# A comprehensive review of spike sorting algorithms in neuroscience

Wentao Quan[1], Youguo Hao[2], Xudong Guo[1], Peng Wang[1], Yukai Zhong[3]

[1]School of Health Science and Engineering, University of Shanghai for Science and Technology, Shanghai 200093, China. [2]Putuo District People's Hospital, Shanghai 200060, China. [3]Yangpu District Kongjiang Hospital, Shanghai 200082, China.

Corresponding author: Youguo Hao.

## Highlights

• The detailed steps of spike sorting algorithm and the different algorithms used in each step are summarized.
• The advantages and disadvantages of each step of spike sorting algorithm are compared.
• The detailed application of deep learning technology in spike sorting is introduced.

## Abstract

Spike sorting plays a pivotal role in neuroscience, serving as a crucial step of separating electrical signals recorded from multiple neurons to further analyze neuronal interactions. This process involves separating electrical signals that originate from multiple neurons, recorded through devices like electrode arrays. This is a very important link in the field of brain-computer interfaces. The objective of spike sorting algorithm (SSA) is to distinguish the behavior of one or more neurons from background noise using the waveforms captured by brain-embedded electrodes. This article starts from the steps of the conventional SSA and divides the SSA into three steps: spike detection, spike feature extraction, and spike clustering. It outlines prevalent algorithms for each phase before delving into two emerging technologies: template matching and deep learning-based methods. The discussion on deep learning is further subdivided into three approaches: end-to-end solution, deep learning for spike sorting steps, and spiking neural networks-based solutions. Finally, it elaborates future challenges and development trends of SSAs.

## Introduction

The control signals currently used in brain-computer interfaces (BCI) mainly include the following five types: action potentials (spike), local field potentials, electrocorticography (ECoG), epidural field potentials, and electroencephalography [1]. The sites where each signal is acquired in the brain are shown in **Figure 1**. **Figure 1** was adapted from [1]. Spike is a key signal reflecting the behavioral activity of the brain, and the study of neuronal information storage and encoding via spike contributes to better development of analog brain computers [2]. There is a close relationship between local field potentials and spike, and their good stability and frequency content provide key advantages for long-term, low-power BCI [3]. ECoG, capturing detailed aspects of actual and imagined body movements, facilitates the development of robust and long-term BCI systems [4]. Compared to ECoG, epidural field potentials are acquired in the epidural space, which greatly reduces the risk of surgical complications, making them suitable for long-term implantation of BCI [5]. though non-invasive and low-cost, offers limited brain activity information, restricting its use to simple communication and
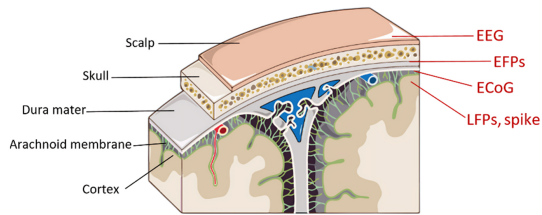
Address correspondence to: **Youguo Hao**, Putuo District People's Hospital, No.1291 Jiangning Road, Putuo, Shanghai 200060, China. Email: youguohao6@163.com.

**Figure 1. The locations in the brain where various signals occur.** EEG, electroencephlography; EFPs, eqidural field potentials; ECoG, electrocorticography; LFPs, local field potentials. This Figure was adapted from [1].
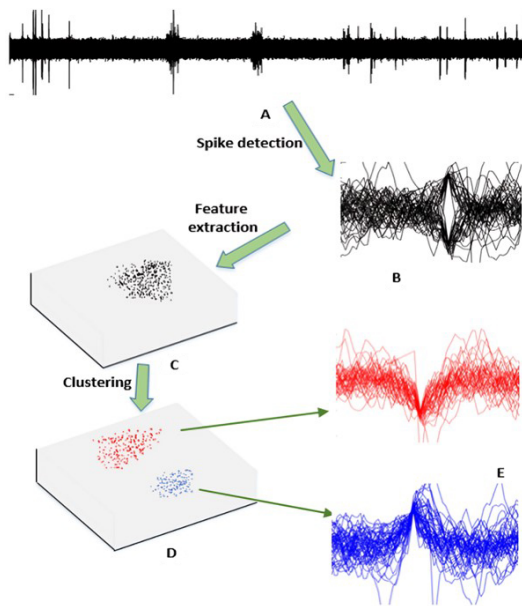


**Figure 2. Basic steps for spike sorting.** (A) A piece of filtered neuronal signal; (B) Spikes are detected, usually using an amplitude threshold of signal; (C) Spike features are extracted to achieve a dimensionality reduction; (D) Clustering operation based on extracted features; (E) The clustering algorithm classifies the waveforms and associates each cluster with a unit. This figure was adapted from [17].

control tasks [6]. Other signals have also been used in BCI, including near-infrared spectroscopy of brain blood flow and functional magnetic resonance imaging. The former's poor spatial resolution and the latter's poor temporal resolution, coupled with high costs, rendering them impractical for widespread application in BCI technologies.

In general, there is a correlation between the level of invasiveness of neural recordings and the quality of signal obtained as well as spatial, spectral, and temporal resolution. The higher the level of invasiveness, the more brain information can be obtained [7]. Spikes, directly emitted by the cells at the implantation site with a frequency range generally between 300Hz and 6000Hz, can reflect the interaction

between high-level brain functions and single neurons [8]. Typically, spikes are captured by inserting high-impedance electrodes with a diameter of several microns into the brain, thereby recording neuronal activity near the electrodes in the range of 50-100 microns [9]. In recent years, BCI using spike has attracted increasing interest from scientific community, with the goal of mapping, assisting, enhancing, and repairing cognitive or sensorimotor functions in human or animal brains. Today, spikes are used in the field of BCI to detect pain signals, control robotic arms or wheelchairs, and treat spinal cord injuries [10-12]. However, due to the large amount of spike data and low signal-to-noise ratio (SNR), it has not been widely used clinically.

Since the 1970s, the accurate extraction, detection, and the classification of spikes has been a hot topic in the neuroscience field, known as spike sorting [13]. This review provides a general overview of spike sorting, introduces corresponding algorithms for different steps, and summarizes and compares the current mainstream algorithms. Finally, current study limitations are summarized, and future developments are projected.

