
Research Article: Theory/New Concepts | Novel Tools and Methods

How do spike collisions affect spike sorting performance?

<https://doi.org/10.1523/ENEURO.0105-22.2022>

Cite as: eNeuro 2022; 10.1523/ENEURO.0105-22.2022

Received: 11 March 2022

Revised: 15 June 2022

Accepted: 23 June 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.eneuro.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Copyright © 2022 Garcia et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1. Manuscript Title (50 word maximum)

How do spike collisions affect spike sorting performance?

2. Abbreviated Title (50 character maximum)

How do spike collisions affect spike sorting performance?

3. List all Author Names and Affiliations in order as they would appear in the published article

Samuel Garcia, Centre de Recherche en Neurosciences de Lyon, CNRS, Lyon, France
 Alessio P. Bucciono, Department of Biosystems Science and Engineering, ETH Zurich, Switzerland
 Pierre Yger, Institut de la Vision, Sorbonne Université, INSERM, Paris

4. Author Contributions: Each author must be identified with at least one of the following: Designed research, Performed research, Contributed unpublished reagents/ analytic tools, Analyzed data, Wrote the paper.
 Example: CS and JS Designed Research; MG and GT Performed Research; JS Wrote the paper

P.Y, S.G and A. B. all contributed equally to the manuscript

5. Correspondence should be addressed to (include email address)

samuel.garcia@cnrs.fr

6. Number of Figures 5

7. Number of Tables

8. Number of Multimedia

9. Number of words for Abstract 200

10. Number of words for Significance

Statement 85

11. Number of words for Introduction 483

12. Number of words for Discussion 522

13. Acknowledgements**14. Conflict of Interest**

A. No (State 'Authors report no conflict of interest')

B. Yes (Please explain)

15. Funding sources

This work was supported by the ETH Zurich Postdoctoral Fellowship 19-2 FEL-17 (APB)

How do spike collisions affect spike sorting performance?

Abstract

1 Recently, a new generation of devices have been developed to record neural activity simultane-
2 ously from hundreds of electrodes with a very high spatial density, both for *in vitro* and *in vivo*
3 applications. While these advances enable to record from many more cells, they also challenge the
4 already complicated process of *spike sorting* (i.e. extracting isolated single-neuron activity from
5 extracellular signals). In this work, we used synthetic ground-truth recordings with controlled lev-
6 els of correlations among neurons to quantitatively benchmark the performance of state-of-the-art
7 spike sorters focusing specifically on spike collisions. Our results show that while modern template-
8 matching based algorithms are more accurate than density-based approaches, all methods, to some
9 extent, failed to detect synchronous spike events of neurons with similar extracellular signals. In-
10 terestingly, the performance of the sorters is not largely affected by the the spiking activity in
11 the recordings, with respect to average firing rates and spike-train correlation levels. Since the
12 performances of all modern spike sorting algorithms can be affected as function of the activity of
13 the recorded neurons, scientific claims on correlations and synchrony should be carefully assessed
14 based on the analysis provided in this paper.
15

16 **keywords:** spike sorting, spike collision, benchmark, overlapping spikes

17 Significance statement

18 High-density extracellular recordings allow experimentalists to get access to the spiking activity of
19 large neuronal population, *via* the procedure of spike sorting. It is widely known that spike sorters
20 are affected by *spike collisions*, i.e., the occurrence of spatio-temporally overlapping events, but a
21 quantitative benchmark is still lacking. In this contribution, we perform systematic comparisons on
22 the performance of many different spike sorters against spike collisions, showing that modern spike
23 sorters, to different degrees, are still affected by synchronous events. Our results suggest that scientific
24 claims on neuron correlations and synchrony should be carefully assessed as they could result from
25 spike sorting errors.

26 Introduction

27 Accessing the activity of large ensemble of neurons is a crucial challenge in neuroscience. In recent years,
28 Multi-Electrode Arrays (MEA) and large silicon probes have been developed to record simultaneously
29 from hundreds of electrodes packed with a high spatial density, both *in vivo* [Angotzi et al., 2019, Jun
30 et al., 2017] and *in vitro* [Berdondini et al., 2009, Frey et al., 2009]. With these devices, each electrode
31 records the extracellular field in its vicinity and can detect the action potentials (or spikes) emitted by
32 the neighboring neurons in the tissue. In contrast to intracellular recording, extracellular recordings do
33 not give a direct and unambiguous access to single neuron activity and one needs to further process the
34 recorded signals to extract the spikes emitted by the different cells around the electrodes. This complex
35 problem of source separation is termed “spike sorting”. While various solutions for small number of
36 channels (tens at max) can be found in the large literature on spike sorting algorithms [Quiroga et al.,

37 2004], these new devices with thousands of channels challenge the *classical* approach to perform spike
38 sorting.

39 Recently, a new generation of spike sorting algorithms have been developed to be able to deal with
40 hundreds (or even thousands) of channels recorded simultaneously (see [Hennig et al., 2019, Lefebvre
41 et al., 2016] for recent reviews). The extent to which these modern spike sorting algorithm recover
42 all the spikes from a neuronal population is still under investigations, and might differ depending on
43 the species, tissue, cell types, activity level. While most of the real ground truth recordings [Neto
44 et al., 2016, Yger et al., 2018] are assessing the performance at the single cell level, in order to obtain
45 an exhaustive assessment of the spike sorting performance at the population level, one must turn to
46 use fully artificial or hybrid dataset [Buccino and Einevoll, 2020, Magland et al., 2020] to properly
47 compare and quantify the performances of the algorithms. But even with such dataset, in most of the
48 studies, errors are only measured as False Positive/Negative rates, and the reasons behind failures of
49 the algorithms are often overlooked.

50 In this study, we focused on a key property of the spike trains, at the core of most of these failures,
51 i.e. their fine temporal correlations. Indeed, temporal correlations are ubiquitous in the brain, and the
52 higher the number of recorded cells because of the increased density of the probes, the more prominent
53 they are. Correlations might have an important role in population coding ([Averbeck et al., 2006] for
54 a review), but correlated activity for nearby cells results, in the extracellular signals, in overlapping
55 activities and thus are harder to identify than isolated spikes. While pioneering work [Pillow et al.,
56 2013] claimed that template-matching based algorithms were more suited to recover overlapping spikes
57 (either in space and/or time), the extent to which they are is not properly defined. In this work, our
58 aim is to estimate how good (or bad) modern spike sorters are in recovering colliding spikes. What
59 are the limits of the sorters, and what are the key parameters of the recordings and/or of the neurons
60 that could influence these numbers?

