i^{T_n} \cap C_j^{T_m}|}{\min(|C_i^{T_n}|, |C_j^{T_m}|)}.$$

This coefficient ranges between 0, if the clusters are completely different (i.e. they do not share any spike), and 1, when they are the same, or one of them is a subset of the other. To avoid considering the same cluster of spikes twice (from different temperature partitions), if the value of $O_{i,j}^{T_n, T_m}$ is not lower than a constant k_0 , only the cluster at the higher temperature is kept. In our implementation, we chose $k_0 = 0.9$. We also observed that the performance of the algorithm is not affected by this parameter choice in a wide range of values (we verified this for k_0 between 0.5 and 0.95).

4) Feature selection criterion

In the new implementation we introduced an automatic data-driven selection of the relevant features to be used for clustering. For each wavelet coefficient, the Lilliefors test returns a test statistic ks_stat . The coefficients with large values of ks_stat tend to be associated with multimodal distributions, and represent the ones that should be selected for clustering. First, we sort the set of ks_stat values in ascending order, leading to the sequence $\{ks_stat_sorted\}$. We empirically found that this sequence tends to increase exponentially. We introduce a smooth estimate of the first derivative of the sequence by computing a difference quotient through a sliding window with a span of 10 samples (i.e. the first quotient is: $(ks_stat_sorted(10) - ks_stat_sorted(1)) / 10$), normalizing by the number of coefficients and the maximum value of ks_stat . Then, we look for the first point where the estimate is larger than 1 for three consecutive samples. This is an estimate of the point colloquially referred to as “the knee of the exponential”, which is indeed where the radius of curvature reaches its minimum. Finally, all the coefficients to the right of the estimated knee are selected, i.e. those where $ks_stat > ks_stat_sorted(knee)$. Examples of the criterion can be seen in Figure 2 and Figure 7A.

We have implemented the automatic feature selection criterion in all cases, except for the real single channel recordings from the human hippocampus (see below), given that in this case, we wanted to compare the automatic sorting with the actions taken by the expert sorters (who did the sorting with a fix number of wavelet coefficients).

Real datasets

1) Single channel: recordings in the human medial temporal lobe

We used recordings in the human medial temporal lobe of 5 patients implanted with depth electrodes for epilepsy diagnosis (Quiñan Quiroga et al. 2008). Each electrode probe had a total of nine microwires at its end, eight active recording channels and one reference. The differential signal from the microwires was amplified by a 64-channel Neuralynx system, filtered between 1 and 9,000 Hz and sampled at 28 kHz.

This dataset was sorted with the former *Wave_clus* implementation (i.e. with a fixed number of 10 wavelets coefficients), which was then manually optimized by an expert. It should be highlighted that with such real data, we do not have ground truth. However, the first goal of the new method is to automatically replicate the steps manually done by the expert to optimize the sorting outcomes. Therefore, the expert’s solution was used as the desired solution to assess the performance of the novel implementation introduced in this work. In other words, the goal with this dataset was to fully automatize the manual optimization

process to bypass the need of the user's intervention. A total of 200 recordings from the human MTL were used (each approximately 15 minutes long), where at least two clusters were identified.

Using the same temperature plot for the expert sorters and the automatic implementation, we could quantify if the latter selected the same clusters (peaks in the temperature plot) as the expert sorters. The True Positive Rate (TPR) was computed as the ratio between the number of clusters identified by the algorithm that corresponded to a cluster selected by the expert and the total number of clusters identified by the expert. In addition, the False Positive Rate (FPR) was calculated by counting the number of clusters identified by the algorithm in the temperature plot that did not correspond to a cluster selected by the expert, divided by the maximum number of false positives generated over all recordings, algorithms and *N_inc* values. This way, for each algorithm and *N_inc* value (ranging from 10 to 60 in steps of 5), we computed the TPR and FPR on the real dataset, allowing us to construct Receiver Operating Characteristic (ROC) curves. In addition, we quantified for each ROC curve the Area Under the Curve (AUC) as a measure of performance of each algorithm.

2) Tetrodes

We also evaluated the performance of the algorithm with tetrode recordings. A tetrode recording from a locust was kindly provided by Ofer Mazar and Gilles Laurent (Perez-Orive et al. 2002). Moreover, we used a set of 20 recordings with tetrodes implanted in 4 patients with intractable epilepsy, using the same procedures as reported in Quiñan Quiroga et al. (2008). Although we do not have ground truth with real data, the solutions from the different algorithms were analyzed by an expert, who assessed whether clusters corresponded to different single units based on spike shape differences and standard criteria, such as cross correlation of spike times and the presence of refractory period violations).

Simulated Datasets

1) Single channel simulations

Two different publicly available simulated datasets were used to evaluate the performance of the proposed method against ground truth. The first one (*Sim1*) was introduced in Quiñan Quiroga et al. (2004) (available at http://www.vis.caltech.edu/~rodri/Wave_clus/Simulator.zip) and it comprises four simulations, each with 3 different neurons, under different levels of background noise. The second set (*Sim2*) was presented in Pedreira et al. (2012) (available at <http://bioweb.me/CPGJNM2012-dataset>), and it includes several simulations done with a varying number of neurons, ranging from 2 to 20. The rationale for using up to 20 neurons was to test the performance of the algorithm in challenging scenarios, especially considering the fact that current spike sorting algorithms tend to detect fewer neurons than they should, based on anatomical and physiological considerations (Henze et al 2000; Pedreira et al 2012). In addition, three independent experts performed manual supervision on this set optimizing the results obtained with the former implementation of *Wave_clus*. We averaged their performance to get a single "Expert" score per simulation.

All the simulated waveforms were clustered using the two different *Wave_clus* implementations and the number of hits and false positives were quantified as in Pedreira et al. (2012) and Niediek et al. (2016). A selected cluster was considered as a hit when more than the 50% of its spikes were correctly identified. Choosing other definitions (i.e. considering more than 70 or 80% of the spikes) gave qualitatively similar results. Selected clusters that were not a hit, were labeled as false positives. Missed clusters were calculated as the number of simulated units minus the number of hits. The multiunit clusters were not considered in the hits calculations but were considered as false positives if less than 50% of their spikes came from the multiunit.

With single channel recordings, we compared the performance of our algorithm to the one given by other recently proposed methods: For performance comparison with other algorithms, we used Klusta (Rossant et al. 2016; [version February 2017](#)) and Combinato (Niediek et al. 2016; [version April 2018](#)). In particular, Klusta is an improved implementation of the former Klustakwik (Harris et al. 2000), a spike sorting algorithm that extracts waveform features through principal components analysis and uses a Gaussian mixture model to perform the clustering of the data. Some of its detection parameters were modified from their default values to achieve a detection performance comparable to the one from *Wave_clus*. Particularly, we used `filter_low` = 300, `filter_high` = 3000, `threshold_strong_std_factor` = 4, `threshold_weak_std_factor` = 3.8, `extract_s_before` = 24, and `extract_s_after` = 40. Its default number of principal components used for feature extraction is 3. However, we observed in our simulations that a larger number of hits was achieved with 10 ($p \sim 10^{-3}$), without differences in terms of the number of false positives, so the results we show for Klusta are based on 10 principal components. Combinato is a recently introduced algorithm that also uses SPC for clustering the data, but includes an iterative re-clustering of large clusters. As the authors tested Combinato with the same simulated dataset we used here, we chose the set of parameters reported in their paper.

2) Tetrodes

We created a set of simulated tetrode recordings using Neurocube (Camuñas-Mesa and Quiroga 2013), a hybrid modelling approach that uses a detailed compartmental model to simulate the contribution of neurons near the recording electrode, and previously recorded spike shapes to generate the background noise (with their amplitude scaled inversely to the squared distance between the neuron and the recording site). To simulate multi-unit activity (with detailed neuron models; 1:4 ratio between interneurons and pyramidal neurons), we used an 8% of active neurons located between 60 and 150 microns away from the tetrode (using a density of 300,000 neurons/mm³), with their spike time intervals drawn from an exponential distribution (mean = 3 Hz). A certain number of pyramidal single units (ranging from 11 to 20) were randomly placed between 10 and 40 microns away from the tetrode with their spike time intervals drawn from an exponential distribution (mean = 3 Hz). The diameter of the electrodes comprising the tetrode was of 20 microns, with 40 microns spacing. Each simulation was 15 minutes long. Waveforms generated by multi-unit neurons that had a peak larger than five times the standard deviation of the noise and were less than 1 ms from the spike of a single unit were subtracted from the signal to reduce the level of closely overlapping spikes. Following spike detection, the spikes associated with the multi-unit activity represented between 50% and 70% of the total number of detected spikes, which is consistent with what is typically observed in real data.

Performance was assessed with the same criteria of hits and false positives used with the single channel dataset. We compared our results with other recently proposed algorithms that were specifically developed to sort multichannel recordings: MountainSort (Chung et al. 2017; [version 0.11.6](#)), Kilosort (Pachitariu et al. 2016; [version February 2018](#)), and SpyKING CIRCUS (Yger et al. 2018; [version 0.6.4](#)). We also included Klusta, for which we used the same parameters as in the single channel simulated set. For the other algorithms we used their default parameters (in SpiKING CIRCUS, we enabled “auto-mode merging” as suggested for a fully automatic implementation, and did not use the spatial whitening option as it led to a large decrease in detection performance).