**Conventional spike sorting algorithm (SSA)**

Spike sorting is an algorithm to detect individual spikes from extracellular neural recordings and classify them based on their shape. The algorithm classifies the detected spikes as originating neurons [14]. The SSA works by dividing different neurons into different clusters, according to their different electrical proximity to the recording electrode and different shapes of dendrites [15]. In theory, a single neuron is represented by a single cluster; however, due to interactions among different neurons, the obtained spike information often overlaps, resulting in a large number of false positives in classification across adjacent channels, thereby complicating accurate categorization [16]. As shown in **Figure 2**, traditional spike sorting has three main steps. **Figure 2** was adapted from [17]. The first step, spike detection, involves identifying spikes at the corresponding time points based on the waveform of each spike. The second step is to extract features corresponding to different spikes and convert them into another feature domain. The final step, spike clustering, is to classify the spikes into different clusters [13].

*Spike detection*

Spike detection is the initial phase of SSA, ded-

icated to isolating individual spikes from continuous neural recordings, setting the groundwork for the subsequent stages of the sorting process. Generally, the duration of a spike event is usually 1-3 ms, and the sampling time of each spike event lasts for 30-90 sampling points to prevent spikes from being repeatedly detected. There are currently three mainstream methods of spike detection, including threshold detection method, nonlinear energy operator method, and wavelet transform product method. Each of these methods has its own advantages and disadvantages [17, 18].

*Amplitude thresholding*

The most straightforward algorithm is to apply an amplitude threshold to the spike signal. Since the filtered signal is easily visualized on top of background noise activity, this method is resource-intensive and computationally fast. However, if the applied threshold is too small, noise fluctuations will cause false positive events, and some spikes will be missed. Therefore, setting an adaptive threshold is more suitable.

The most reasonable way is finding a multiple $k$ of the noise $\sigma_n$ as the threshold, that is $Thr=k\sigma_n$. The general value of $k$ is between 3 and 5. $\sigma_n$ can also be other values. A threshold estimation formula based on the signal median is given [19],

$$Thr = 4 \times median(\frac{|x|}{0.6745}) \quad (1)$$

where 0.6745 comes from the inverse function of the standard normal distribution cumulative function, which has a value of 0.75. This method processes the background noise according to the standard Gaussian distribution. Even though the noise distribution in real world may deviate from the Gaussian distribution, it has been confirmed that this median-based estimation is more accurate than the standard deviation estimation [20]. However, Equation (1) is not suitable for all situations.

Another study proposed an improved formula [21],

$$Thr = \frac{median\{\frac{|x|}{0.6745}\}\sqrt{2\log_2(n)}}{1.6} \quad (2)$$

where $n$ is the length of the detected data. The addition of $n$ into Equation (2) reduces the mixed noise interference introduced by the col-

lection environment, making it more suitable for the adaptive threshold algorithm [21]. However, Equation (2) was later proved to be unsuitable for detection of high noise levels [22]. In this context, Equation (3) was proposed, as follows:

$$Thr = median(\frac{|x|}{0.6745})\sqrt{\frac{\log_2(n)}{n}(\frac{\ln n}{\ln 2})^3 q} \quad (3)$$

Equation (3) reduces the interference caused by noise fluctuations, making the set threshold more robust. It has been experimentally confirmed that under higher noise levels, the detection accuracy of Equation (3) is significantly higher than Equation (1) and (2) [22].

*Nonlinear energy operator (NEO)*

The NEO is a more powerful detection method, leveraging both frequency and amplitude information. This capability renders NEO more effective than the threshold method in low SNR situations [23]. For the continuous sampling signal $x(n)$, its discrete time NEO $\varphi$ is shown in Equation (4). In addition, some studies suggest using smoothing windows to convolve NEO time series, which can reduce the probability of false positives.

$$\varphi[x(n)] = x^2(n) - x(n+1) \times x(n-1) \quad (4)$$

$$Thr_{NEO} = \frac{8}{N} \sum_{n=1}^{N} \varphi[x(n)] \quad (5)$$

Spike is the component with the highest energy in the entire collected signal. Compared with the noise signal, the energy of the spike signal increases instantaneously. This method is ideal for on-chip implementation due to its good transient and computational properties [23]. Malik et al. provided a method for NEO threshold judgment, as shown in Equation (5), whereby a value surpassing this threshold is classified as a spike [24].

*Wavelet transform product*

Wavelet transform has also been proposed for spike detection due to its ability to integrate information from both the time and frequency domains [25]. Compared with the first two methods, the wavelet transform method does not regard the background noise as a standard Gaussian distribution [26]. The principle of this method mainly exploits the fact that the wavelet function is a long "peak" waveform. When sliding wavelet decomposition is used, its essence
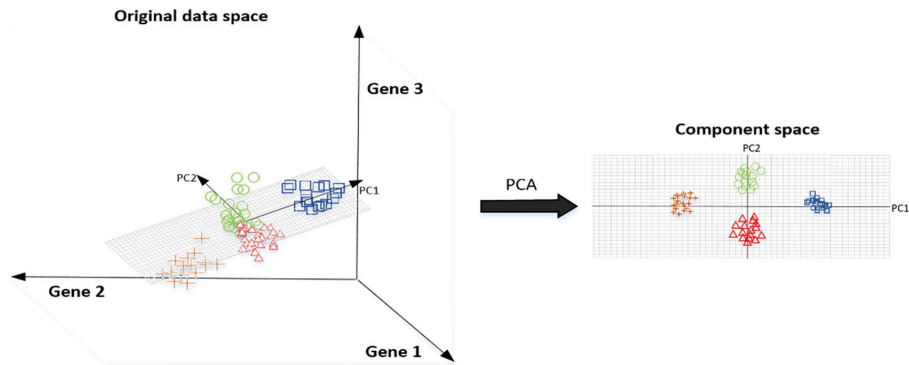
**Figure 3. The process of PCA dimensionality reduction.** PCA, principal component analysis. This Figure was adapted from [42].

is a "template matching" process to evaluate the similarity between the wavelet function and spike. The mathematical definition of wavelet transform is shown in Equation (6), where $\psi$ is the wavelet function, $\tau$ is the translation scale, and α is the scale factor.