61 Materials and Methods

62 All the code used to generate the figures is available at <https://spikeinterface.github.io/>.

63 Simulated datasets

64 We used the MEArec simulator [Buccino and Einevoll, 2020] to generate 30-minutes long synthetic
65 ground truth recordings. In brief, MEArec uses biophysically detailed multicompartment models to
66 simulate the extracellular action potentials, or so called "templates". For this study, we used 13 cell
67 models from layer 5 of a juvenile rat somatosensory cortex [Markram et al., 2015, Ramaswamy et al.,
68 2015] to get a dictionary of biologically plausible templates. Given this database, we took the layout
69 of a NeuroNexus probe (A1x32-Poly3-5mm-25s-177-CM32 with 32 electrodes in three columns and
70 hexagonal arrangement, a x- and y-pitch of 18 μm and 22 μm , respectively, and an electrode radius of
71 7.5 μm), and randomly positioned 20 cells in the vicinity of the probe, so that every simulated neuron
72 has a unique *template* (i.e. average extracellular action potential). Templates are then combined with
73 spike trains and slightly modulated in amplitude to add physiological variability. Additive uncorrelated
74 Gaussian noise is finally added to the traces. We generated simulated recordings with 20 neurons
75 randomly positioned in front of the probe, a noise level of 5 μV and a sampling rate of 32 kHz. To
76 obtain more robust results, we generated 5 recording per conditions with various random seeds. The
77 spike times were kept unchanged, but the positions and the templates of the 20 neurons were changed
78 in each of the individual recording. This allowed us to populate the distribution of cosine similarities
79 between pairs.

80 Generating spike trains with controlled correlations

81 To generate the recordings with various firing rates and correlations levels, we used the mixture pro-
82 cess method described in [Brette, 2009]. Since by default the method generates controlled cross-

83 correlograms with a decaying exponential profile, we modified it to generate cross-correlograms with
 84 a Gaussian profile, in order to have more synchronous firing for small lags. By setting three different
 85 rate levels (5, 10 and 15 Hz) and three different correlation levels (0, 10 and 20 %) this gave rise to 9
 86 conditions, so to 45 recordings in total (5 recordings per conditions, see above).

87 Template similarity

88 We define the template for neuron i as $\mathbf{T}_i \in \mathbb{R}^{T \times C}$, with T representing the number of samples and
 89 C the number of channels. After *flattening* the template by concatenating the signals from different
 90 channels ($\mathbf{T}_i^f \in \mathbb{R}^{T \cdot C}$), the similarity between two neurons i and j is quantified via the cosine similarity
 91 defined as follows:

$$similarity = \frac{\mathbf{T}_i^f \cdot \mathbf{T}_j^f}{\|\mathbf{T}_i^f\| \|\mathbf{T}_j^f\|} = \cos(\theta) \quad (1)$$

92 where θ is the angle between the two $(T \cdot C)$ -dimensional vectors \mathbf{T}_i^f and \mathbf{T}_j^f . The cosine similarity
 93 is therefore bounded between -1 (templates are anti-parallel) and 1 (templates are parallel). A cosine
 94 similarity of 0 means that the templates are orthogonal.

95 Spike sorters

96 All the spike sorters used in this study were run using the SpikeInterface framework [Buccino et al.,
 97 2020], with default parameters. The following are the exact versions that we used for the different spike
 98 sorters: Tridesclous (1.6.4), Spyking-circus (1.0.9) [Yger et al., 2018], HerdingSpikes (0.3.7) [Hilgen
 99 et al., 2017], Kilosort (v1, 2, or 3) [Pachitariu et al., 2016], YASS (2.0) [Lee et al., 2020], IronClust
 100 (5.9.8) [Chung et al., 2017], HDSort (1.0.3) [Diggelmann et al., 2018]. The desktop machine used has
 101 36 Intel Xeon(R) Gold 5220 CPU @ 2.20GHz, 200Go of RAM and a Quadro RTX 5000 with 16Gb of
 102 RAM as a GPU.

103 Spike sorting comparison

104 All the quantitative metrics between the results of the spike sorting software and the ground-truth
 105 recording were made via the SpikeInterface toolbox.

106 When comparing a spike sorting output to the ground-truth spiking activity, first an agreement
 107 score between each pair of ground-truth and sorted spike trains is computed as:

$$score_{ij} = \frac{\#n_{matches}}{\#n_{igt} + \#n_{jsorted} - \#n_{matches}}$$

108 where $\#n_{igt}$ and $\#n_{jsorted}$ are the numbers of spikes in the i -th ground-truth spike train and the
 109 j -th sorted spike trains, respectively. $\#n_{matches}$ is the number of spikes within 0.4 ms between the
 110 two spike trains.

111 Once scores for all pairs are computed, an hungarian assignment is used to match ground-truth
 112 units to sorted units [Buccino et al., 2020]. All spikes from matched spike trains are then labeled as:
 113 true positive (TP), if the spike is found both in the ground-truth and the sorted spike train; false
 114 positive (FP), if the spike is found in the sorted spike train, but not in the ground-truth one; and false
 115 negative (FN), if the spike is only found in the ground-truth spike train.

116 After labeling all matched spikes, we can now define these unit-wise performance metrics for each
 117 ground-truth unit that has been matched to a sorted unit:

$$accuracy = \frac{\#TP}{\#TP + \#FP + \#FN} \quad (2)$$

$$precision = \frac{\#TP}{\#TP + \#FP} \quad (3)$$

$$recall = \frac{\#TP}{\#TP + \#FN} \quad (4)$$

118 The global accuracy, precision, and recall values shown in Figure 2D are the average values of the
119 performance metrics computed by unit.