RESULTS

Description of the new sorting implementation

The key advantage of the new Wave_clus implementation is that it automatically selects clusters from different temperatures. This allows identifying units with different firing rates and spike shape characteristics that differ in their density and location within the feature space. To illustrate this point, Figure 1A shows a segment of real data (high-passed filtered), where the spikes associated with different units can be observed. Based on the different data partitions generated by the SPC algorithm, which are shown in the temperature plot on the bottom left of Figure 1A, an expert user selected 4 clusters (3 single units and 1 multi-unit) at different temperatures. To achieve this automatically, a three-step procedure is used (see Materials and Methods for details). First, putative clusters are selected at different temperatures. It should be noted that a similar idea of finding peaks at different temperatures in the temperature plot has also been introduced in Niediek et al. (2016). Second, peaks from very high temperatures are eliminated, to avoid overclustering. Figure 1B shows an example of the identification of T_B , with the grey shaded area indicating all the temperatures that are disregarded and which may have contributed with spurious peaks. Third, the algorithm quantifies the overlap between pairs of clusters to avoid double counting when considering the same spikes at different temperatures. Figure 1C shows an example with three peaks, which in this case were obtained because the larger cluster at $T = 0.06$ was split in two at $T = 0.07$. The overlapping coefficient between the clusters at both temperatures was equal to 1 in both cases—thus correctly showing that these three peaks correspond to the same cluster split in two—and therefore the clusters at $T = 0.07$ were retained and the one at $T = 0.06$ was discarded.

A final improvement compared to the former implementation is that the new algorithm automatically selects the number of relevant features (i.e. wavelet coefficients) to be used for clustering. Previously, we selected a fixed number of coefficients—10 out of a total of 64 wavelet coefficients—which are the ones differing the most from a Gaussian distribution, according to the Lilliefors test. Figure 2 shows the number of coefficients selected with the proposed method on each simulation of dataset *Sim2* (see Materials and Methods). This dataset comprises 5 simulations with N single units, with N ranging from 2 to 20. We observe that the number of selected coefficients increases with the number of neurons in the simulation (Spearman correlation, $\rho = 0.42$, $p \sim 10^{-5}$). In other words, with few neurons,

relatively few coefficients will be enough to separate the associated clusters, but as more neurons are present in a recording, it is necessary to consider more coefficients. However, if we use all 64 coefficients, we will also consider those that capture just noise, thus decreasing clustering performance (see below). The proposed criterion automatically selects the coefficients in a data-driven way. Three particular examples can be seen in the inset of the figure, showing the dynamics of the statistic of the Lilliefors test and the estimated “knee of the exponential” (see Materials and Methods), from which all coefficients to the right were selected.

Performance with real single channel recordings

We first evaluated the performance of the new method with real data collected from recordings in the human medial temporal lobe (see Materials and Methods). The data set comprises 200 recordings with at least two units identified by an expert user. We used N_{inc} values from 10 to 60 in steps of 5. This way, as shown in Figure 3A, we could compute an ROC curve for each sorting implementation we tested (see Materials and Methods).

With the former *Wave_clus* implementation (dashed line), increasing N_{inc} led to an increase of the true positive rate (TPR), with a reduction of the false positive rate (FPR), until $N_{inc} = 25$; then the TPR stabilized at 75%–80%. When the proposed peak selection method was included (dotted line), low values of N_{inc} were associated with high TPRs, but also with high FPRs. Increasing N_{inc} reduced the number of false positives (at the expense of also reducing the TPR), though the FPR remained relatively large when compared with the former *Wave_clus* implementation. The addition of the inclusion criterion led to the ROC being shifted to the left (dash-dotted line), but with a reduction in the TPR (although still with TPRs 10% higher than with the former *Wave_clus* implementation) due to overclustering. Finally, the addition of the regime border detection reduced the chance of overclustering (solid line), boosting the TPR an extra 5% while maintaining low FPRs. Figure 3B shows the statistics on the Area Under the Curve (AUC) associated with the different sorting implementations. An AUC closer to 1 is related to a better performance. The addition of each of the proposed improvements lead to a significant increase in the AUC with the new implementation including the peak selection, inclusion, and regime border detection criteria being the one associated with the best performance. In Figure 3C, we further evaluated the contribution of each step in terms of the number of total errors, i.e. misses + false positives. From the former *Wave_clus* implementation, the addition of the peaks selection alone actually led to a significant worsening of the performance due to the high number of false positives generated (notice how the ROC in Figure 3A is shifted to the right). However, when the inclusion criterion was also added, the number of errors dropped significantly, outperforming both peaks and the former *Wave_clus* implementation. Finally, the further addition of the regime border detection led to a further improvement in performance, with the resulting new implementation exhibiting significantly lower number of error than any of the other alternatives.

This dataset allowed us to evaluate the performance of the improvements associated with the actions taken by the experts and to assess the contribution of each of them. As a result, we found that by choosing $N_{inc} = 20$ with the new implementation we automatically detected almost 95% of the clusters manually detected by the expert sorters with an average

of less than one false positive. This represents a 15% increase in TPR with a similar low FPR when compared against the former *Wave_clus* implementation.

Performance with simulated single channel recordings

Results with dataset Sim1

The advantage of using simulated data is that it provides ground truth to evaluate the performance of a spike sorting implementation. First we used the dataset introduced in Quian Quiroga et al. (2004), where three neurons were simulated in each recording with different waveforms (in the simulations named “Easy”, waveforms were very different from each other, whereas in the ones labeled “Difficult” they were very similar) and with different levels of background noise. For each simulation, we evaluated the number of hits (see Materials and Methods) with the former and new *Wave_clus* implementations using a variable number of wavelet coefficients (with the new implementation using the peak selection, inclusion, and regime border detection criteria introduced in the section “Description of the new sorting implementation”). Table 1 shows that, as shown in Quian Quiroga et al. (2004), the former *Wave_clus* implementation using 10 WC (wavelet coefficients) was able to cluster all three simulated units in all but one case (the hardest one in terms of noise level and waveform similarity). However, when no dimensionality reduction was performed, i.e., when all 64 WC were used, the performance dropped significantly (former 10 WC vs. 64 WC, paired sign test, $p = 0.008$). This drop is particularly evident for the cases with large noise levels, as several coefficients captured noise features, compromising the performance of the clustering algorithm. The new implementation exhibited a similar performance with 10 WC and also showed an analogous drop in performance when all 64 WC were used (new 10 WC vs. 64 WC, paired sign test, $p = 0.016$). Finally, when the selection of a variable number of WC was used, the performance was the same as with 10 WC. As discussed below, this is due to the fact that for relatively low number of neurons 10 WC already resulted in nearly optimal performance.

Figure 4 shows the number of automatically selected coefficients for each simulation. First, the “easy” simulations have more coefficients selected than the “difficult” ones; i.e. as the waveforms are easier to separate, more coefficients give information that differentiates between the different waveforms. Note that for each set of simulations with varying degree of noise, more coefficients were selected for the cases with relatively low noise levels (i.e. simulations labeled ‘n05’). Particularly, for the simulations “Easy1” (performed with 8 different levels of noise), there was a strong correlation between noise level and the number of coefficients selected (Spearman correlation, $\rho = 0.96$, $p \sim 10^{-4}$). The other simulations have only 4 levels of noise, which is not sufficient to properly estimate the correlation. This trend can be explained by the fact that when the noise level is low, several coefficients capture differences between spike shapes, whereas with larger noise levels, more coefficients capture noise features that do not contribute to the separation of clusters and are therefore (correctly) not selected. Importantly, the new implementation with a variable number of WC has no a priori upper bound on the number of coefficients selected. Yet, Table 1 and Figure 4 show that the number chosen did not go too high for the high noise conditions, preventing the drop of performance observed when all coefficients were selected (64WC).

Results with dataset Sim2

Dataset Sim 1 consisted of 3 simulated neurons of different spike shapes and with varying levels of noise. To test the performance of the algorithm in more challenging scenarios, we also used the dataset Sim 2, which consisted of 95 simulations, each with 2 to 20 neurons (see Materials and Methods). Three independent experts performed a manual optimization of the solution provided by the former implementation of *Wave_clus*. Panels A and C in Figure 5 show the performance of the algorithms in terms of hits and misses. As expected, the experts significantly improved the performance over the automatic solution from the former *Wave_clus* implementation, which suffered from the selection of a single temperature in the temperature plot. When using a fixed number of wavelet coefficients (10), the unsupervised performance of the proposed method was not significantly different (both in terms of misses and false positives) to the one obtained by the experts (with the former *Wave_clus* implementation), thus showing that the ad hoc improvements introduced to mimic the expert's actions were successful to deliver a fully unsupervised algorithm. It should be noticed that using only the peak selection step (without the inclusion and regime border detection) significantly deteriorated the performance of the algorithm due to a large number of false positives (as it was the case with the real data presented in Figure 3). Conversely, when the new implementation included the new feature selection criterion (i.e. using a variable number of wavelet coefficients), a further significant increase in the number of hits was achieved, outperforming the results obtained by the experts. In fact, the performance obtained by the experts saturates at about 8 hits, whereas the new algorithm managed to correctly isolate an average of 14 out of 20 simulated neurons from a single channel. Moreover, the use of variable coefficients led to significant improvements over the case with 10 WC, in terms of both misses ("2 to 10", $p \sim 10^{-3}$; "11 to 20", $p \sim 10^{-12}$) and false positives ("2 to 10", not significant; "11 to 20", $p \sim 10^{-3}$).