$$W(\alpha, \tau) = \int_{-\infty}^{\infty} x(t) \frac{1}{\alpha^{\frac{1}{2}}} \psi(\frac{t - \tau}{\alpha}) dt \quad (6)$$

$$P(n) = \prod_{j=j_{max}-2}^{j_{max}} | w(2^j, n) | \quad (7)$$

It identifies the signal segment as a spike when there is a relatively high similarity between the signal segment and the wavelet function. There are many choices for wavelet functions in spike detection, such as Haar wavelet, db4 wavelet, and Biorthogonal wavelet. In applications, discrete wavelet transform is generally used, where the value of α is $2^j(j=1,2,3...)$. First, the wavelet coefficient $w$ ($\alpha$, $\tau$) corresponding to each $j$ value for each time point is calculated, and all time points until the set maximum $j$ value are summed up, termed as $j_{max}$. Next, the value of $P$ ($n$) is calculated, which is essentially the product of the wavelet coefficients across three consecutive scales, as shown in Equation (7). Similar to the method described above, using a smoothing window to convolve the $P$ ($n$) sequence can reduce spurious spike due to cross terms and background noise [27, 28].

Among the above three spike detection technologies, since the first two have relatively low computing resource requirements for hardware implementation, they are often favored for online implementation. When the SNR is high and computing resources are limited, the threshold detection method is more suitable. Although the NEO has a simple structure, it necessitates

data storage and introduces delays. The first two methods are prone to false positives and false negatives as the noise level changes. Without considering the amount of computing resources, the wavelet transform method can yield the best results [29, 30].

*Spike feature extraction*

Feature extraction can be performed from three aspects: time domain features, transformation domain features, and dimensionality reduction features. For the time domain features, actual operation can be performed based on the geometric features and derivative features of the spike [31, 32]. For transformation domain features, techniques such as Fourier transform, Hilbert transform, and wavelet transform are employed [33-35]. For dimensionality reduction features, principal component analysis (PCA), linear discriminant analysis (LDA), and Laplacian eigenmaps are utilized [36-38]. After feature extraction, 2 to 3 key features are generally selected and used as inputs in the clustering algorithm to ensure that the next step of the clustering algorithm can be effectively implemented [39, 40]. Several mainstream spike feature extraction algorithms are introduced below.

*PCA*

PCA is the most commonly used method in current feature extraction. Its essence is a dimensionality reduction method that projects a data set with multiple relevant features into a coordinate system with fewer relevant features. All detected spike time samples are used as input *x(n)*, and the covariance matrix is eigenvalue decomposed, where the eigenvectors represent the projection direction of the original data. The algorithm mainly includes three steps: (1) creating the covariance matrix of the data set;

(2) calculating all eigenvalues of the covariance matrix and retaining the first k eigenvalues in descending order of variance; (3) transforming the data points through *k* feature vectors [41].

**Figure 3** shows the dimensionality reduction on the original data through PCA. **Figure 3** was adapted from [42]. The disadvantage of this method is its tendency to lose information as the complexity of the collected spike data increases. It is stated that when a higher **k** value is selected, the variance of the higher component may be affected by the background noise, so this method is not suitable for feature extraction where the spike information is rich [13].

*Features of wavelet transform*

Wavelet transform is an integral transform that decomposes a specific signal into the sum of different wavelet signals. The relevant equations of wavelet transform are presented in Section 2.1.3. Through the low-frequency and high-frequency wavelet coefficients, it can accurately depict the spike shape [43]. Generally, wavelet transform convolves the original signal with wavelet functions of different parameters derived from the mother wavelet, which can quantify the details of the signal at different resolutions [35]. Haar wavelet and Daubechies wavelet are the two most commonly used mother wavelets in analyzing neurophysiological signals due to their orthogonality. They are able to identify the characteristics of a spike through a few wavelet coefficients without making a priori assumptions about the shape of the spike [19, 44, 45].

The features that can be extracted from the time spectrum mainly include extreme value features and energy distribution features. The former is the most prevalent feature, involve extracting the peak amplitude across each scale and correlating them with the corresponding scale parameters and translation parameters to construct a comprehensive feature vector. The latter requires calculating the energy distribution of each scale space from the energy spectrum, and then identifying the corresponding scale parameters and translation parameters to form the feature vector [46].

*LDA*

LDA is a linear projection technology that seeks a suitable projection direction within a low-dimensional space, so as to maximize the distance between different categories and minimize the distance between the same categories [47]. LDA first calculates the covariance matrix

of the data and performs eigenvalue decomposition of the covariance matrix. Then, a projection matrix is formed based on the obtained eigenvectors and eigenvalues to map the original data into a new low-dimensional space [48]. LDA needs to clearly minimize the inter-class variance and maximize the inter-class variance, which are shown in Equation (8) and Equation (9), respectively. Here, $x_i$ is the *ith* data point in the *kth* cluster *ck*, *μk* represents the average of data points in the *kth* cluster, *μ* represents the average of all data points, and *n* is the total number of data.

$$S_W = \sum_{k=1}^{K} \sum_{x_i \in c_k}^{n_k} (x_i - \mu_k)(x_i - \mu_k)^T \qquad (8)$$

$$S_b = \frac{\sum_{k=1}^{K} n_k (\mu_k - \mu)(\mu_k - \mu)^T}{n} \qquad (9)$$

LDA has demonstrated its effectiveness in separating spike features, especially when the SNR is low [48].

Unlike PCA, which simplifies high-dimensional data into two or three dimensions without significant information loss and facilitates visualization and classification identification, LDA optimizes dimensionality reduction to avoid overfitting. Wavelet transform can extract detailed features in the spike without too much information loss. However, PCA does not account for the categorical relationships between data points and relies solely on variance for determining principal components, leading to potential information loss. The classification effectiveness of LDA can easily be compromised when there is a great overlap between the data distribution of different categories. Wavelet transform is easily affected by noise, thus selecting suitable parameters for the wavelet basis function is essential [49-51].