120 Using the unit metrics and the output of the matching procedure, we can further classify each
121 sorted unit as:

122 **well detected:** sorted units with an accuracy ≥ 0.8

123 **false positive:** sorted units that are not matched to any ground-truth unit and have a score < 0.2

124 **redundant:** sorted units that are not the best match to a ground-truth unit but have a score ≥ 0.2

125 **overmerged:** sorted units with a score ≥ 0.2 with more than one ground-truth unit

126 In order to generate the spike lag versus recall figures (e.g. Figures 3-6) we expanded the SpikeIn-
127 terface software with several novel comparison methods and visualization widgets. In particular, we
128 extended the ground-truth comparison class to the `CollisionGTComparison`, which computes per-
129 formance metrics by spike lag. In addition to the agreement score computation and the matching
130 described in the previous paragraphs, this method first detects and flags all “synchronous spike events”
131 in the ground-truth spike trains. Two spikes from two separate units are considered to be a “syn-
132 chronous spike event” if their spike times occur within a time lag of 2 ms. The synchronous events
133 are then binned in 11 bins spanning the $[-2, 2]$ ms interval and the *collision recall* is computed for
134 each bin. With a similar principle, we implemented the `CorrelogramGTComparison` to compute the
135 lag-wise relative errors in cross-correlograms between ground-truth units and spike sorted units.

136 Results

137 Generation of the ground-truth recordings

138 To test how robust the recently developed spike sorting pipelines are against spike collisions [Chung
139 et al., 2017, Hilgen et al., 2017, Lee et al., 2020, Pachitariu et al., 2016, Yger et al., 2018], we generated
140 synthetic datasets using the `MEAREC` simulator [Buccino and Einevoll, 2020] (see Methods). More
141 precisely, we took the layout of a NeuroNexus probe with 32 electrodes in three columns and hexagonal
142 arrangement, and randomly positioned 20 cells in the vicinity of the probe (see Figure 1A), so that
143 every simulated neuron has a unique *template* (i.e. average extracellular action potential). Figure 1B
144 shows three sample templates with respectively low, almost null, and high similarity. The similarity
145 between templates is computed as the cosine similarity of the flattened signals (see Methods) and the
146 random generation of the positions and cell types of the simulated neurons (and thus of the templates)
147 gives rise to the similarity matrix displayed in see Figure 1C. This similarity, as expected, decreases
148 with the distance between the neurons, computed either from the ground-truth positions of the cells
149 from the simulation or estimated as the barycenters of the templates (Figure 1D). The more negative
150 the similarity is, the more templates are “in opposition”; the more positive it is, the more templates
151 are “similar”. A similarity close to 0 means that templates do not overlap and are strongly orthogonal,
152 i.e. dissimilar. Importantly, the simulations allowed us to cover rather uniformly the space of cosine
153 similarities between templates, which will be used to assess the performance of spike sorters during
154 collisions (Figure 1E).

155 To generate the spike trains, we first used a simple approach that forced all the neurons to fire as
156 independent Poisson sources at a fixed and homogeneous firing rate of 5 Hz. To make the simulation
157 more biologically plausible, we pruned all spikes breaking a refractory period violation of 4 ms. The
158 resulting auto- and cross-correlograms for three sample units are shown in Figure 1F (auto-correlograms
159 are in green on the diagonal), while Figure 1G and H display the average (red line) and standard

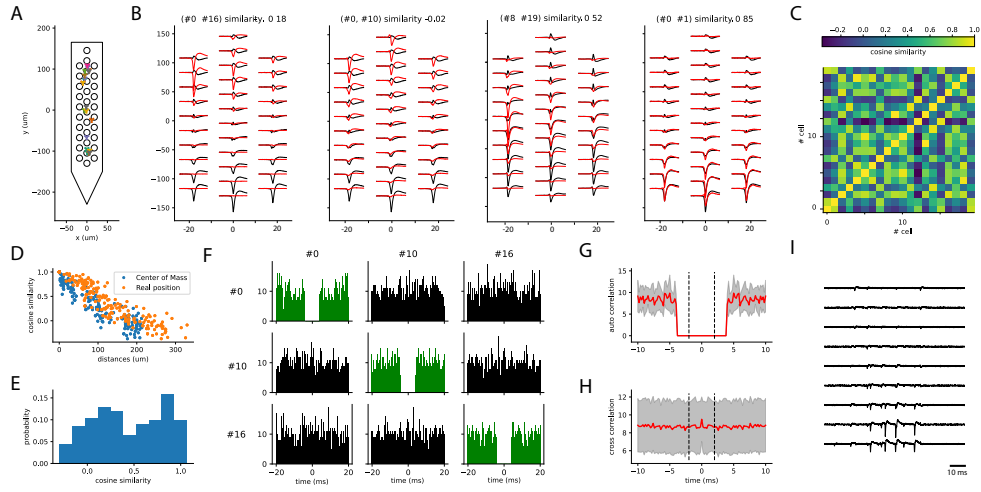


Figure 1: **Generation of the synthetic recordings.** **A)** 20 cells are randomly placed in front of a 32-channel NeuroNexus probe layout. The plot shows the location of each cell for one recording. **B)** Sample templates generated by neurons that are close too each other (#0 and #1) or far apart (#2 and #3). **C)** Cosine similarity matrix between all pairs of templates for a sample recording. **D)** Cosine similarity as function of the distance between the neurons, either using the real position from the simulations (orange circles), or the estimated barycenter of the templates (blue circles). **E)** Histogram of the cosine similarity distribution from one of the simulated recordings. **F)** Cross- and auto- correlograms for three sample spike trains. **G)** Average auto-correlograms of all units (red line, gray area represents the standard deviation). **H)** Average cross-correlogram over all pairs of neurons (red line, gray area represents the standard deviation around the mean). **I)** Sample traces from 10 channels of one synthetic recording.

160 deviation (grey shaded area) auto- and cross-correlation among all units, respectively. A sample
 161 snippet of the generated traces from one recording is shown in Figure 1I, for a subset of 10 channels
 162 out of 32. Due to the independence of the Poisson sources, both the average cross-correlograms
 163 (Figure 1G) and auto-correlograms – outside the ± 4 ms used as refractory period – (Figure 1H) are
 164 flat.

165 **Global performance of the spike sorters**

166 In order to assess the global performances of the sorting procedure on our synthetic datasets, we
 167 generated 5 recordings with various random seeds and averaged the results. Figure 2 summarizes the
 168 main findings. First, we noticed that, as seen in Figure 2A, the run time was roughly constant across
 169 sorters, except for HDSort, with its higher run time. The number of well detected units is similar
 170 among sorters, as shown in Figure 2B, but it is worthwhile noticing that Kilosort 3 is the only sorter
 171 producing many false positive and redundant units (see Methods for classification of units). Kilosort 2
 172 and HDSort also identify more false positive than well detected units. Importantly, we did not perform
 173 any curation of the spike sorting output, but we consider the raw output of each sorter as is.