Comparison to other single channel algorithms

Figure 5 also shows the comparison of performance with other recently introduced algorithms: Klusta and Combinato (see Materials and Methods). All the different implementations were statistically compared using paired sign tests. As the simulations contained more neurons, there was also more room for differentiating the performance between algorithms. For this reason, panels C and D in Figure 5 evaluate the performance (using paired sign tests) separating the cases with a small ("2 to 10") and large ("11 to 20") number of neurons. When comparing the performance of Klusta and the experts, Klusta showed more hits in the set "11 to 20" ($p \sim 10^{-4}$). In turn, Combinato showed no significant differences in the number of hits compared to the experts (although $p = 5.4 \times 10^{-2}$ in the set "11 to 20"). Klusta and Combinato showed similar performance, with Klusta being better in the set "11 to 20" ($p = 2.3 \times 10^{-2}$). The new implementation of *Wave_clus* with a variable number of wavelet coefficients, gave a further significant increase in the number of hits, outperforming all other methods with small and large number of neurons (Figure 5C).

Next, we considered false positives. As shown in panels B and D in Figure 5, Klusta and Combinato showed a significant increase in the number of false positives when compared to all other algorithms, regardless of whether the number of simulated neurons was small or

large. In the set “11 to 20”, the new *Wave_clus* implementation achieved a significantly smaller number of false positives than all other algorithms, with reductions of more than 50% when compared to Klusta and Combinato.

Finally, we analyzed the computational cost of the different algorithms by using recordings with 20,000 spikes ($\pm 1\%$). Combinato detected 13,500 spikes, whereas the other algorithms detected 15,000 (in all cases, nearly all the single unit spikes were detected). Combinato shows larger computing times than the former *Wave_clus* implementation. This is mainly due to the fact that every cluster identified with more than 1,000 spikes was subjected to a second run of SPC. The new *Wave_clus* with 10 wavelet coefficients (i.e. including the steps to get a fully automatic implementation and the improvements used for spike interpolation and wavelet decomposition) led to a large drop in computing time (~ 20 seconds on average). When the new feature extraction was included, computing times were increased by 6 seconds on average, but were still smaller than those from the former *Wave_clus* and Combinato. The difference with the new implementation with 10 wavelet coefficients is fully accounted by the time taken by SPC to cluster the spikes in a higher dimensional space (in this set, an average of 20 wavelet coefficients were used in the new *Wave_clus*). However, the small additional computing time was accompanied by a reduction of misses and false positives by almost 50%. When compared with Klusta and Combinato, reductions in misses and false positives were even larger, and with less computing time. It should be noticed that if recordings with more spikes were used, the computing time will not increase much, as 20,000 spikes can be used for clustering with SPC and the remaining ones can be assigned afterwards by template matching.

Performance with real tetrode recordings

Results with an example tetrode recording in the locust

We evaluated the performance of the proposed algorithm using real data recorded with tetrodes. Figure 6A shows the 7 clusters isolated by the proposed algorithm from a tetrode recording in a locust, which display clear differences in the spike shapes in the different channels. In contrast, the former automatic *Wave_clus* implementation only isolated 2 clusters, whereas following manual supervision it was possible to isolate 6 clusters (C5 and C7 remained merged in this case). Only 20 wavelet coefficients were automatically selected in the proposed implementation, in contrast to the 40 selected in the former implementation. Therefore, the 20 additional coefficients can be carrying information about noise, preventing the SPC from separating clusters C5 and C7.

Results with tetrode recordings in the human hippocampus

Next, we analyzed a set of tetrode recordings in the human hippocampus. Figure 6B shows one example from that set, where the new *Wave_clus* implementation was able to isolate 10 clusters. The former automatic *Wave_clus* led to 2 clusters of poorly isolated units (one with “small” spikes and another with “large” spikes); following manual supervision it was possible to increase the number of isolated clusters to 7. The new implementation automatically selected 57 wavelet coefficients, with the 17 additional coefficients allowing *Wave_clus* to isolate more clusters. Figure 6C presents the analysis over all 20 tetrode recordings.

Although it is not possible to have ground truth with this data, we observe that the new implementation led to a larger number of identified putative single units (according to the expert opinion), not only in comparison to the former automatic *Wave_clus* implementation, but also with respect to the former implementation optimized using manual supervision (but still using a fixed number of features). This latter difference can be attributed to the automatic selection of a variable number of wavelet coefficients.

Performance with simulated tetrode recordings

To quantify performance with tetrodes, we created 10 simulations including between 11 and 20 single units (see Materials and Methods). The top part of Figure 7A shows a diagram of the simulated tetrode, where each electrode had 20 microns diameter with a separation of 40 microns between them. The bottom part of Figure 7A shows the sequence of sorted values of the statistic *ks_stat* obtained from the Lilliefors test for the simulation with 11 single units. A total of 72 coefficients were automatically selected, which enabled the new *Wave_clus* to correctly isolate all 11 units without any false positives (the waveforms for the resulting clusters are shown in Figure 7B).

Comparison to other algorithms

To compare the performance of the different algorithms over the simulated dataset, we analyzed the number of hits/misses and false positives, as previously done for the single channel dataset. All algorithms detected a similar number of spikes, capturing above 90% of the ones generated by the single units. Figure 7C shows that the proposed algorithm had a significant lower number of misses compared to Klusta, SpyKING CIRCUS, and Kilosort ($p = 0.002$, $p = 0.004$, and $p = 0.039$, respectively), and a similar performance compared to MountainSort (the difference was not significant). However, as shown in Figure 7D, differences with MountainSort and with all other algorithms were significant when considering the number of false positives. To quantify the overall performance, we computed the total number of errors. Figure 7E shows that the proposed algorithm had significantly fewer errors than all other algorithms (with at least half the number of errors compared to the algorithm with the second-best performance). Interestingly, there were no significant differences between MountainSort, Kilosort, and SpyKING CIRCUS ($p > 0.34$).

DISCUSSION

We presented a fully automatic method that performed significantly better than a former *Wave_clus* implementation and also compared to other recently proposed unsupervised algorithms. Most of the modifications introduced are actually inspired by the actions performed by experts optimizing the automatic solutions given by the former implementation. In addition, we introduced a criterion for selecting a variable number of wavelet coefficients, which allowed discriminating more units, but keeping a low number of features when a small number of units was present in the recording or when the background noise was high. In fact, using all the coefficients under low signal to noise ratios poses a problem for clustering methods in high-dimensional spaces, which is associated with the “curse of dimensionality” (Bishop 2006). Importantly, the improved performance achieved by the new dimensionality

reduction was done with an effective implementation that lead to computing times that where smaller than those of other algorithms.

Using data from real single channel recordings, the new algorithm successfully detected, in a fully unsupervised way, about 95% of the clusters obtained by expert sorters, who manually optimized the solution of the former *Wave_clus* implementation. This was accomplished while maintaining a low number of false positives. When tested with single channel simulated data, the new implementation with 10 wavelet coefficients achieved the same performance as the one from the expert sorters, showing that the modifications introduced accomplished a successful automatization. However, when the variable number of wavelet coefficients was used, the proposed algorithm was able to correctly detect more than the 8 neurons identified by expert sorters using the former *Wave_clus* (while maintaining a low number of false positives). The performance with the new *Wave_clus* implementation was also significantly better than the one obtained with other algorithms: Klusta and Combinato.

A good performance of the proposed algorithm was also observed with real and simulated data from tetrode recordings, and was significantly better than the one obtained with several recently introduced algorithms designed for multichannel recordings: Klusta, MountainSort, Kilosort, and SpyKING CIRCUS. In addition, the implementation of *Wave_clus* is modular, in the sense that the feature extraction and clustering steps are decoupled. Therefore, it is possible to, for example, maintain the proposed method to select wavelet coefficients while replacing SPC with a different clustering algorithm, such as the ones used by MountainSort (Chung et al. 2017) or SpyKING CIRCUS (Rodriguez and Laio 2014). Alternatively, a different method for feature extraction might be introduced while maintaining the SPC and the steps introduced here to provide a fully automatic solution.