*Spike clustering*

Spike clustering refers to organizing points in the feature space into clusters. Each cluster is associated with a different neuron, separated by defined boundaries between different clusters [52]. Optimal clustering techniques should minimize user intervention and achieve efficient and accurate results. It is generally believed that the changes within the cluster are caused by the noise superimposed on the real spike waveform [53]. The background noise obeys the Gaussian distribution. Therefore, most clustering algorithms are built based
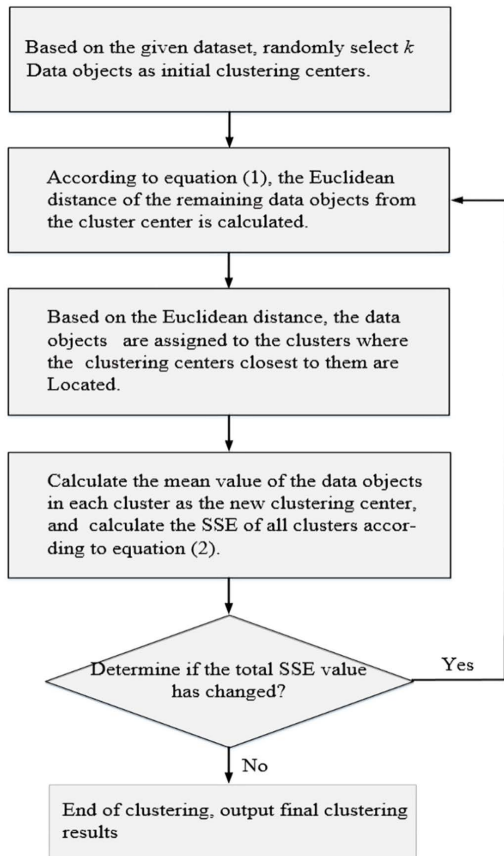
Figure 4. K-means algorithm flow chart. SSE, sum of squared error. This figure was adapted from [53].

on the Gaussian model. Most SSAs perform clustering by fitting a Gaussian mixture model, modeling the feature density curve as a sum of Gaussians, or fitting a mixture of t distributions. Several clustering methods commonly used in recent years are introduced.

*K-means*

K-means clustering is one of the prevalent unsupervised learning techniques. It primarily uses the Euclidean distance $d(x,c_i)$ as an indicator to measure the similarity between data objects. Its calculation formula is shown in Equation (10), where $x$ is the data object, $c_i$ is the *ith* cluster center, and $m$ is the dimension of the data. The similarity is inversely proportional to the distance between data objects.

$$d(x,C_i) = \sqrt{\sum_{j=1}^{m}(x_j - C_{ij})^2} \qquad (10)$$

$$SSE = \sum_{i=1}^{k}\sum_{x \in C_i}|d(x,C_i)|^2 \qquad (11)$$

The algorithm needs to specify the initial number of clusters k and k initial cluster centers in advance. Then, based on the similarity between the data object and the cluster center, the location of the cluster center is continuously updated, and the sum of squared error of the cluster is reduced, as shown in Equation (11). When the sum of squared error no longer changes or the objective function converges, clustering ends and the final result is obtained [54, 55]. The K-means algorithm is a continuously iterative process, and its algorithm flow is shown in **Figure 4. Figure 4** was adapted from [53].

*Mean shift*

Mean shift is a non-parametric clustering method first proposed by Fukunaga and Hostetler [56]. It treats data samples as empirical representations of probability density functions in a *d*-dimensional feature space. Dense regions in the *d*-dimensional feature space is used to represent local maxima of the underlying distribution [57]. The algorithm first estimates the kernel density function, as shown in Equation (12), where *k* represents the symmetrical kernel function, which usually is a Gaussian function.

$$f(x) = \sum_{i} k\left(\frac{\|x-x_i\|^2}{h^2}\right) \qquad (12)$$

After each data point is assigned to the corresponding *f(x)*, the algorithm moves the data point to the KDF mapping according to the gradient function in Equation (13) [58].

$$\nabla f(x) = \frac{2c_{k,d}}{nh^{d+2}}\left[\sum_{i=1}^{n}k'(\|\frac{x-x_i}{h}\|)\right]\left[\frac{\sum_{i=1}^{n}x_i k'(\|\frac{x-x_i}{h}\|^2)}{\sum_{i=1}^{n}k'(\|\frac{x-x_i}{h}\|^2)} - x\right] \qquad (13)$$

The focus of the mean shift algorithm is primarily to improve calculation efficiency and robustness [15]. In recent research, Yang et al. proposed a method that combines the mean shift algorithm with the optical flow method, which reduces the cumulative error and thereby improves computational efficiency [59].

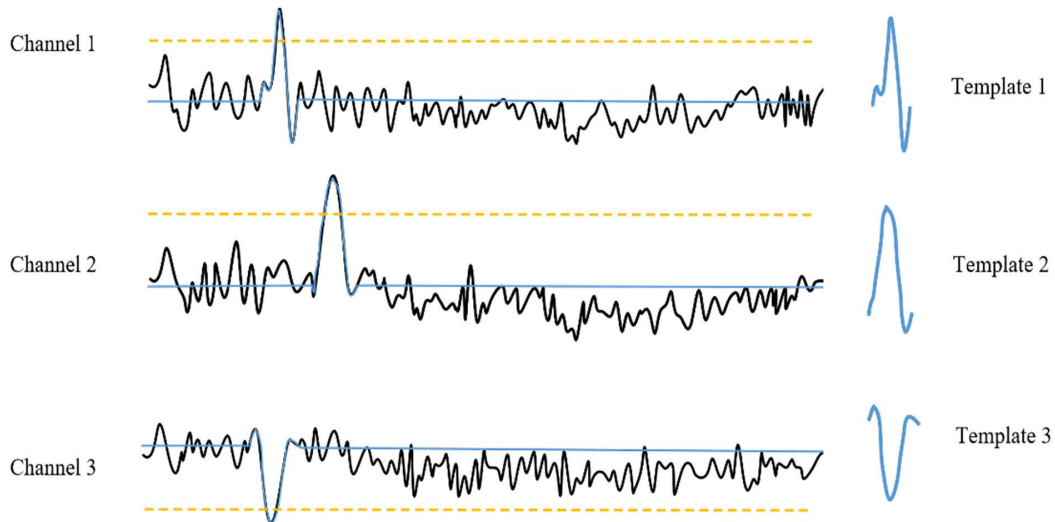*Super-paramagnetic clustering (SPC)*

**Figure 5. Templates are obtained using threshold detection method.** This figure was adapted from [15].

SPC is an unsupervised classification technique based on clustering self-tenancy of data density and interaction energy. Drawing inspiration from the interactions observed among superparamagnetic particles under varying temperatures, SPC distinguishes itself from conventional algorithms by excluding the use of clustering distribution functions. First, it uses Equation (14) to mathematically simulate the interaction between $x_i$ and $x_j$ at different points, where a is the average distance between nearest neighbor points. Then, it assigns a "Potts spin" state variable s from to q to each point $x_i$, where q is usually chosen to be 20. Finally, it uses the Wolff algorithm to perform N Monte Carlo iterations at different temperatures [19, 60]. The state value $p_{ij}$ is updated according to Equation (15), where T represents the temperature value.