174 To check whether all sorters correctly *discovered* all templates, we computed the cosine similarity
 175 between the ground-truth templates from the simulations and the ones found by the sorters, comparing
 176 such a metric with the accuracy. By doing so, we wanted to rule out the fact that the sources of the
 177 errors could primarily be due to problems in the clustering. Indeed, if the spike sorting algorithms are

178 properly behaving, they should find templates very similar to the ground-truth ones. As it can be seen
 179 in Figure 2C, all sorters are on average finding the correct templates, with the notable exception of
 180 YASS (in grey) and to some less extent HDSort (in red). The average cosine similarity between found
 181 and ground-truth templates is larger than 0.97 for most template-matching based sorters (Spyking-
 182 circus, Kilosort 1/2/3, IronClust, Tridesclous), so we can safely assume that most of the errors are
 183 not due to the clustering step. Moreover, the overall accuracy of most of the spike sorters is relatively
 184 high (~ 0.95), except for HDSort and HerdingSpikes which yield lower scores (Figure 2D). However,
 185 this averaged number does not tell us anything regarding the nature of these errors. While this error
 186 rate might seem low, it is likely that it is crucial, since it can potentially originate from the collisions,
 187 and thus from the correlations among neurons.

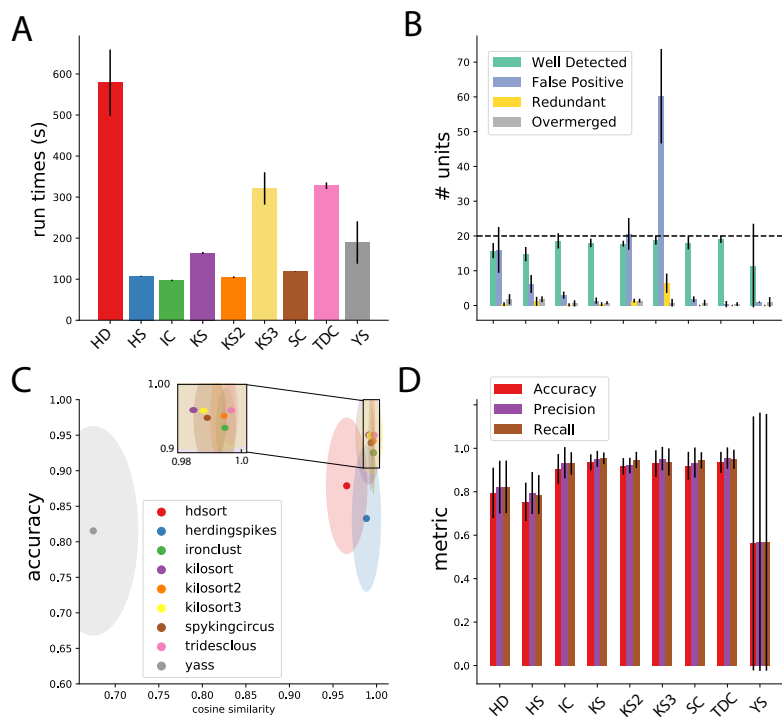


Figure 2: **Spike sorting performance.** **A)** Average run times over 5 different recordings (see Methods) for all the tested sorters. Errors bars indicate the standard deviation over multiple recordings. **B)** Average number of cells found by the sorters that are either well detected, redundant, overmerged or considered as false positive (see Methods). Error bars indicates standard deviation over multiple recordings. **C)** The average cosine similarity between templates found by the sorters and ground-truth templates, as function of the accuracy for the given neurons. Ellipses shows standard error of the means in cosine similarity (x-axis) and accuracy (y-axis). **D)** Average metrics (accuracy, precision, recall, see Methods) for all the sorters. Error bars show standard deviation over multiple recordings.

188 Spike sorting performance is affected by spike collisions

189 Using fully synthetic recordings with exhaustive ground truth, we can look at how good individual
190 spike sorters perform specifically with respect to spatio-temporal collisions. To do so, we computed
191 the *collision recall* (see Methods) as a function of the lag between two spikes, for a given pair of
192 neurons. By averaging over multiple pairs of ground-truth neurons with similar template similarity
193 (and over multiple recordings, see Methods), we can obtain a picture of how accurate the sorters
194 are specifically with respect to the spike time lags and the similarities between templates. Figure 3
195 displays the collision recall per sorter as a function of the lag (x-axis), colored by the similarity between
196 templates. Each panel shows the performance of a different spike sorter. One can immediately see
197 that only few sorters are able to accurately resolve lag correlations that are close to zero, even when
198 templates are strongly orthogonal (low cosine similarity). This is the case for Kilosort 1 and 2, and for
199 Spyking-circus, all of which use a template-matching procedure that should theoretically explain this
200 behavior. It is worthwhile noting that the decrease in performance for Kilosort 3 is surprising, since the
201 authors confirmed the software is using the exact same template-matching procedure than in previous
202 versions. This means that errors are likely originating either from subtle variations in the preprocessing
203 steps, and/or in the clustering that has been changed and thus might lead to slight differences in the
204 templates. However, while performances are still good for Kilosort 1 and 2 even when the average cosine
205 similarity between pairs is increased, they slightly degrade for Spyking-circus. Density-based sorters
206 (HerdingSpikes and IronClust), on the other hand, do not have a spike collision resolution strategy
207 and this is reflected by their overall poorer performance. It is interesting to notice that Tridesclous,
208 HDSort, YASS, and Kilosort 3, also using a template-matching based procedure to resolve the spikes,
209 are not properly resolving the temporal correlations even for dissimilar templates. Different template-
210 matching strategies are probably the cause of the differences among sorters. For example, HDSort and
211 HerdingSpikes do not implement any strategy for spike collision resolution [Diggelmann et al., 2018]
212 and that is reflected in the quick degradation of performance as template similarity increases. Kilosort
213 uses a GPU-based implementation of the k-SVD algorithm [Aharon et al., 2006], used in matching
214 learning as a dictionary learning algorithm for creating a dictionary for sparse representations. By
215 doing so, it performs a reconstruction of the extracellular traces by optimizing both the templates and
216 the spike times, which is an enhancement compared to what is done in Spyking-circus and Tridesclous.
217 This might explain the boost in performance especially striking for templates with high similarity
218 (*similarity* > 0.8).

