

How do spike collisions affect spike sorting performance?

Samuel Garcia^{¶*1}, Alessio P. Buccino^{¶†2}, and Pierre Yger^{¶‡3}

¹Centre de Recherche en Neurosciences de Lyon, CNRS, Lyon, France

²Department of Biosystems Science and Engineering, ETH Zurich, Switzerland

³Institut de la Vision, Sorbonne Université, INSERM, Paris

¶ These authors contributed equally to this work.

Abstract

Recently, a new generation of devices have been developed to record neural activity simultaneously from hundreds of electrodes with a very high spatial density, both for *in vitro* and *in vivo* applications. While these advances enable to record from many more cells, they also dramatically increase the amount of overlapping “synchronous” spikes (colliding in space and/or in time), challenging the already complicated process of *spike sorting* (i.e. extracting isolated single-neuron activity from extracellular signals). In this work, we used synthetic ground-truth recordings to quantitatively benchmark the performance of state-of-the-art spike sorters focusing specifically on spike collisions. Our results show that while modern template-matching based algorithms are more accurate than density-based approaches, all methods, to some extent, failed to detect synchronous spike events of neurons with similar extracellular signals. Interestingly, the performance of the sorters is not largely affected by the spiking activity in the recordings, with respect to average firing rates and spike-train correlation levels.

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

1 Introduction

Accessing the activity of large ensemble of neurons is a crucial challenge in neuroscience. In recent years, Multi-Electrode Arrays (MEA) and large silicon probes have been developed to record simultaneously from hundreds of electrodes packed with a high spatial density, both *in vivo* [14, 2] and *in vitro* [10, 4]. With these devices, each electrode records the extracellular field in its vicinity and can detect the action potentials (or spikes) emitted by the neighboring neurons in the tissue. In contrast to intracellular recording, extracellular recordings do not give a direct and unambiguous access to single neuron activity and one needs to further process the recorded signals to extract the spikes emitted by the different cells around the electrodes. This complex problem of source separation is termed “spike sorting”. While various solutions for small number of channels (tens at max) can be found in the large literature on spike sorting algorithms [22], these new devices with thousands of channels challenge the *classical* approach to perform spike sorting.

Recently, a new generation of spike sorting algorithms have been developed to be able to deal with hundreds (or even thousands) of channels recorded simultaneously (see [16, 12] for recent reviews). The

*samuel.garcia@crns.fr

†alessio.buccino@bsse.ethz.ch

‡pierre.yger@inserm.fr

30 extent to which these modern spike sorting algorithm recover all the spikes from a neuronal population
31 is still under investigations, and might differ depending on the species, tissue, cell types, activity level.
32 While most of the real ground truth recordings [28, 19] are assessing the performance at the single cell
33 level, in order to obtain an exhaustive assessment of the spike sorting performance at the population
34 level, one must turn to use fully artificial or hybrid dataset [17, 6] to properly compare and quantify
35 the performances of the algorithms. But even with such dataset, in most of the studies, errors are only
36 measured as False Positive/Negative rates, and the reasons behind failures of the algorithms are often
37 overlooked.

38 In this study, we focused on a key property of the spike trains, at the core of most of these failures,
39 i.e. their fine temporal correlations. Indeed, temporal correlations are ubiquitous in the brain, and the
40 higher the number of recorded cells because of the increased density of the probes, the more prominent
41 they are. Correlations might have an important role in population coding ([3] for a review), but
42 correlated activity for nearby cells results, in the extracellular signals, in overlapping activities and
43 thus are harder to identify than isolated spikes. While pioneering work [21] claimed that template-
44 matching based algorithms were more suited to recover overlapping spikes (either in space and/or
45 time), the extent to which they are is not properly defined. In this work, our aim is to estimate
46 how good (or bad) modern spike sorters are in recovering colliding spikes. What are the limits of the
47 sorters, and what are the key parameters of the recordings and/or of the neurons that could influence
48 these numbers?

49 2 Results

50 2.1 Simulated recordings

51 To test how robust the recently developed spike sorting pipelines are against spike collisions [28, 20, 8,
52 13, 15], we generated synthetic datasets using the MEArec simulator [6] (see Methods). More precisely,
53 we took the layout of a NeuroNexus probe (A1x32-Poly3-5mm-25s-177-CM32 with 32 electrodes in
54 three columns and hexagonal arrangement, a x- and y-pitch of 18 μm and 22 μm , respectively, and
55 an electrode radius of 7.5 μm), and randomly positioned 20 cells in the vicinity of the probe (see
56 Figure 1A), so that every simulated neuron has a unique *template* (i.e. average extracellular action
57 potential). Figure 1B shows three sample templates with respectively low, almost null, and high
58 similarity. The similarity between templates is computed as the cosine similarity of the flattened
59 signals (see Methods) and the random generation of the positions and cell types of the simulated
60 neurons (and thus of the templates) gives rise to the similarity matrix displayed in see Figure 1C. This
61 similarity, as expected, decreases with the distance between the neurons, computed either from the
62 ground-truth positions of the cells from the simulation or estimated as the barycenters of the templates
63 (Figure 1D). The more negative the similarity is, the more templates are “in opposition”; the more
64 positive it is, the more templates are “similar”. A similarity close to 0 means that templates do not
65 overlap and are strongly orthogonal, i.e. dissimilar. Importantly, the simulations allowed us to cover
66 rather uniformly the space of cosine similarities between templates, which will be used to assess the
67 performance of spike sorters during collisions (Figure 1E).