$$J_{ij} = \frac{1}{k}\exp[-\frac{|x_i - x_j|^2}{2a^2}] \qquad (14)$$

$$P_{ij} = 1 - \exp(-\frac{J_{ij}}{T}\delta_{s_i,s_j}), \quad \delta_{s_i,s_j} = \begin{cases} 1 & if\ s_i = s_j \\ 0 & otherwise \end{cases} \quad (15)$$

$$G_{ij} = \frac{(q-1)c_{ij}+1}{q} \qquad (16)$$

s will be iteratively updated according to the value of $p_{ij}$ until it becomes stable. Equation (16) introduces the spin-spin correlation function, where $c_{ij}$ is a number belonging to the interval

[0, 1]. An appropriate threshold ε can be selected. When $G_{ij} > ε$ , different data points $x_i$ and $x_j$ belong to the same cluster.

*Fuzzy C-means (FCM)*

FCM combines the K-means algorithm with fuzzy logic, which has good processing effects on data with high complexity. FCM is an iterative algorithm whose execution steps are similar to K-means. First, the number of clusters k needs to be determined. Its value can be fixed or adaptive. Next, it randomly initializes the value $μ_{ij}$ of the cluster member of each point $x_i$ relative to the cluster center $C_j$, and then calculates the cluster center according to Equation (17) [15, 61].

$$C_j \quad \frac{\sum_{i\ 1}^{D} μ_{ij}^{m} x_i}{\sum_{i\ 1}^{D} μ_{ij}^{m}} \qquad (17)$$

$$μ_{ij} = \frac{1}{\sum_{k=1}^{N}(\frac{||x_i - c_j||}{||x_i - c_k||})^{\frac{2}{m-1}}} \qquad (18)$$

where D is the number of data points, and N is the number of clusters, and m is the fuzzy parameter. The values of cluster members are then calculated and updated using Equation (18), where m is a parameter greater than 1. When m approaches 1, the ambiguity decreases, and the algorithm result approaches the

K-means algorithm.

Of the four algorithms introduced above, K-means has the simplest structure and is easy to implement, excelling in efficiency with large datasets. Mean shift has good adaptability to data with linear structure. SPC stands out for good adaptability to automatically determine the number of clusters. FCM is distinguished by its good robustness in data processing and has good adaptability to situations where data points belong to multiple clusters. However, each algorithm has its limitations: K-means is sensitive to noise and outliers, and thus does not have good robustness; Mean shift has high computational complexity on large-scale data sets; SPC is complex and can be affected by temperature, and the convergence speed of FCM is relatively slow [62-65].

**Novel SSAs**

Novel SSAs do not follow the three steps of traditional algorithms. For example, using neural system simulation and deep shrinkage automatic encoding methods to perform spike sorting [66, 67]. Current mainstream algorithms can perform two or more steps simultaneously, thereby improving algorithm efficiency. These algorithms are mainly divided into two categories, one is template matching and the other is deep learning.

*Template matching*

Template matching identifies the position of a specific template within a signal sequence, facilitating accurate and efficient locate specific spikes in spike sorting. It takes the normalized-template-matching method proposed as an example [68]. Firstly, conventional SSAs are used to extract the initial template, such as threshold detection, as shown in **Figure 5**. The figure was adapted from [15]. Subsequently, the template was used to move against the neural recording to assess waveform similarity that determined by calculating the relationship between the extracellular voltage signal and the average spike waveform of each candidate single unit, which here called cross-correlation. Let the average spike waveform of all channels be $\mu_i$, where $\mu_i=[\mu_{i,1},\mu_{i,2},...\mu_{i,N}]$, $\mu_{i,n}$ is at the average spike wave in $n$ channels. Let the $L=f_s k$, where $f_s$ is the sampling frequency and $k$ is the length of the custom time window, then $\mu_{i,n}$ is a $1 \times L$ vector. In the same way, the neural recording signal is defined as a vector $u_{i,n}$ that changes with time, thus obtaining $V(t)=[u(t)_1,u(t)_2,...u(t)_N]$. Based on the above, the cross-correlation equation can be obtained:

$$C_i(t) = V(t) \times \mu_i^T \qquad (19)$$

Based on cross-correlation, the two steps of spike detection and feature extraction can be accomplished, and then K-means and other methods can be employed to complete clustering and the entire SSA step. However, since neural recording signals are usually non-stationary, template matching methods require template recalibration over a period of time.

Recent advancements in template matching for peak sorting include the utilization of particle swarm optimization algorithm and the multi-template matching method to optimize the template [69, 70]. For the optimized template, experimental results demonstrated that it could significantly reduce the occurrence of false positives and false negatives. Another study combined templates with convolutional neural networks (CNNs), employing CNNs to train on real spikes and overlapping spikes, which allows for correct spike classification and minimizes the occurrence of repeated spikes [71].

*Deep learning*

Deep learning is a popular technology that has been used in various fields for several years. It involves the application of machine learning through deep artificial neural networks (ANN), which has more layers, parameters and neurons. The use of deep learning methods to solve the peak sorting problem has become a popular research direction, and there are three main strategies: (1) end-to-end solution, (2) deep learning for spike sorting steps, and (3) spiking neural network (SNN)-based solutions [72-74].

*End-to-end solution*

The end-to-end solution uses the filtered neural recording as the input of the deep neural network and directly outputs the spike sequences produced by different neurons. This approach eliminates tedious preprocessing steps by starting from raw data as input. For instance, researchers used CNN and long short-term memory to perform peak sorting in an end-to-end manner [75]. Here, CNN is used for spike detection, and long short-term memory is used to receive the output stream of CNN and process the corresponding information. Other researchers proposed a medium 1D-CNN model tailored for single-channel spike sorting [76]. Although the end-to-end solution is a simple and

efficient way, its internal parameters require a large amount of data for training, and this method is not suitable for complex spikes.