219 Generation of controlled spike collision simulated data

220 The results shown in the previous section have been obtained only in a particular regime of activity,
221 with all neurons firing independently as Poisson sources with an average firing rate of 5 Hz. However,
222 neurons usually do not fire independently of each other, but rather have intrinsic correlations, also
223 depending on different brain areas, brain states, and species. In addition, the average firing rates can
224 also largely vary depending on brain areas. As an example, it is well known that Purkinje cells in the
225 cerebellum have a very high firing rate [Sedaghat-Nejad et al., 2021] that networks tends to synchronize
226 their activity either in slow waves during sleep [?], or during pathological activity (such as epileptic
227 seizures [?]). Therefore, assessing how performances may vary during different conditions is important
228 to generalize our observations.

229 In order to study how spike sorting is affected by correlations and firing rates, we used a mixture pro-
230 cedure [Brette, 2009] that allowed us to control precisely the shape of the auto- and cross-correlograms
231 for the injected spike trains. More precisely, we decided to explore in a systematic manner three rate
232 levels (5, 10 and 15 Hz), and three correlation levels (0, 10, and 20 %). Note that the 5 Hz firing rate
233 with 0 % correlation corresponds to the scenario displayed in Figures 2-3.

234 Figure 4 shows the average of cross- and auto-correlograms and the spike trains of a recording where
235 cells are firing as independent Poisson sources at 5 Hz in panels A-C (and thus with 0 % correlation,
236 as shown by the flat average cross-correlograms in Figure 4A) and at 15 Hz with 20 % correlation

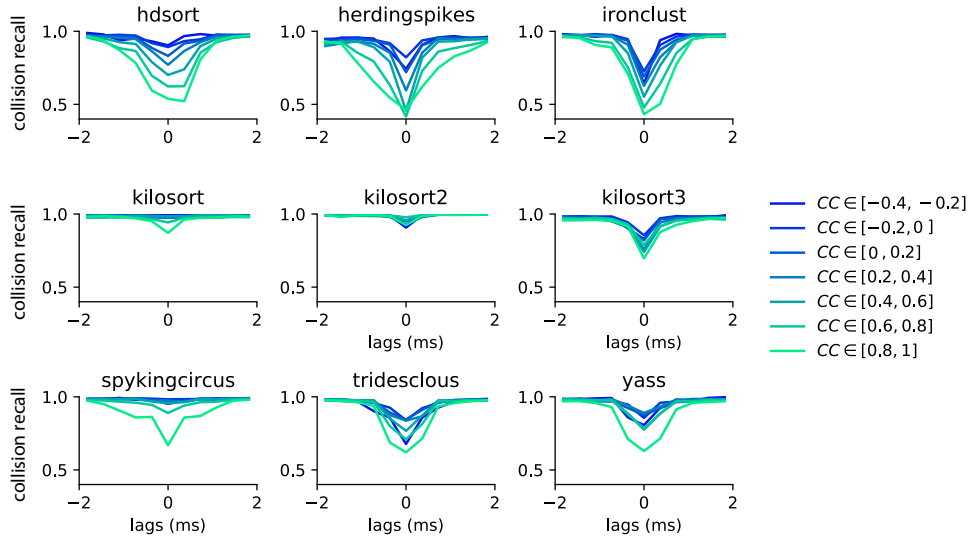


Figure 3: **Collision recall per sorter.** Error (quantified as the collision recall, see Methods) for various sorters and for all possible lags (between -2 and 2 ms), as function of the similarity between the pairs of templates (color code). All curves are averaged over multiple pairs and multiple recordings (see Methods).

237 (Figure 4D-F). Even though experimental recordings would contain a broader spectrum of firing rates
 238 and correlations, here we focus on assessing how different firing regimes affect spike sorting performance
 239 in a controlled setting. By varying these conditions, we wanted to challenge the internal clustering step
 240 of the spike sorting algorithms and see how generalizable are the results we observed in the previous
 241 section. One would expect that the increased density of spikes (both in terms of firing rates and of
 242 synchrony) should degrade the performance of the spike sorters by affecting both the clustering step
 243 and the template-matching step, which in turn would degrade the resolution of spike collisions. It is
 244 worthwhile noting that all the rates and correlation levels are homogeneous among neurons and only
 245 the templates are different.

246 **Do correlations and firing rates affect spike sorting of spike collisions?**

247 To assess whether firing rate and spike train correlation affect spike sorting performance, we selected
 248 all unit pairs with a similarity greater than 0.5. We first averaged the recall curves for all template
 249 similarities (i.e. we averaged the curves with similarity greater than 0.5 shown in Figure 3).

250 In Figure 5A we show the recall with respect to the spike lags averaged over all 9 configurations (3
 251 firing rates x 3 correlations) for each sorter. The thick line represents the mean recall and the shaded
 252 area is the standard deviation over different rate-correlation configuration. All sorters, except YASS,
 253 appear to have a very consistent curve (low standard deviation) over different configurations and do
 254 not seem affected by changes in average firing rates and correlations in the spike trains. YASS' large
 255 standard deviation can be explained by looking at individual recall curves at different rate-correlation
 256 regimes (Figure 6 - yellow lines): the spike sorting performance degrades with increasing firing rates,
 257 but it does not seem to be strongly affected by increased correlation rates. However, we should stress
 258 that since the collision recall is a relative measure, the same value for a larger number of spikes (when