Although the new implementation of *Wave_clus* showed good performance with real data, there is still room for improvement. When performing long-term recordings, the stability of the spike waveforms can be affected by different causes, such as electrode drift or changes in the recording conditions. Different solutions have been proposed to tackle this issue (Rey et al. 2015b). In fact, Combinato (Niediek et al. 2016) offers a novel approach to reliably track neurons over long periods of time, and it is a subject of further investigation to develop and compare different methodologies to optimally track neurons over days. Another problem that is observed in real data is due to overlapping spikes (Ekanadham et al. 2014; Franke et al. 2015; Rey et al. 2015b). This issue is normally diminished by simultaneously recording from different channels located close to each other (so that the overlap is not seen in all of them). The new *Wave_clus* implementation could be used in this context, although further research is required to find a good way for grouping the information from different channels and then properly combining the results. Finally, quality metrics can be incorporated to the sorting algorithm (Joshua et al. 2007; Rey et al. 2015b; Harris et al. 2016) and spatial pre-processing can be considered for large electrode arrays (Musial et al. 2002; Huang and Miller 2005).

The development of automatic and reliable spike sorting algorithms is becoming critical, given that within the next 10 years we will likely witness the number of recording sites going up to thousands (Stevenson and Kording 2011; Alivisatos et al. 2013). The amount and complexity of the data to be produced by the next generation of probes is too large to be

handled by researchers in a supervised way. In this context, the only viable option to fully take advantage of technological developments is to accompany them with the development of easy-to-use and properly validated tools for fully automatic spike sorting (Einevoll et al. 2012; Harris et al. 2016). In this respect, we have here presented a fully automatic algorithm that not only matched but also outperformed the performance of supervised algorithms and outperformed, as well, the performance of previously proposed unsupervised implementations.

ACKNOWLEDGEMENTS

We specially thank Ofer Mazor, Gilles Laurent for kindly providing exemplary multichannel recordings. This research was supported by the Medical Research Council (G1002100) and the Human Frontiers Research Program.

DISCLOSURES

No conflicts of interest, financial or otherwise, are declared by the author(s).

REFERENCES

- Alivisatos AP, Andrews AM, Boyden ES, Chun M, Church GM, Deisseroth K, Donoghue JP, Fraser SE, Lippincott-Schwartz J, Looger LL, Masmanidis S, McEuen PL, Nurmikko AV, Park H, Peterka DS, Reid C, Roukes ML, Scherer A, Schnitzer M, Sejnowski TJ, Shepard KL, Tsao D, Turrigiano G, Weiss PS, Xu C, Yuste R, Zhuang X. Nanotools for Neuroscience and Brain Activity Mapping. *ACS Nano* 7: 1850–1866, 2013.
- Berényi A, Somogyvári Z, Nagy AJ, Roux L, Long JD, Fujisawa S, Stark E, Leonardo A, Harris TD, Buzsáki G. Large-scale, high-density (up to 512 channels) recording of local circuits in behaving animals. *J Neurophysiol* 111: 1132–1149, 2014.
- Bishop CM. *Pattern recognition and machine learning*. New York: Springer, 2006.
- Blanche TJ, Spacek MA, Hetke JF, Swindale NV. Polytrodes: High-Density Silicon Electrode Arrays for Large-Scale Multiunit Recording. *J Neurophysiol* 93: 2987–3000, 2005.
- Blatt M, Wiseman S, Domany E. Superparamagnetic clustering of data. *Phys Rev Lett* 76: 3251, 1996.
- Blatt M, Wiseman S, Domany E. Data Clustering Using a Model Granular Magnet. *Neural Comput* 9: 1805–1842, 1997.
- Buzsáki G. Large-scale recording of neuronal ensembles. *Nat Neurosci* 7: 446–451, 2004.
- Camuñas-Mesa LA, Quiroga RQ. A Detailed and Fast Model of Extracellular Recordings. *Neural Comput* 25: 1191–1212, 2013.
- Chung JE, Magland JF, Barnett AH, Tolosa VM, Tooker AC, Lee KY, Shah KG, Felix SH, Frank LM, Greengard LF. A Fully Automated Approach to Spike Sorting. *Neuron* 95: 1381–1394.e6, 2017.

- 739 Csicsvari J. Massively Parallel Recording of Unit and Local Field Potentials With Silicon-
740 Based Electrodes. *J Neurophysiol* 90: 1314–1323, 2003.
- 741 De Falco E, Ison MJ, Fried I, Quiñero R. Long-term coding of personal and universal
742 associations underlying the memory web in the human brain. *Nat Commun* 7: 13408, 2016.
- 743 Domany E, Blatt M, Gdalyahu Y, Weinshall D. Superparamagnetic clustering of data:
744 application to computer vision. *Comput Phys Commun* 121: 5–12, 1999.
- 745 Einevoll GT, Franke F, Hagen E, Pouzat C, Harris KD. Towards reliable spike-train
746 recordings from thousands of neurons with multielectrodes. *Curr Opin Neurobiol* 22: 11–17,
747 2012.
- 748 Ekanadham C, Tranchina D, Simoncelli EP. A unified framework and method for automatic
749 neural spike identification. *J Neurosci Methods* 222: 47–55, 2014.
- 750 Fee MS, Mitra PP, Kleinfeld D. Variability of extracellular spike waveforms of cortical
751 neurons. *J Neurophysiol* 76: 3823–3833, 1996.
- 752 Franke F, Pröpper R, Alle H, Meier P, Geiger JRP, Obermayer K, Munk MHJ. Spike sorting
753 of synchronous spikes from local neuron ensembles. *J Neurophysiol* 114: 2535–2549, 2015.
- 754 Frey U, Egert U, Heer F, Hafizovic S, Hierlemann A. Microelectronic system for high-
755 resolution mapping of extracellular electric fields applied to brain slices. *Biosens Bioelectron*
756 24: 2191–2198, 2009.
- 757 Gold C, Henze DA, Koch C, Buzsáki G. On the origin of the extracellular action potential
758 waveform: A modeling study. *J Neurophysiol* 95: 3113–3128, 2006.
- 759 Harris KD, Henze DA, Csicsvari J, Hirase H, Buzsáki G. Accuracy of Tetrode Spike
760 Separation as Determined by Simultaneous Intracellular and Extracellular Measurements. *J*
761 *Neurophysiol* 84: 401–414, 2000.
- 762 Harris KD, Quiñero R, Freeman J, Smith SL. Improving data quality in neuronal
763 population recordings. *Nat Neurosci* 19: 1165–1174, 2016.
- 764 Henze DA, Borhegyi Z, Csicsvari J, Mamiya A, Harris KD, Buzsáki G. Intracellular Features
765 Predicted by Extracellular Recordings in the Hippocampus In Vivo. *J Neurophysiol* 84: 390–
766 400, 2000.
- 767 Homer ML, Nurmikko AV, Donoghue JP, Hochberg LR. Sensors and Decoding for
768 Intracortical Brain Computer Interfaces. *Annu Rev Biomed Eng* 15: 383–405, 2013.
- 769 Huang Y, Miller JP. Phased array processing for spike discrimination. *Neurocomputing* 65–
770 66: 507–516, 2005.
- 771 Joshua M, Elias S, Levine O, Bergman H. Quantifying the isolation quality of extracellularly
772 recorded action potentials. *J Neurosci Methods* 163: 267–282, 2007.
- 773 Lambacher A, Vitzthum V, Zeitler R, Eickenscheidt M, Eversmann B, Thewes R, Fromherz
774 P. Identifying firing mammalian neurons in networks with high-resolution multi-transistor
775 array (MTA). *Appl Phys A* 102: 1–11, 2011.
- 776 Lewicki MS. A review of methods for spike sorting: the detection and classification of neural
777 action potentials. *Netw Comput Neural Syst* 9: R53–R78, 1998.

778 Litke AM, Bezayiff N, Chichilnisky EJ, Cunningham W, Dabrowski W, Grillo AA, Grivich M,
 779 Grybos P, Hottowy P, Kachiguine S, Kalmar RS, Mathieson K, Petrusca D, Rahman M, Sher
 780 A. What does the eye tell the brain?: Development of a system for the large-scale recording
 781 of retinal output activity. *IEEE Trans Nucl Sci* 51: 1434–1440, 2004.

782 Musial P., Baker S., Gerstein G., King E., Keating J. Signal-to-noise ratio improvement in
 783 multiple electrode recording. *J Neurosci Methods* 115: 29–43, 2002.

784 Niediek J, Boström J, Elger CE, Mormann F. Reliable Analysis of Single-Unit Recordings
 785 from the Human Brain under Noisy Conditions: Tracking Neurons over Hours. *PLOS ONE*
 786 11: e0166598, 2016.

787 Pachitariu M, Steinmetz N, Kadir S, Carandini M, Harris KD. Kilosort: realtime spike-sorting
 788 for extracellular electrophysiology with hundreds of channels. *bioRxiv*. .

789 Pedreira C, Martinez J, Ison MJ, Quiñ Quiroga R. How many neurons can we see with
 790 current spike sorting algorithms? *J Neurosci Methods* 211: 58–65, 2012.

791 Perez-Orive J, Mazor O, Turner GC, Cassenaer S, Wilson RI, Laurent G. Oscillations and
 792 sparsening of odor representations in the mushroom body. *Science* 297: 359–365, 2002.