*Deep learning for spike sorting steps*

Section 2 explains spike sorting as a process comprised three independent steps. Some of these steps can be improved by using deep learning methods. The most common one is to use deep learning for spike detection and denoising. Extracellular nerve signal recording contains a variety of components originating from the signal recording equipment or the measured organism itself, such as the activity of nearby neuronal cells. These noises can affect the detection accuracy. Lecoq et al. used original noise samples to train a spatiotemporal nonlinear difference model, thereby achieving the effect of eliminating independent noise sources in the input signal [77]. Okreghe et al. proposed a deep spike detection technique using CNN in conjunction with K-means, exhibiting robust learning capabilities for high-dimensional vectors and effective removal of artifacts in the acquisition channel [78].

*SNN-based solutions*

SNN, a relatively new type of neural network, exhibits functions more akin to neural networks in living organisms when compared to other conventional neural networks. SNN uses a temporal encoding scheme to encode information into spike sequences. It has more powerful information representation capabilities than ANN, especially for complex temporal or spatiotemporal data. In hardware implementation, SNN can achieve fast and large-scale parallel information processing [79, 80]. In the application of spike sorting, Mukhopadhyay et al. proposed a spike sorting method that combines SNN and K-means to achieve training and classification of on-chip systems through a two-step shared training scheme [81]. Werner et al. proposed a method to implement spike sorting using a two-layer SNN. In this technology, there is no need for preliminary waveform detection. The first layer is used to filter the original signal, and the second layer is used for detection and feature extraction. Both layers of neurons follow the Leaky Integrate-and-Fire model [82]. The disadvantage of this technology is that the parameters of Leaky Integrate-and-Fire need to be manually adjusted to adapt to different shapes of waveforms. Later, Bernert et al. optimized it and confirmed the high efficiency of this technology in high-density electrode settings [83].

## Conclusion

This review provides a detailed introduction of various steps involved in SSA, covering the predominant algorithms associated with each step and their current research status. Additionally, a detailed description of the application of deep learning technology in the field of spike sorting is provided. With the development of BCI, there is an escalating demand for enhanced neural signal acquisition channels and signal quality, which also poses challenges to the development of SSAs.

The foremost challenge is to solve the problem of signal non-stationarity. Neural signals change over time and can be affected by multiple factors, leading to temporal variations in spike shapes, especially in systems that require manual intervention [84, 85]. The second challenge is overlapping spikes, which significantly reduce the performance of conventional clustering algorithms [86]. The most commonly accepted solutions nowadays are Bayesian statistical techniques and template matching introduced in section 3.1 and the ANN. However, the above methods are computationally intricate and demand substantial hardware capabilities. Therefore, future research should aim to simplify and enhance the efficiency of algorithms designed to tackle overlapping spikes. The third challenge is the determination of the training data set and evaluation scheme. Most current spike sorting systems are evaluated using synthetic datasets with real labels. The simplest method is to add modeled Gaussian noise to the spike waveform. Due to the interactions between neuronal cells, it is very challenging to form realistic multi-channel synthetic data.

## References

[1]    Slutzky MW, Flint RD. Physiological properties of brain-machine interface input signals. J Neurophysiol 2017;118(2):1329-1343.

[2]    Wang JH, Cui S. Neuronal Signal Encoding and Storage as Principles of Brain Function. Prog

Biochem Biophys 2016;43:367-373.

[3] Jackson A, Hall TM. Decoding Local Field Potentials for Neural Interfaces. IEEE Trans Neural Syst Rehabil Eng 2017;25(10):1705-1714.

[4] Schalk G, Leuthardt EC. Brain-computer interfaces using electrocorticographic signals. IEEE Rev Biomed Eng 2011;4:140-154.

[5] Branco MP, Geukes SH, Aarnoutse EJ, et al. Nine decades of electrocorticography: A comparison between epidural and subdural recordings. Eur J Neurosci 2023;57(8):1260-1288.

[6] Wang Y, Nakanishi M, Zhang D. EEG-Based Brain-Computer Interfaces. Adv Exp Med Biol 2019;1101:41-65.

[7] Markowitz DA, Wong YT, Gray CM, et al. Optimizing the Decoding of Movement Goals from Local Field Potentials in Macaque Cortex. J Neurosci 2011;31(50):18412-22.

[8] Franke F, Natora M, Boucsein C, et al. An online spike detection and spike classification algorithm capable of instantaneous resolution of overlapping spikes. J Comput Neurosci 2010;29(1-2):127-148.

[9] Quiroga RQ. What is the real shape of extracellular spikes? J Neurosci Methods 2009;177:194-198.

[10] Hu SL, Zhang QS, Wang J, et al. A Real-Time Rodent Neural Interface for Deciphering Acute Pain Signals from Neuronal Ensemble Spike Activity. 2017 51st Asilomar Conference on Signals, Systems, and Computers 2017;93-97.

[11] Odoherty JE, Merolla PA, Even CN. Brain-machine interface (BMI) with user interface (UI) aware controller. US patent US11630516. 2023/04/18.

[12] Oweiss K, Aghogolzadeh M. Multiscale intra-cortical neural interface system. US patent US2011307079. 2011/12/15.

[13] Lewicki MS. A review of methods for spike sorting: the detection and classification of neural action potentials. Network 1998;9:R53-R78.

[14] Lefebvre B, Yger P, Marre O. Recent progress in multi-electrode spike sorting methods. J Physiol Paris 2016;110:327-335.

[15] Zhang T, Azghadi MR, Lammie C, et al. Spike sorting algorithms and their efficient hardware implementation: a comprehensive survey. J Neural Eng 2023;20 (2).

[16] Marre O, Amodei D, Deshmukh N, et al. Mapping a Complete Neural Population in the Retina. J Neurosci 2012;32:14859-14873.

[17] Hennig MH, Hurwitz C, Sorbaro M. Scaling Spike Detection and Sorting for Next-Generation Electrophysiology. Adv Neurobiol 2019;22:171-184.

[18] Rey HG, Pedreira C, Quian Quiroga R. Past, present and future of spike sorting techniques. Brain Res Bull 2015;119:106-117

[19] Quiroga RQ, Nadasdy Z, Ben-Shaul Y. Unsupervised spike detection and sorting with wavelets and superparamagnetic clustering. Neural Comput 2004;16:1661-1687.