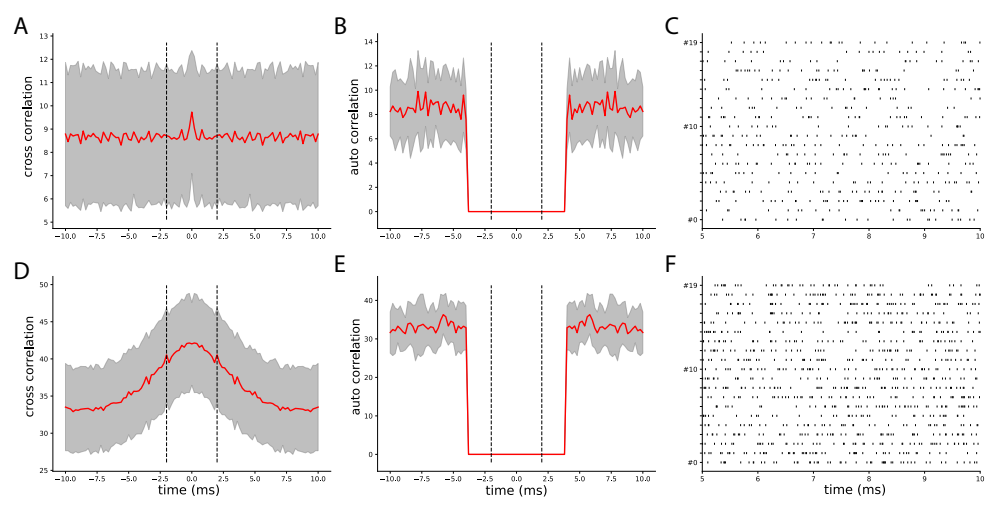


Figure 4: **Controlling spike trains correlations and firing rates.** **A)** Average cross-correlograms between all pairs of distinct neurons firing as independent Poisson sources at 5Hz (red curve, gray area represents the standard deviation) **B)** Same as **A**, but for auto-correlograms. **C)** Rater plot showing the activity of the uncorrelated neurons firing at 5Hz. **D-F)** Same as **A-C**, but for a rate of 15 Hz and 20 % correlation.

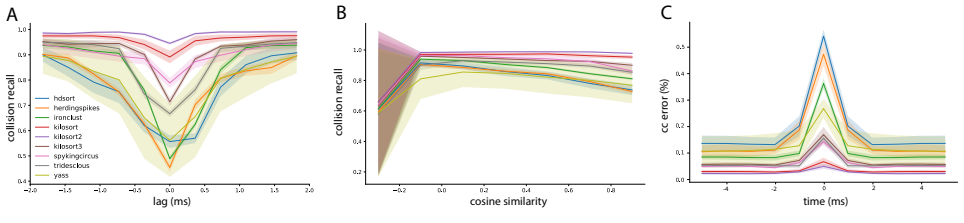


Figure 5: **Spike sorting performance for different conditions.** **A)** Average collision recall over the 9 conditions shown in Figure 6 (3 firing rate levels and 3 correlation levels) as function of the lag between spikes, for pairs of cells with cosine similarity higher than 0.5. The shaded area shows the standard deviation over the conditions. **B)** Similarly as **A**, the average collision recall as function of the cosine similarity between pairs of cells. **C)** Mean relative error between the ground-truth cross-correlograms and the estimated ones, for all sorters, averaged over all pairs with a similarity higher than 0.5

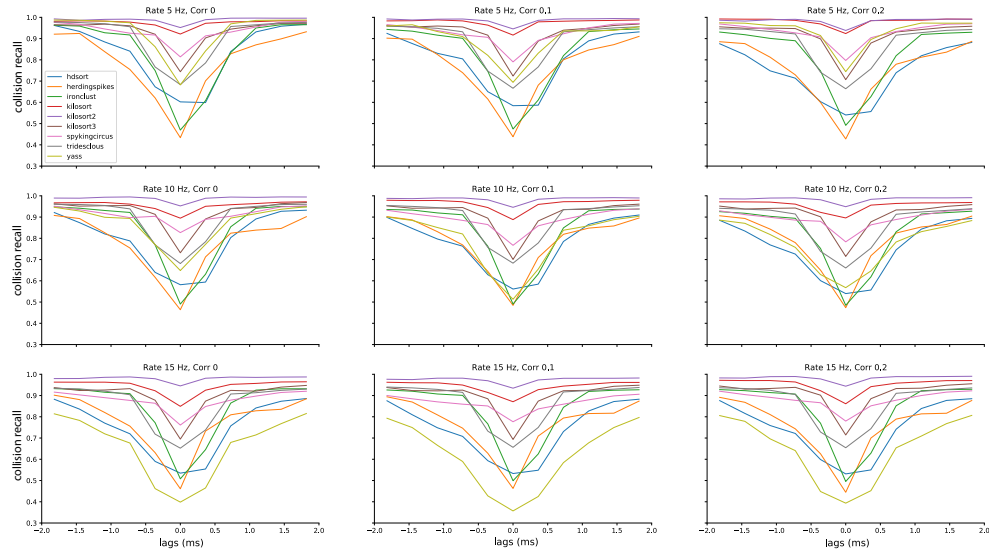


Figure 6: Average performances of the spike sorters as function of the temporal lags. Each panel shows the average collision recall for template pairs with a similarity above 0.5 for a different condition, in terms of firing rate and correlation levels.

259 firing rate is increased) means that overall, there are more misses for all sorters.
 260 Similar considerations can be done by looking at the average recall with respect to template simi-
 261 larity (Figure 5B). To construct this plots, we integrated the curves in Figure 3 over lags for different
 262 cosine similarities. Also in this case, the curves appear consistent (low standard deviation) with the
 263 exception of YASS, for which recall is reduced with increased firing rate regimes (Figure 7 - yellow
 264 lines). It is worth noticing that when the cosine similarity becomes negative, all the sorters perform
 265 very poorly in properly resolving the overlaps. This could be explained by the fact that when a pair
 266 of templates is anti-parallel (for example in the left panel of Figure 1A), a subset of electrodes might
 267 show a negative signal for one template and a positive signal from the other (due to return currents
 268 in the dendritic signals [Gold et al., 2009]). Effectively, when a spike collision between the two spikes
 269 occur, this would lower the amplitude of the negative peak, which could reduce the detectability of
 270 the spike.
 271 The collision recall metric is mostly useful to obtain a quantitative insight on the behavior of the
 272 spike sorting algorithms, but how do these errors transpose in practical situations? To assess this, we
 273 measure the relative error (in percentage) between the ground-truth cross-correlograms and the ones
 274 computed from the spike sorting outputs. We then averaged these error curves among all recordings
 275 and experimental conditions (firing rates and synchrony levels). As shown in Figure 5, the error in
 276 the estimated cross-correlogram can be as large as more than 50% for small lags, and for some spike
 277 sorting algorithms such as HDSort, HerdingSpikes or IronClust. Moreover, it is also worth noticing
 278 the baseline error rate is not the uniform across sorters. From this metric, we can again conclude that
 279 template-matching based spike sorting algorithms such as Kilosort (1, 2, and 3), Spyking-circus or
 280 Tridesclous are much better to resolve fine temporal correlations among neurons.