793 Quiñ Quiroga R. Spike sorting. *Scholarpedia* 2: 3583, 2007.

794 Quiñ Quiroga R, Mukamel R, Isham EA, Malach R, Fried I. Human single-neuron
 795 responses at the threshold of conscious recognition. *Proc Natl Acad Sci* 105: 3599–3604,
 796 2008.

797 Quiñ Quiroga R, Nadasdy Z, Ben-Shaul Y. Unsupervised spike detection and sorting with
 798 wavelets and superparamagnetic clustering. *Neural Comput* 16: 1661–1687, 2004.

799 Quiñ Quiroga R, Panzeri S. Extracting information from neuronal populations: information
 800 theory and decoding approaches. *Nat Rev Neurosci* 10: 173–185, 2009.

801 Redish AD, Battaglia FP, Chawla MK, Ekstrom AD, Gerrard JL, Lipa P, Rosenzweig ES,
 802 Worley PF, Guzowski JF, McNaughton BL, Barnes CA. Independence of Firing Correlates of
 803 Anatomically Proximate Hippocampal Pyramidal Cells. *J Neurosci* 21: RC134–RC134, 2001.

804 Rey HG, Ison MJ, Pedreira C, Valentin A, Alarcon G, Selway R, Richardson MP, Quiñ
 805 Quiroga R. Single-cell recordings in the human medial temporal lobe. *J Anat* 227: 394–408,
 806 2015a.

807 Rey HG, Pedreira C, Quiñ Quiroga R. Past, present and future of spike sorting techniques.
 808 *Brain Res Bull* 119: 106–117, 2015b.

809 Rodriguez A, Laio A. Clustering by fast search and find of density peaks. *Science* 344:
 810 1492–1496, 2014.

811 Rossant C, Kadir SN, Goodman DFM, Schulman J, Hunter MLD, Saleem AB, Grosmark A,
 812 Belluscio M, Denfield GH, Ecker AS, Tolia AS, Solomon S, Buzsáki G, Carandini M, Harris
 813 KD. Spike sorting for large, dense electrode arrays. *Nat Neurosci* 19: 634–641, 2016.

814 Stevenson IH, Kording KP. How advances in neural recording affect data analysis. *Nat*
 815 *Neurosci* 14: 139–142, 2011.

816 Wood F, Black MJ, Vargas-Irwin C, Fellows M, Donoghue JP. On the variability of manual
817 spike sorting. *IEEE Trans Biomed Eng* 51: 912–918, 2004.

818 Yger P, Spampinato GL, Esposito E, Lefebvre B, Deny S, Gardella C, Stimberg M, Jetter F,
819 Zeck G, Picaud S, Duebel J, Marre O. A spike sorting toolbox for up to thousands of
820 electrodes validated with ground truth recordings in vitro and in vivo. *eLife* 7: e34518, 2018.

821

FIGURE CAPTIONS

Figure 1: Improvements of the proposed method. A. Example of an extracellular recording from the human right entorhinal cortex. The black trace in the top panel shows the high frequency content of the signal (between 300 and 3000 Hz in this example), where there is a superposition of several effects. Neurons located more than ~ 150 microns away from the tip of the electrode contribute to the background noise, so their spikes cannot be detected. Closer neurons (between 50 and 150 microns away from the tip of the electrode) generate spikes larger than the background noise, but they cannot be isolated into different units, thus being associated with the multiunit activity (cluster 1). Finally, nearby neurons (less than ~ 50 microns) have even larger spikes, and sorting algorithms allow us to assign the recorded spikes to the different neurons that generated them (clusters 2–4). On the bottom-left panel, the temperature map shows all the different partitions generated by the SPC. Each line is associated with the size of the k th cluster ranked by size at each temperature. The filled circles denote the location of the clusters depicted on the right, which have been selected by an expert user who manually optimized the solution of the former *Wave_clus*. In the top panel, the time of occurrence of each spike is marked with a triangle associated with the four isolated clusters. B. Example of the detection of the border of the superparamagnetic regime. T_B was identified at $T = 0.16$, so the partitions at $T_i \geq T_B$ were discarded (gray shaded area). C. Example of the inclusion criterion when a cluster was split in two at a higher temperature. The overlap coefficient between the marked clusters at $T = 0.07$ and the marked cluster at $T = 0.06$ was equal to one, so the clusters at $T = 0.07$ are retained, whereas the one at $T = 0.06$ is discarded.

Figure 2: New criterion for feature extraction applied on the dataset *Sim2*. On the left, the number of selected wavelet coefficients is shown for each simulation, with the simulations being sorted according to the number of simulated single units. There is a clear correlation between the number of chosen coefficients and the number of simulated units (Spearman correlation, $\rho = 0.42$, $p \sim 10^{-5}$). On the right, three particular examples are shown, depicting the sequence of *ks_stat*, i.e. the statistic associated with the Lilliefors test. The shaded grey area covers the coefficients that were selected on each particular example. Notice the change in the slope of the curve at the point marked by the vertical dashed line.

Figure 3: Performance of the proposed improvements with a real dataset. A. For each algorithm implementation, the parameter N_{inc} took values from 10 to 60 in steps of 5 (as N_{inc} is increased, we moved on the figure from right to left). In each case, the true and false positive rates were computed, allowing us to plot mean and SEM across the 200 recordings analyzed. For the inclusion criterion, we used $k_0 = 0.9$, and the border of the superparamagnetic regime was computed using $Thr_{border} = 0.4$. The new *Wave_clus* implementation (solid line) includes the peak selection, inclusion, and regime border detection criteria. The black circle marks its performance with $N_{inc} = 20$, which is the final value chosen for the new implementation. B. Mean and standard error of the mean of the Area Under the Curve across the 200 recordings. Paired sign tests were used to evaluate the statistical difference across implementations. C. Mean and standard error of the mean of the total number of errors (i.e. misses + false positives) across the 200 recordings. Paired sign tests were used to evaluate the statistical difference across implementations.

Figure 4: Selected wavelet coefficients for each simulation in the simulated dataset *Sim1*. Each row represents a simulation with a given level of noise (e.g. E1_n10 depicts the simulation from the set “Easy1” with a noise level of 10). The number of selected coefficients is reduced as the noise level is increased. Particularly, for the simulations Easy1, created with 8 different noise levels, there was a significant correlation (Spearman correlation, $\rho = 0.96$, $p \sim 10^{-4}$).

Figure 5: Performance of the proposed algorithm with the simulated dataset *Sim2* (single channel). A. Number of hits as a function of the number of simulated neurons for the different algorithm implementations. Mean and SEM across the 5 simulations with each number of simulated neurons is shown. The former *Wave_clus* was used in its automatic form or followed by manual supervision by an expert sorter. The curve for the new implementation with a fixed number of 10 wavelet coefficients is largely overlapping with the one from the experts. The new proposed method was also compared with Klusta and Combinato, recently introduced algorithms that are suitable for sorting single channel recordings. B. Same as in A but for false positives. C. Mean and SEM of misses for each algorithm. The analysis was done for the subsets with a small (“2 to 10”) and large (“11 to 20”) number of neurons. Paired sign tests were used to evaluate the statistical difference across implementations. All the implementations showed better performance than the unsupervised former *Wave_clus*, whereas the proposed algorithm with a variable number of wavelet coefficients was significantly better than all other implementations. D. Same as C but with false positives. Significant differences were found in the subset with small number of neurons when comparing with Klusta and Combinato, but with more neurons, the proposed algorithm with a variable number of wavelet coefficients was significantly better than all other implementations. E. Computing time of each algorithm using simulated single channel recordings with 20,000 ($\pm 1\%$) spikes each. The new proposed algorithm has smaller computing times than Klusta and Combinato, while achieving a much better performance, both in terms of hits and false positives.