[20] Quiroga RQ. Spike sorting. Curr Biol 2012;22:R45-R46.

[21] Wang J, Guo TX, Lu YS, et al. Heuristic Threshold Detection Algorithm for Spike. Comput Eng Appl 2022;58:192-196.

[22] Zhao DC, Jiao SY, Chen H, et al. Heuristic adaptive threshold detection method for neuronal spikes. IET Signal Process 2023;e12214.

[23] Mukhopadhyay S, Ray GC. A new interpretation of nonlinear energy operator and its efficacy in spike detection. IEEE Trans Biomed Eng 1998;45:180-187.

[24] Malik MH, Saeed M, Kamboh AM. Automatic Threshold Optimization in Nonlinear Energy Operator Based Spike Detection. 2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC) 2016;774-777.

[25] Kim KH, Kim SJ. A wavelet-based method for action potential detection from extracellular neural signal recording with low signal-to-noise ratio. IEEE Trans Biomed Eng 2003;50:999-1011.

[26] Nenadic Z, Burdick JW. Spike detection using the continuous wavelet transform. IEEE Trans Biomed Eng 2005;52:74-87.

[27] Shen Q, Liu XY, Jiang DZ. Using modulus maximum pair of wavelet transform to detect spike wave of epileptic EEG. Proceedings of the 20th Annual International Conference of the IEEE Engineering in Medicine and Biology Society. Vol.20 Biomedical Engineering Towards the Year 2000 and Beyond (Cat. No.98CH36286). IEEE, 1998;1543-1545.

[28] Yuning Y, Kamboh A, Andrew JM. Adaptive threshold spike detection using stationary wavelet transform for neural recording implants. 2010 IEEE Biomedical Circuits and Systems Conference (BioCAS 2010) 2010;9-12.

[29] Paralikar KJ, Rao CR, Clement RS. New approaches to eliminating common-noise artifacts in recordings from intracortical microelectrode arrays: Inter-electrode correlation and virtual referencing. J Neurosci Methods 2009;181:27-35.

[30] Rakshit A, Okkesim S, Hofmann UG, et al. Assessing Ambiguity Of Spike Detection By Thresholding. 2023 11th International IEEE/EMBS Conference on Neural Engineering (NER) 2023;1-4.

[31] Zamani M, Demosthenous A. Feature

Extraction Using Extrema Sampling of Discrete Derivatives for Spike Sorting in Implantable Upper-Limb Neural Prostheses. IEEE Trans Neural Syst Rehabil Eng 2014;22:716-726.

[32] Kamboh AM, Mason AJ. Computationally Efficient Neural Feature Extraction for Spike Sorting in Implantable High-Density Recording Systems. IEEE Trans Neural Syst Rehabil Eng 2013;21:1-9.

[33] Yang CH, Yuan Y, Si JN. Robust spike classification based on frequency domain neural waveform features. J Neural Eng 2013;10:066015.

[34] Zhu JD, Lin CF, Chang SH, et al. Analysis of spike waves in epilepsy using Hilbert-Huang transform. J Med Syst 2015;39:170.

[35] Mallat SG. A theory for multiresolution signal decomposition: the wavelet representation. IEEE Trans Pattern Anal Mach Intell 1989;11:674-693.

[36] Zhao DC, Jiao SY, Li L, et al. Neuron spike classification method, involves obtaining original spike signal and preprocess obtained spike signal, and reducing dimensionality of spike data by principal component analysis to obtain eigenvalues and eigenvectors, and map spike points to feature space constructed by the eigenvectors. China patent CN115844422A. 2023/03/28.

[37] Chen W, Shen CP, Chiu MJ, et al. Epileptic EEG Visualization and Sonification Based on Linear Discriminate Analysis. 2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC) 2015;4466-4469.

[38] Liu F, Zhang WJ, Gu SC. Local linear Laplacian eigenmaps: A direct extension of LLE. Pattern Recogn Lett 2016;75:30-35.

[39] Wang T. Characteristic of Neural Signal Feature for Spike Sorting and Detection. J Phys Conf Ser 2021;2078:012042.

[40] Huang L, Gan L, Ling BW. A Unified Optimization Model of Feature Extraction and Clustering for Spike Sorting. IEEE Trans Neural Syst Rehabil Eng. 2021;29:750-759.

[41] Wheeler BC, Heetderks WJ. A comparison of techniques for classification of multiple neural signals. IEEE Trans Biomed Eng 1982;29:752-759.

[42] Visualization of PCA principal component analysis. Available at: https://zhuanlan.zhihu.com/p/28909807

[43] Lieb F, Stark HG, Thielemann C. A stationary wavelet transform and a time-frequency based spike detection algorithm for extracellular recorded data. J Neural Eng 2017;14:036013.

[44] Hulata E, Segev R, Ben-Jacob E. A method for spike sorting and detection based on wavelet packets and Shannon's mutual information. J Neurosci Methods 2002;117:1-12.

[45] Letelier JC, Weber PP. Spike sorting based on discrete wavelet transform coefficients. J Neurosci Methods 2000;101:93-106.

[46] Wang PSP, Yang JW. A Review of Wavelet-Based Edge Detection Methods. Int J Pattern Recognit Artif Intell 2012;26(07):1255011.

[47] Sifaou H, Kammoun A, Alouini MS. High-dimensional Linear Discriminant Analysis Classifier for Spiked Covariance Model. J Mach Learn Res 2020;21:4508–4531.

[48] Keshtkaran MR, Yang Z. Noise-robust unsupervised spike sorting based on discriminative subspace learning with outlier handling. J Neural Eng 2017 ;14:036003.

[49] Chah E, Hok V, Della-Chiesa A, et al. Automated spike sorting algorithm based on Laplacian eigenmaps and k-means clustering. J Neural Eng 2011;8:016006.

[50] Quiroga RQ, Panzeri S. Extracting information from neuronal populations: information theory and decoding approaches. Nat Rev Neurosci 2009;10:173-185.

[51] Rossant C, Kadir SN, Goodman DFM, et al. Spike sorting for large, dense electrode arrays. Nat Neurosci 2016;19:634-641.

[52] Carlson D, Carin L. Continuing progress of spike sorting in the era of big data. Curr Opin Neurobiol 2019;55:90-96.