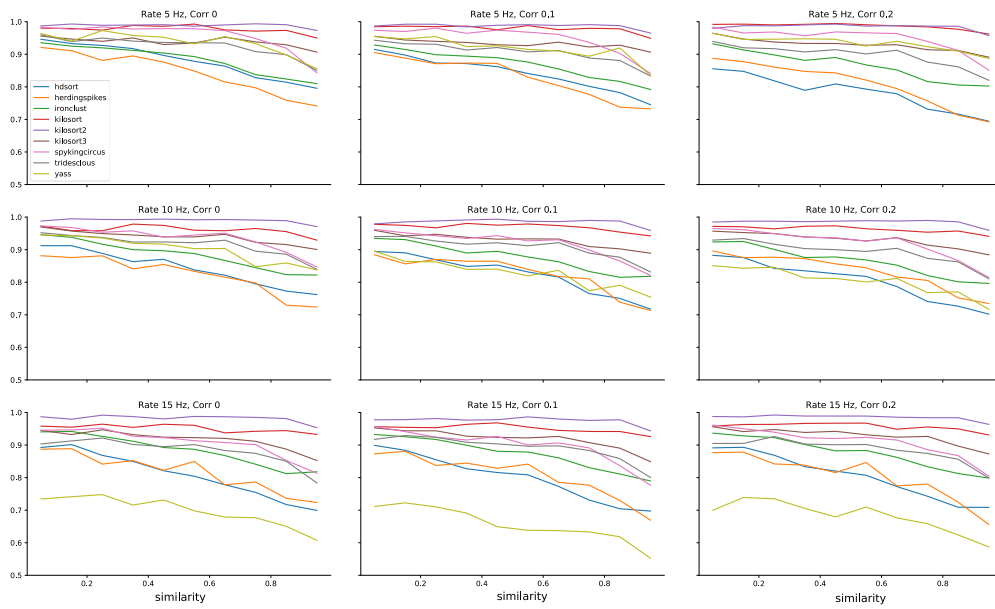


Figure 7: Average performances of the spike sorters as function of the template similarity. Each panel shows the average collision recall over all lags in $[-2, 2]$ ms for a different condition, in terms of firing rate and correlation levels.

281 Discussion

282 In this study, we showed in a systematic and quantitative manner how spatio-temporal correlations
283 can be underestimated during spike sorting. Using synthetic datasets, we compared a large diversity of
284 modern spike sorters and showed how they behaved as function of the similarity between the templates
285 and the temporal lags between spikes. As expected, the closer the spikes are in time, the harder is
286 it, for all sorter, to properly resolve the overlaps. However, more interestingly, the more similar the
287 templates are, the higher the failures are. These failures are striking especially for spike sorters that are
288 not relying on template-matching based approaches (HerdingSpikes, IronClust). For the ones using a
289 template-matching based approach (Kilosort, Spyking-circus, Tridesclous, HDSort), the problem is less
290 pronounced (with the exception of HDSort) but still present, and therefore this phenomenon should
291 be taken into account when making claims about the synchrony.

292 To our surprise, the global behavior of the spike sorters did not depend much on the overall
293 firing rate and/or the correlation levels. This allows us to generalize the findings and we think that
294 the quantitative results shown here could be translated to various *in vitro* or *in vivo* recordings from
295 different brain regions and species. As shown in Figure 5, while the variability over different conditions
296 is rather high for some algorithms, template-matching based algorithms tend to be rather robust and
297 overall better in resolving spike collisions. This is a very encouraging sign towards a unified and
298 reproducible automated solution for spike sorting [Buccino et al., 2020, Magland et al., 2020], agnostic
299 of the recording conditions.

300 The results shown in the paper were obtained with purely artificial recordings, since we need
301 exhaustive information on the ground-truth spiking activity of all neurons to quantitatively compare
302 and benchmark different spike sorters. However, it would be interesting to generalize these observations
303 with real recordings, assuming one would have a proper ground truth at the population level. Indeed,
304 such a ground truth is needed to compute the *collision recall* and see how sorters behave as function
305 of lags and similarities between templates. To our knowledge, such a ground truth does not exist
306 [Diggelmann et al., 2018, Neto et al., 2016, Yger et al., 2018]. While one could try to generate an
307 “approximated” ground truth by combining the output of several spike sorters with an *ensemble* spike
308 sorting approach (as in [Buccino et al., 2020]), the disagreements among sorters are currently so high
309 that this process is hard if not impossible, if one wants to sample from a large number of pairs.

310 While missing spikes for very dissimilar templates and small lags is problematic, the errors made for
311 very similar templates may be less frequent depending on the probe layout and neuronal preparation.
312 Indeed, such errors strongly depend on the distribution of template similarities between all pairs of
313 recorded cells, and this distribution might differ from recording to recording. For example, in the
314 retina [Wässle, 2004] one would expect highly synchronous cells, of the same functional type, to be far
315 apart from each other because of an intrinsic tiling of the visual space. Such properties are unknown *in*
316 *vivo* or in cortical structures, but might bias the distribution of template similarities between nearby
317 neurons, and thus modify the estimation of collision recalls.

318 References

- 319 M. Aharon, M. Elad, and A. Bruckstein. rm K-SVD: An Algorithm for Designing Overcomplete
320 Dictionaries for Sparse Representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322,
321 November 2006.
- 322 Gian Nicola Angotzi, Fabio Boi, Aziliz Lecomte, Ermanno Miele, Mario Malerba, Stefano Zucca,
323 Antonino Casile, and Luca Berdondini. Sinaps: An implantable active pixel sensor cmos-probe for
324 simultaneous large-scale neural recordings. *Biosensors and Bioelectronics*, 126:355–364, 2019.
- 325 Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations, population coding
326 and computation. *Nature reviews neuroscience*, 7(5):358–366, 2006.