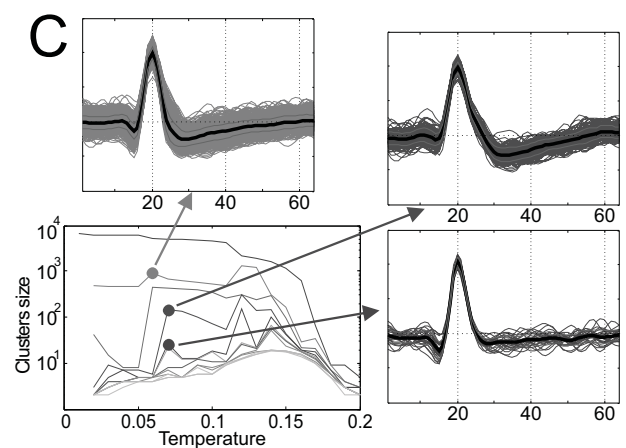
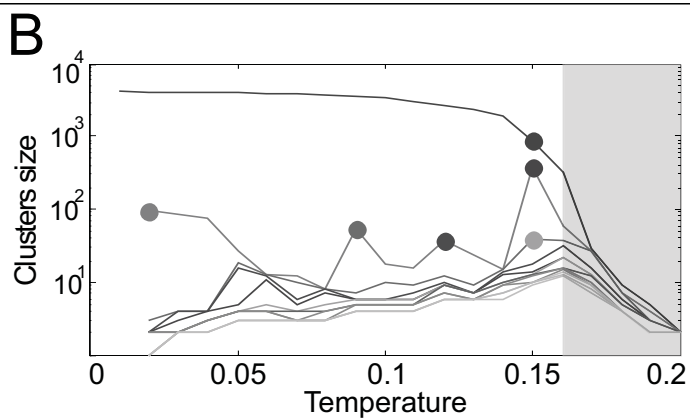
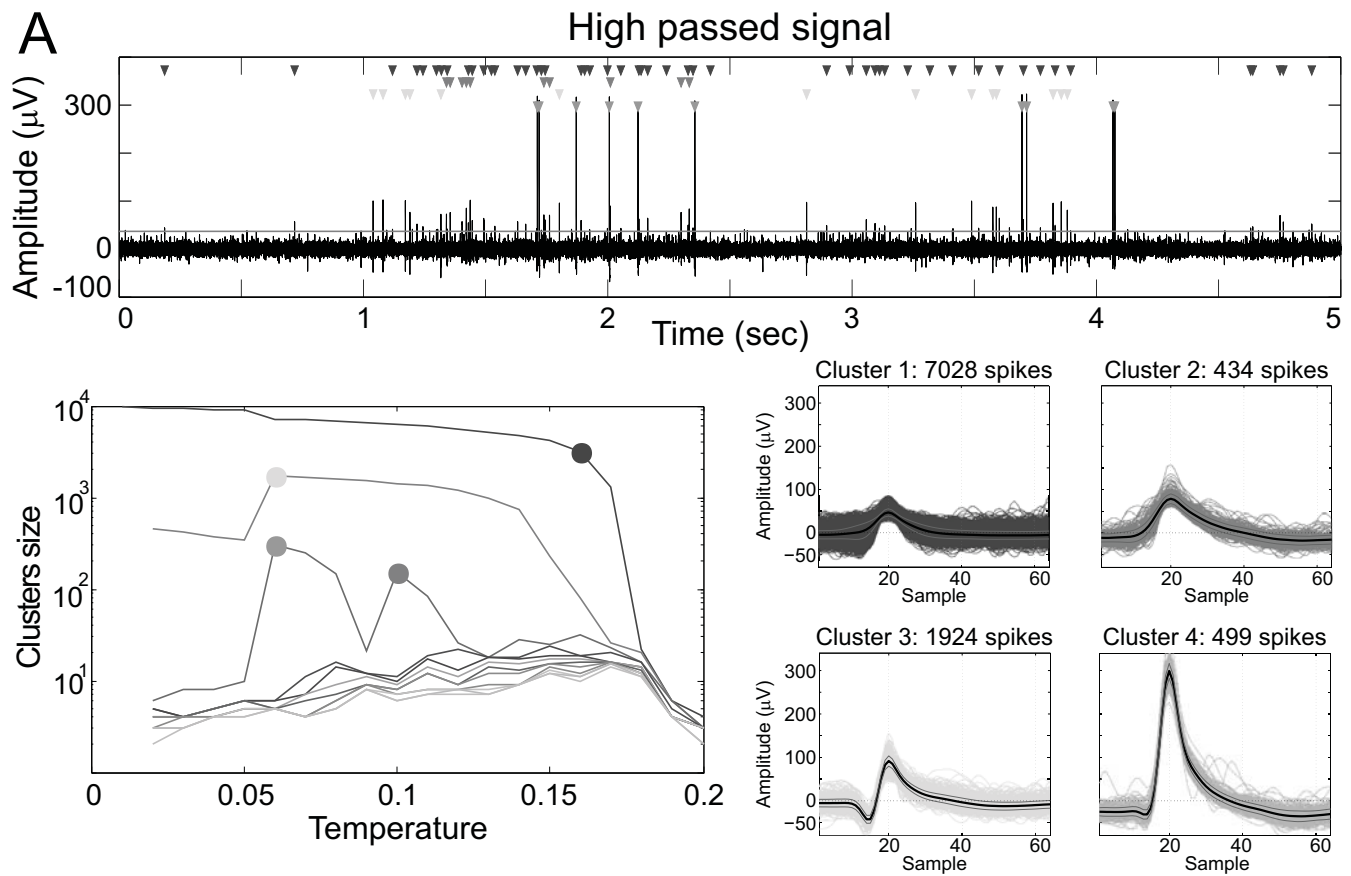
Figure 6: Performance of the proposed algorithm with real data from tetrodes. A. In a tetrode recording from a locust, the proposed algorithm was able to isolate 7 clusters (which would be associated to putative single units according to an expert opinion; see Methods for details). Spikes for each class on each recording channel are overlaid (total number of spikes per cluster shown in brackets), with the thick black line representing the mean waveform and the lighter ones associated with the SEM. Horizontal dotted lines represent the voltage $y = 0$. Ticks on the x axis are placed every 10 samples, whereas the separation in the y axis is 1, 0.5, 0.2, and 1 for channels 1 to 4 (in arbitrary units). B. In a tetrode recording from the human hippocampus, the proposed algorithm was able to isolate 10 clusters. Ticks on the x axis are placed every 20 samples, whereas the separation in the y axis is 50 μV for all channels. C. The analysis of a set of 20 tetrode recordings in the human hippocampus showed that the proposed algorithm was able to isolate significantly more putative units than the former implementation, even when followed by manual supervision from expert sorters. Mean \pm SEM and boxplots (median \pm 1st/3rd quartiles) are shown, as well as the results from paired sign tests. Inset on the top left shows the number of clusters isolated on each recording by the new-automatic and former-manual implementations.

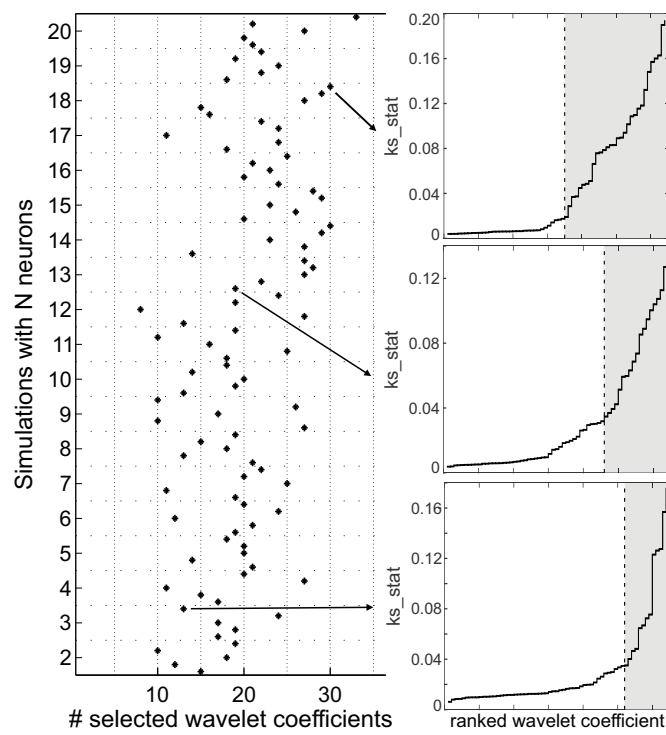
Figure 7: Performance of the proposed algorithm with the simulated tetrode dataset. A. A diagram of the simulated tetrode is shown on the top; each electrode had a 20 microns