[53] Pouzat C, Mazor O, Laurent G. Using noise signature to optimize spike-sorting and to assess neuronal classification quality. J Neurosci Methods 2002;122:43-57.

[54] Yang J, Zhao C. Survey on K-Means Clustering Algorithm. Comput Eng Appl 2019;55:7-14, 63.

[55] Jain AK. Data clustering: 50 years beyond K-means. Pattern Recogn Lett 2010;31:651-666.

[56] Fukunaga K, Hostetler LD. The estimation of the gradient of a density function, with applications in pattern recognition. IEEE T Inform Theory 1975;21:32-40.

[57] Nguyen T, Khosravi A, Hettiarachchi I, et al. Classification of Neural Action Potentials using Mean Shift Clustering. 2014 IEEE International Conference on Systems, Man, and Cybernetics (SMC) 2014;1247-1252.

[58] Comaniciu D, Meer P. Mean shift: A robust approach toward feature space analysis. IEEE Trans Pattern Anal Mach Intell 2002;24:603-619.

[59] Yang X, Li Q, He C. Improved Meanshift Tracking Algorithm Based on Optical flow. 2023 IEEE 6th Eurasian Conference on Educational Innovation (ECEI) 2023;159-162.

[60] Wang PK, Chen CH, Pun SH, et al. Parallel architecture to accelerate superparamagnetic clustering algorithm. Electron Lett 2020;56:701-704.

[61] Inan ZH, Kuntalp M. A study on fuzzy C-means

clustering-based systems in automatic spike detection. Comput Biol Med 2007;37:1160-1166.

[62] Hamed HNA, Saleh AY, Shamsuddin SMH, et al. Multi-objective K-means evolving spiking neural network model based on differential evolution. 2015 International Conference on Computing, Control, Networking, Electronics and Embedded Systems Engineering (ICCNEEE) 2015;379-383.

[63] Yang Z, Zhao Q, Liu WT, et al. Energy Based Evolving Mean Shift Algorithm for Neural Spike Classification. 2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society 2009;966-969.

[64] Ott T, Kern A, Stoop R. Faster spike sorting with belief propagation. 2006 International Symposium on Nonlinear Theory and its Applications 2006;1027-1030.

[65] Lin R, Lin J, Wang J, et al. Diagnosing transformer voiceprint fault based on fuzzy C-means clustering algorithm comprises using fuzzy C-means clustering algorithm to classify different voiceprint fault data, and using classification results to determine type of fault. China patent CN112149569A. 2020/12/29.

[66] Yu H, Qi Y, Pan G. NeuSort: an automatic adaptive spike sorting approach with neuromorphic models. J Neural Eng 2023;20.

[67] Radmanesh M, Rezaei AA, Jalili M, et al. Online spike sorting via deep contractive autoencoder. Neural Networks 2022;155:39-49.

[68] Laboy-Juárez KJ, Ahn S, Feldman DE. A normalized template matching method for improving spike detection in extracellular voltage recordings. Sci Rep 2019;9:12087.

[69] Wu DP, Shi HC, Jiang LR, et al. BECT Spike Detection Algorithm Based on Optimal Template Matching and Morphological Feature Selection. IEEE Trans Circuits Syst II Express Briefs 2022;69:2366-2370.

[70] Jiang TJ, Wu DP, Gao F, et al. Improved Spike Detection Algorithm Based on Multi-Template Matching and Feature Extraction. IEEE Trans Circuits Syst II Express Briefs 2022;69:249-253.

[71] Zhang L, Gao D, Wang M. Sorting Overlapping Spikes Based on Log-Mel Spectrogram and Convolutional Neural Networks. 2023 6th International Conference on Artificial Intelligence and Big Data (ICAIBD) 2023;482-485.

[72] Wu T, Zhao W, Keefer E, et al. Deep Compressive Autoencoder for Action Potential Compression in Large-Scale Neural Recording. J Neural Eng 2018;15:066019.

[73] Wouters J, Kloosterman F, Bertrand A. A data-driven spike sorting feature map for resolving spike overlap in the feature space. J Neural Eng 2021;18.

[74] Buccino AP, Garcia S, Yger P. Spike sorting: new trends and challenges of the era of high-density probes. Prog Biomed Eng 2022;4:022005.

[75] Rácz M, Liber C, Németh E, et al. Spike detection and sorting with deep learning. J Neural Eng 2020;17:016038.

[76] Li ZH, Wang YT, Zhang N, et al. An Accurate and Robust Method for Spike Sorting Based on Convolutional Neural Networks. Brain Sci 2020;10:835.

[77] Lecoq J, Oliver M, Siegle JH, et al. Removing independent noise in systems neuroscience data using DeepInterpolation. Nat Methods 2021;18:1401-1408.

[78] Okreghe CO, Zamani M, Demosthenous A. A Deep Neural Network-Based Spike Sorting With Improved Channel Selection and Artefact Removal. IEEE Access 2023;11:15131-15143.

[79] Wang XW, Lin XH, Dang XC. Supervised learning in spiking neural networks: A review of algorithms and evaluations. Neural Networks 2020;125:258-280.

[80] Farsa EZ, Ahmadi A, Maleki MA, et al. A Low-Cost High-Speed Neuromorphic Hardware Based on Spiking Neural Network. IEEE Trans Circuits Syst II Express Briefs 2019;66:1582-1586.

[81] Kumar Mukhopadhyay A, Chakrabarti I, Basu A, et al. Power efficient Spiking Neural Network Classifier based on memristive crossbar network for spike sorting application. arXiv:1802.09047.

[82] Werner T, Vianello E, Bichler O, et al. Spiking Neural Networks Based on OxRAM Synapses for Real-Time Unsupervised Spike Sorting. Front Neurosci 2016;10:474.

[83] Bernert M, Yvert B. An Attention-Based Spiking Neural Network for Unsupervised Spike-Sorting. Int J Neural Syst 2019;29:1850059.

[84] Fee MS, Mitra PP, Kleinfeld D. Variability of extracellular spike waveforms of cortical neurons. J Neurophysiol 1996;76:3823-3833.

[85] Snider RK, Bonds AB. Classification of non-stationary neural signals. J Neurosci Methods 1998;84:155-166.

[86] Harris KD, Henze DA, Csicsvari J, et al. Accuracy of tetrode spike separation as determined by simultaneous intracellular and extracellular measurements. J Neurophysiol 2000;84:401-414.