- 327 Luca Berdondini, Kilian Imfeld, Alessandro Maccione, Mariateresa Tedesco, Simon Neukom, Milena
328 Koudelka-Hep, and Sergio Martinoia. Active pixel sensor array for high spatio-temporal resolution
329 electrophysiological recordings from single cell to large scale neuronal networks. *Lab on a Chip*, 9
330 (18):2644–2651, 2009.
- 331 Romain Brette. Generation of correlated spike trains. *Neural computation*, 21(1):188–215, 2009.
- 332 Alessio P Buccino, Cole L Hurwitz, Samuel Garcia, Jeremy Magland, Joshua H Siegle, Roger Hurwitz,
333 and Matthias H Hennig. Spikeinterface, a unified framework for spike sorting. *Elife*, 9:e61834, 2020.
- 334 Alessio Paolo Buccino and Gaute Tomas Einevoll. Mearec: a fast and customizable testbench simulator
335 for ground-truth extracellular spiking activity. *Neuroinformatics*, pages 1–20, 2020.
- 336 Jason E Chung, Jeremy F Magland, Alex H Barnett, et al. A fully automated approach to spike
337 sorting. *Neuron*, 95(6):1381–1394, 2017.
- 338 Roland Diggelmann, Michele Fiscella, Andreas Hierlemann, and Felix Franke. Automatic spike sorting
339 for high-density microelectrode arrays. *Journal of neurophysiology*, 120(6):3155–3171, 2018.
- 340 Urs Frey, U Egert, F Heer, S Hafizovic, and Andreas Hierlemann. Microelectronic system for high-
341 resolution mapping of extracellular electric fields applied to brain slices. *Biosensors and Bioelec-
342 tronics*, 24(7):2191–2198, 2009.
- 343 Carl Gold, Cyrille C Girardin, Kevan AC Martin, and Christof Koch. High-amplitude positive spikes
344 recorded extracellularly in cat visual cortex. *Journal of neurophysiology*, 102(6):3340–3351, 2009.
- 345 Matthias H Hennig, Cole Hurwitz, and Martino Sorbaro. Scaling spike detection and sorting for
346 next-generation electrophysiology. *In Vitro Neuronal Networks*, pages 171–184, 2019.
- 347 Gerrit Hilgen, Martino Sorbaro, Sahar Pirmoradian, Jens-Oliver Muthmann, Ibolya Edit Kepiro, Si-
348 mona Ullo, Cesar Juarez Ramirez, Albert Puente Encinas, Alessandro Maccione, Luca Berdondini,
349 et al. Unsupervised spike sorting for large-scale, high-density multielectrode arrays. *Cell reports*, 18
350 (10):2521–2532, 2017.
- 351 James J Jun, Nicholas A Steinmetz, Joshua H Siegle, Daniel J Denman, Marius Bauza, Brian Barbarits,
352 Albert K Lee, Costas A Anastassiou, Alexandru Andrei, Çağatay Aydın, et al. Fully integrated silicon
353 probes for high-density recording of neural activity. *Nature*, 551(7679):232, 2017.
- 354 JinHyung Lee, Catalin Mitelut, Hooshmand Shokri, Ian Kinsella, Nishchal Dethe, Shenghao Wu, Kevin
355 Li, Eduardo B Reyes, Denis Turcu, Eleanor Batty, et al. Yass: Yet another spike sorter applied to
356 large-scale multi-electrode array recordings in primate retina. *bioRxiv*, 2020.
- 357 Baptiste Lefebvre, Pierre Yger, and Olivier Marre. Recent progress in multi-electrode spike sorting
358 methods. *Journal of Physiology-Paris*, 110(4):327–335, 2016.
- 359 Jeremy Magland, James J Jun, Elizabeth Lovero, Alexander J Morley, Cole Lincoln Hurwitz,
360 Alessio Paolo Buccino, Samuel Garcia, and Alex H Barnett. Spikeforest, reproducible web-facing
361 ground-truth validation of automated neural spike sorters. *Elife*, 9:e55167, 2020.
- 362 Henry Markram, Eilif Muller, Srikanth Ramaswamy, et al. Reconstruction and simulation of neocortical
363 microcircuitry. *Cell*, 163(2):456–492, 2015.
- 364 Joana P Neto, Gonçalo Lopes, João Frazão, Joana Nogueira, Pedro Lacerda, Pedro Baião, Arno Aarts,
365 Alexandru Andrei, Silke Musa, Elvira Fortunato, et al. Validating silicon polytrodes with paired
366 juxtacellular recordings: method and dataset. *Journal of neurophysiology*, 116(2):892–903, 2016.

- 367 Marius Pachitariu, Nicholas A Steinmetz, Shabnam N Kadir, et al. Fast and accurate spike sorting
368 of high-channel count probes with kilosort. In *Advances in Neural Information Processing Systems*,
369 pages 4448–4456, 2016.
- 370 Jonathan W Pillow, Jonathon Shlens, EJ Chichilnisky, and Eero P Simoncelli. A model-based spike
371 sorting algorithm for removing correlation artifacts in multi-neuron recordings. *PloS one*, 8(5):
372 e62123, 2013.
- 373 R Quian Quiroga, Zoltan Nadasdy, and Yoram Ben-Shaul. Unsupervised spike detection and sorting
374 with wavelets and superparamagnetic clustering. *Neural computation*, 16(8):1661–1687, 2004.
- 375 Srikanth Ramaswamy, Jean-Denis Courcol, Marwan Abdellah, et al. The neocortical microcircuit
376 collaboration portal: a resource for rat somatosensory cortex. *Front Neural Circuits*, 9, 2015.
- 377 Ehsan Sedaghat-Nejad, Mohammad Amin Fakharian, Jay Pi, Paul Hage, Yoshiko Kojima, Robi Soet-
378 edjo, Shogo Ohmae, Javier F Medina, and Reza Shadmehr. P-sort: an open-source software for
379 cerebellar neurophysiology. *bioRxiv*, 2021.
- 380 H. Wässle. Parallel processing in the mammalian retina. *Nat Rev Neurosci*, 5(10):747–757, Oct 2004.
- 381 Pierre Yger, Giulia LB Spampinato, Elric Esposito, Baptiste Lefebvre, Stéphane Deny, Christophe
382 Gardella, Marcel Stimberg, Florian Jetter, Guenther Zeck, Serge Picaud, et al. A spike sorting
383 toolbox for up to thousands of electrodes validated with ground truth recordings in vitro and in
384 vivo. *Elife*, 7:e34518, 2018.