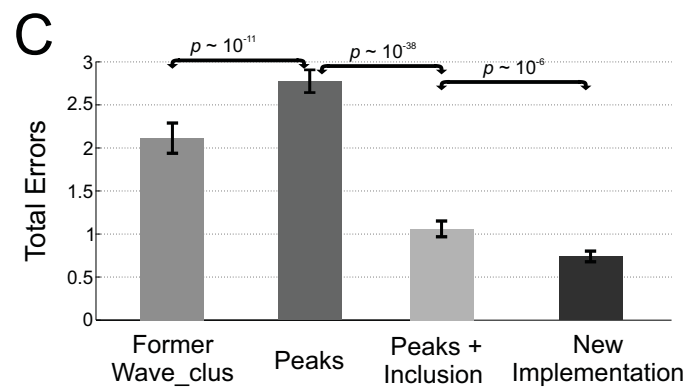
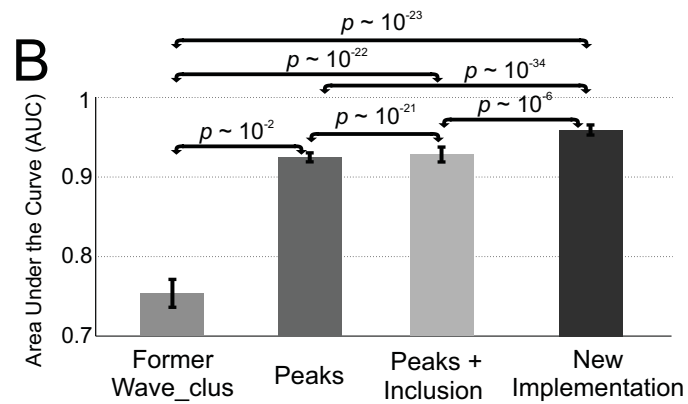
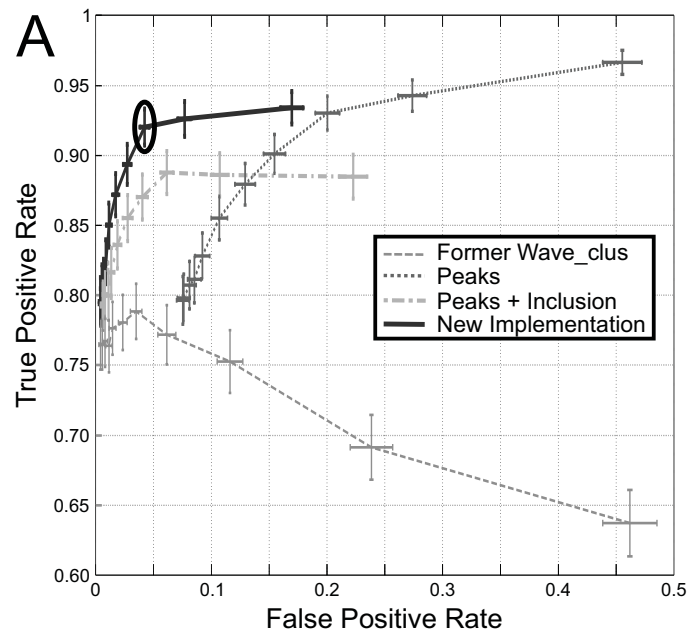
diameter, with a 40 microns separation between them. On the bottom, the coefficient selection for the simulation with 11 single units is shown (same conventions as in Figure 2). In this case, a total of 72 coefficients were selected. B. Results from the exemplar simulation with 11 single units are shown, as they were all properly isolated by the new implementation of *Wave_clus*. Spikes for each class on each recording channel are overlaid (total number of spikes per cluster shown in brackets), with the thick black line representing the mean waveform and the lighter ones associated with the SEM. Ticks on the x axis are placed every 20 samples, whereas the separation in the y axis is 50 uV for all channels. C. Mean and SEM of misses for each algorithm. Paired sign tests were used to evaluate the statistical difference across algorithms. D. Same as C but with false positives. The proposed algorithm with a variable number of wavelet coefficients was significantly better than all other implementations. E. Total number of errors, i.e. misses + false positives, for the different algorithms. The new *Wave_clus* showed significantly less errors than all other algorithms, with an average of at least half the number of the second-best.

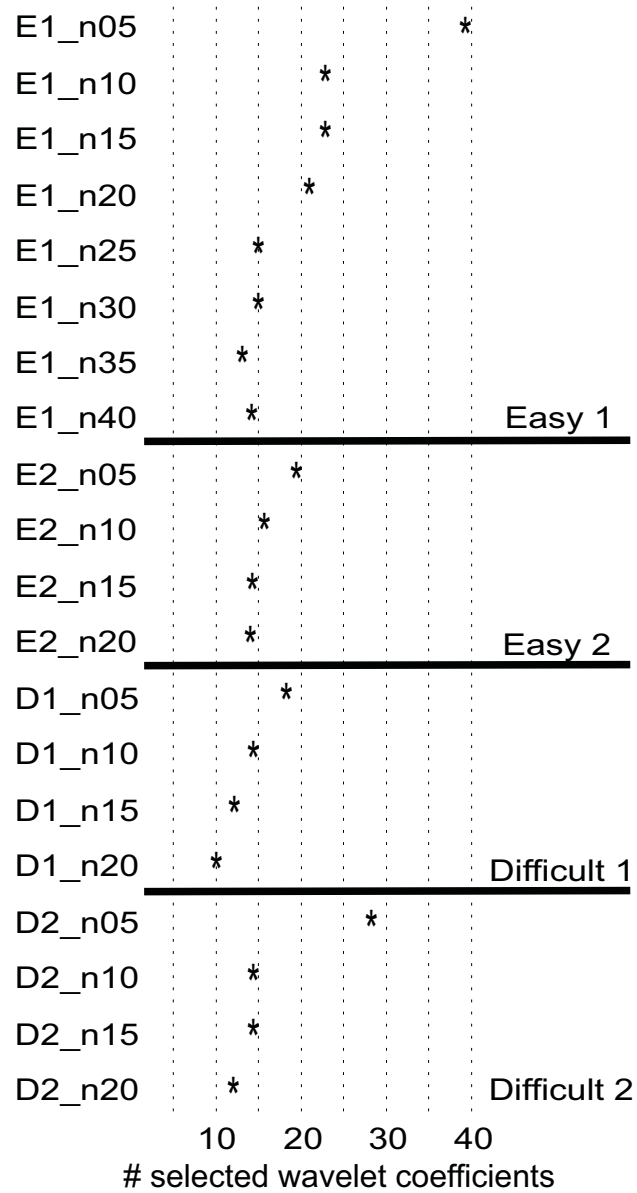
TABLE CAPTION

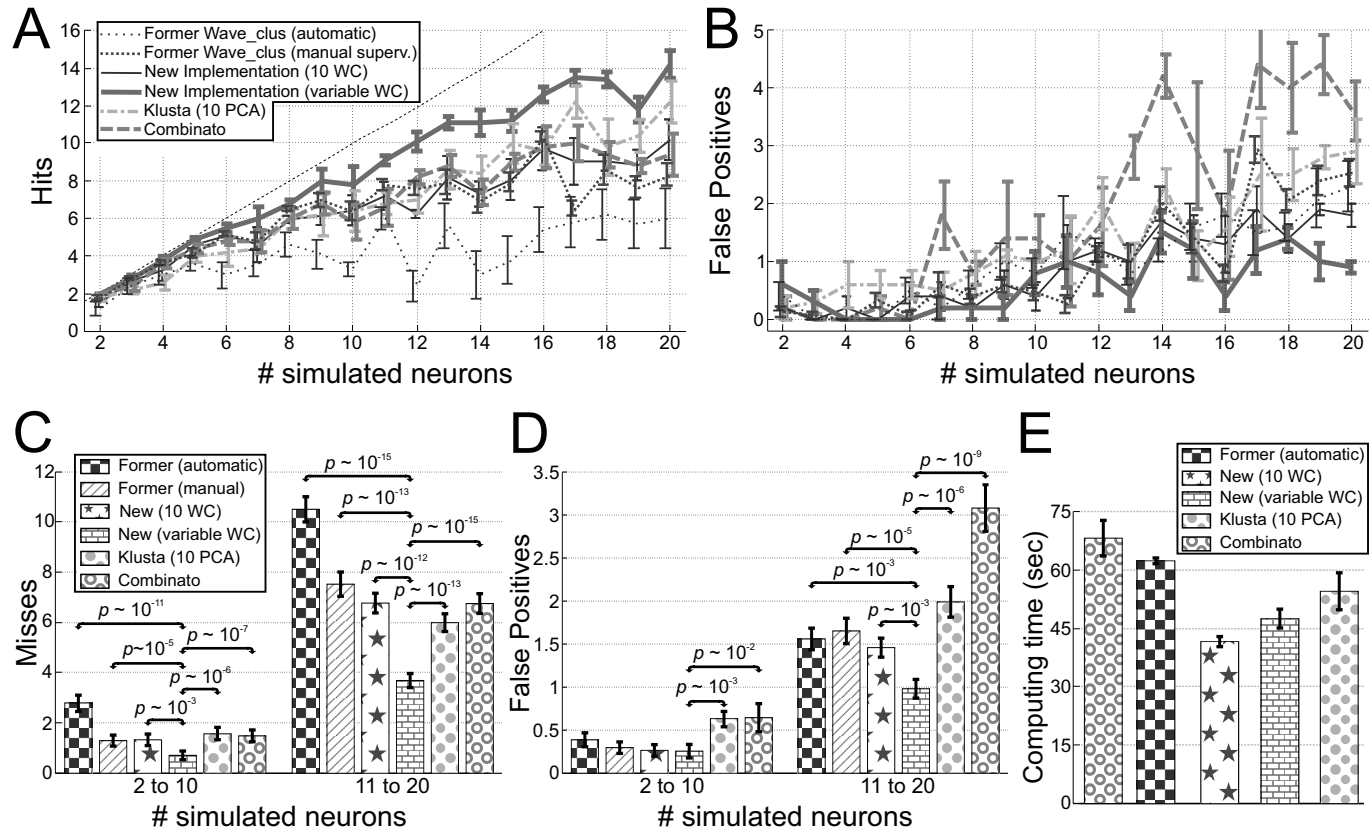
Table 1: Performance of the proposed implementation with the simulated dataset *Sim1*. Each row denotes a simulation of 3 different neurons with different levels of background noise. When the waveforms were different enough they are labelled as “Easy”, whereas when they were similar they were labelled as “Difficult”. The peak amplitude of each spike was normalized to 1 and the standard deviation of the background noise was varied from 0.05 to 0.2 (or 0.4 for “Easy1”). The table shows the number of hits for each algorithm implementation. In the former *Wave_clus* implementation, the performance with 64 wavelet coefficients (WC) was significantly worse than with 10 WC (paired sign test, $p = 0.008$). The same behavior was also seen with the new implementation with a fixed number of WC (10 WC vs. 64 WC, paired sign test, $p = 0.016$), but no difference was seen between the new implementation with 10 or variable number of WC.

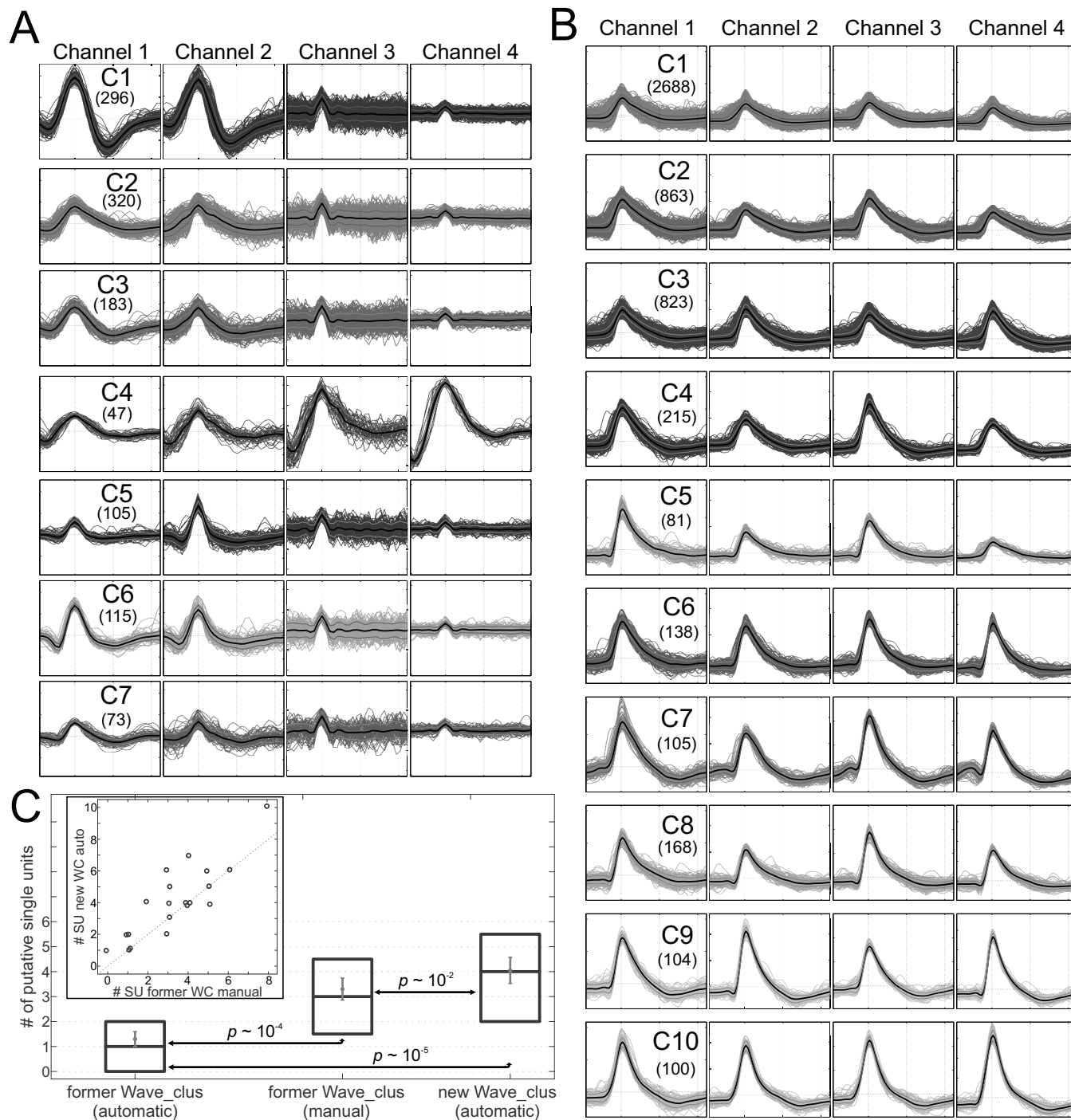


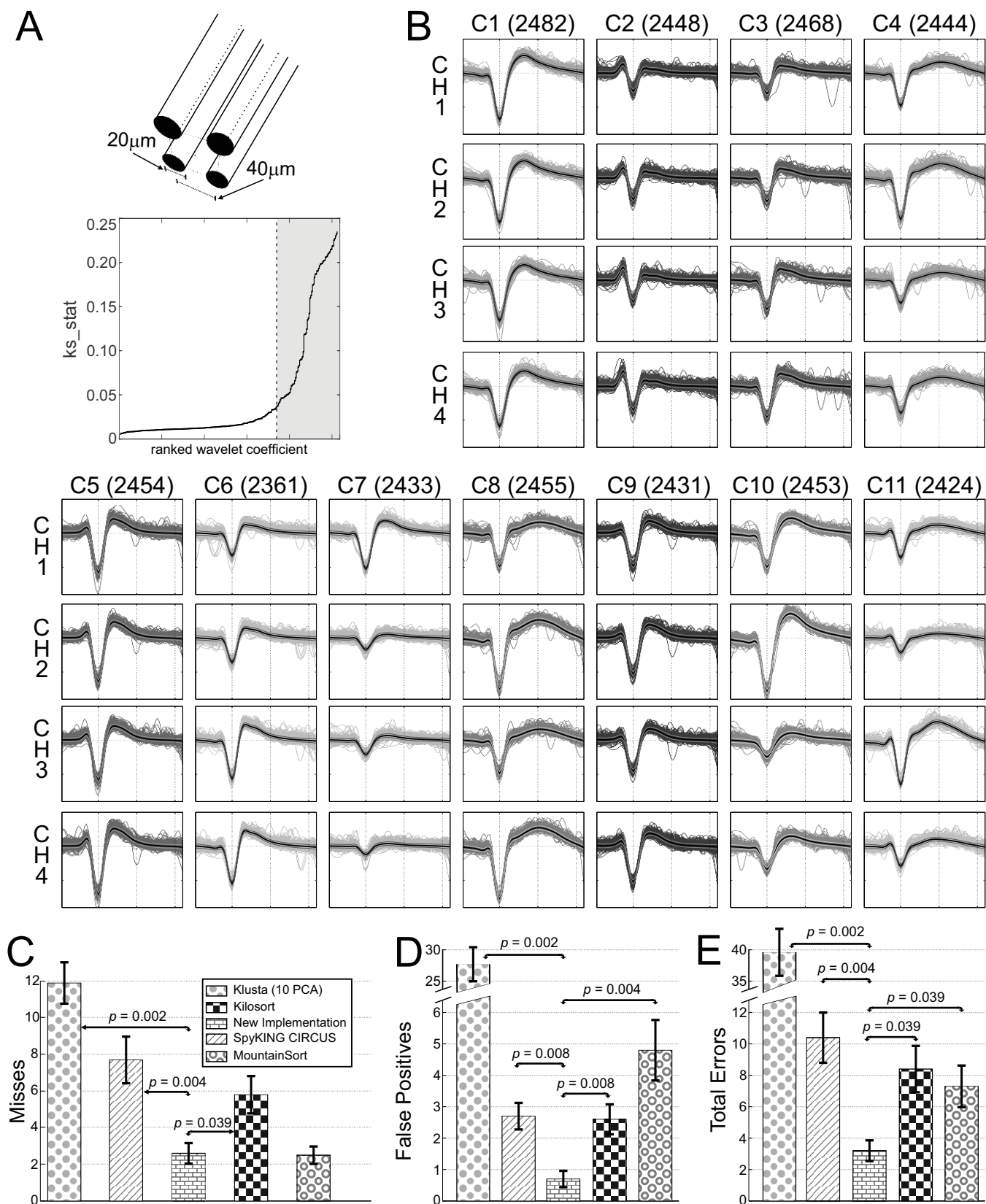












Simulation	Former Wave_clus		New Implementation		
	# coefficients		# coefficients		
	10	64	10	64	variable
Easy1_noise005	3	3	3	3	3
Easy1_noise010	3	3	3	3	3
Easy1_noise015	3	3	3	3	3
Easy1_noise020	3	3	3	3	3
Easy1_noise025	3	3	3	3	3
Easy1_noise030	3	3	3	3	3
Easy1_noise035	3	2	3	3	3
Easy1_noise040	3	1	3	1	3
Easy2_noise005	3	3	3	3	3
Easy2_noise010	3	3	3	3	3
Easy2_noise015	3	3	3	3	3
Easy2_noise020	3	1	3	1	3
Diffi1_noise005	3	3	3	3	3
Diffi1_noise010	3	3	3	3	3
Diffi1_noise015	3	2	3	0	3
Diffi1_noise020	3	0	3	0	3
Diffi2_noise005	3	3	3	3	3
Diffi2_noise010	3	2	3	2	3
Diffi2_noise015	3	2	3	1	3
Diffi2_noise020	2	0	2	0	2