68 To generate the spike trains, we first used a simple approach that forced all the neurons to fire as
69 independent Poisson sources at a fixed and homogeneous firing rate of 5 Hz. To make the simulation
70 more biologically plausible, we pruned all spikes breaking a refractory period violation of 4 ms. The
71 resulting auto- and cross-correlograms for three sample units are shown in Figure 1F (auto-correlograms
72 are in green on the diagonal), while Figure 1G and H display the average (red line) and standard
73 deviation (grey shaded area) auto- and cross-correlation among all units, respectively. A sample
74 snippet of the generated traces from one recording is shown in Figure 1I, for a subset of 10 channels
75 out of 32. Due to the independence of the Poisson sources, both the average cross-correlograms
76 (Figure 1G) and auto-correlograms – outside the ± 4 ms used as refractory period – (Figure 1H) are
77 flat.

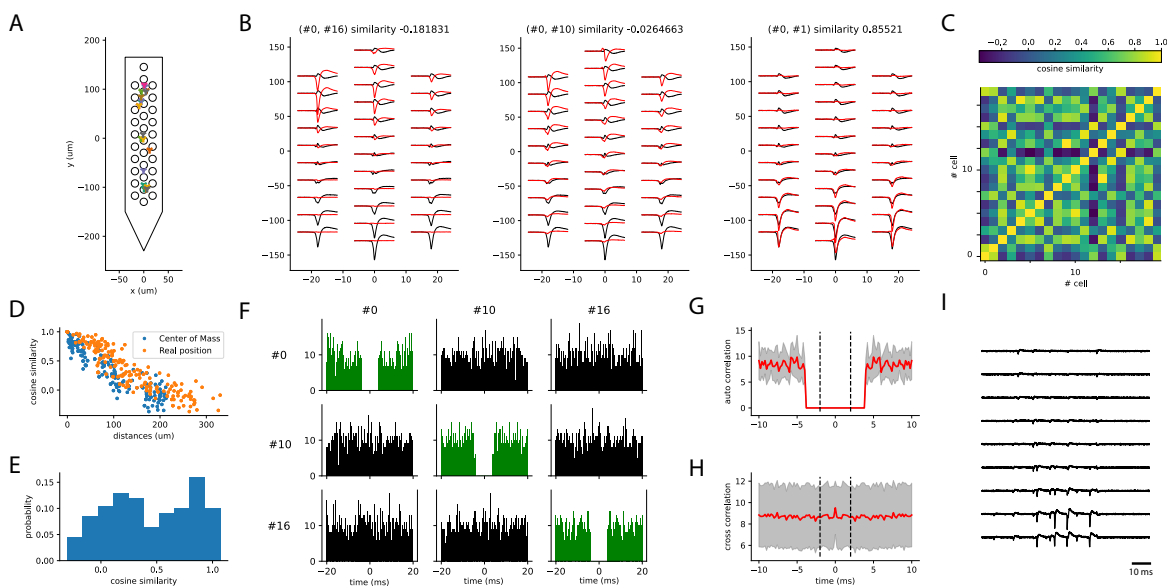


Figure 1: **Generation of the synthetic recordings.** **A)** 20 cells are randomly placed in front of a 32-channel NeuroNexus probe layout. The plot shows the location of each cell for one recording. **B)** Sample templates generated by neurons that are close too each other (#0 and #1) or far apart (#2). **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.

78 2.2 Global performance of the spike sorters

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

87 To check whether all sorters correctly *discovered* all templates, we computed the cosine similarity
 88 between the ground-truth templates from the simulations and the ones found by the sorters, comparing
 89 such a metric with the accuracy. As it can be seen in Figure 2C, all sorters are on average finding the
 90 correct templates, with the notable exception of YASS (in grey) and to some less extent HDsort (in
 91 red). Nevertheless, the overall accuracy of most of the spike sorters is relatively high (~ 0.95), except
 92 for HDsort and Herdingspikes which yield lower scores (Figure 2D). However, this averaged number
 93 does not tell us anything regarding the nature of these errors. While this error rate might seem low,
 94 it is likely that it is crucial, since it can potentially originate from the collisions, and thus from the
 95 correlations among neurons.

96 2.3 Spike sorting performance is affected by spike collisions

97 Using fully synthetic recordings with exhaustive ground truth, we can look at how good individual
 98 spike sorters perform specifically with respect to spatio-temporal collisions. To do so, we computed
 99 the *collision recall* (see Methods) as a function of the lag between two spikes, for a given pair of
 100 neurons. By averaging over multiple pairs of ground-truth neurons with similar template similarity
 101 (and over multiple recordings, see Methods), we can obtain a picture of how accurate the sorters
 102 are specifically with respect to the spike time lags and the similarities between templates. Figure 3
 103 displays the collision recall per sorter as a function of the lag (x-axis), colored by the similarity between
 104 templates. Each panel shows the performance of a different spike sorter. One can immediately see
 105 that only few sorters are able to accurately resolve lag correlations that are close to zero, even when
 106 templates are strongly orthogonal (low cosine similarity). This is the case for Kilosort 1 and 2, and for
 107 Spyking-circus, all of which use a template-matching procedure that should theoretically explain this

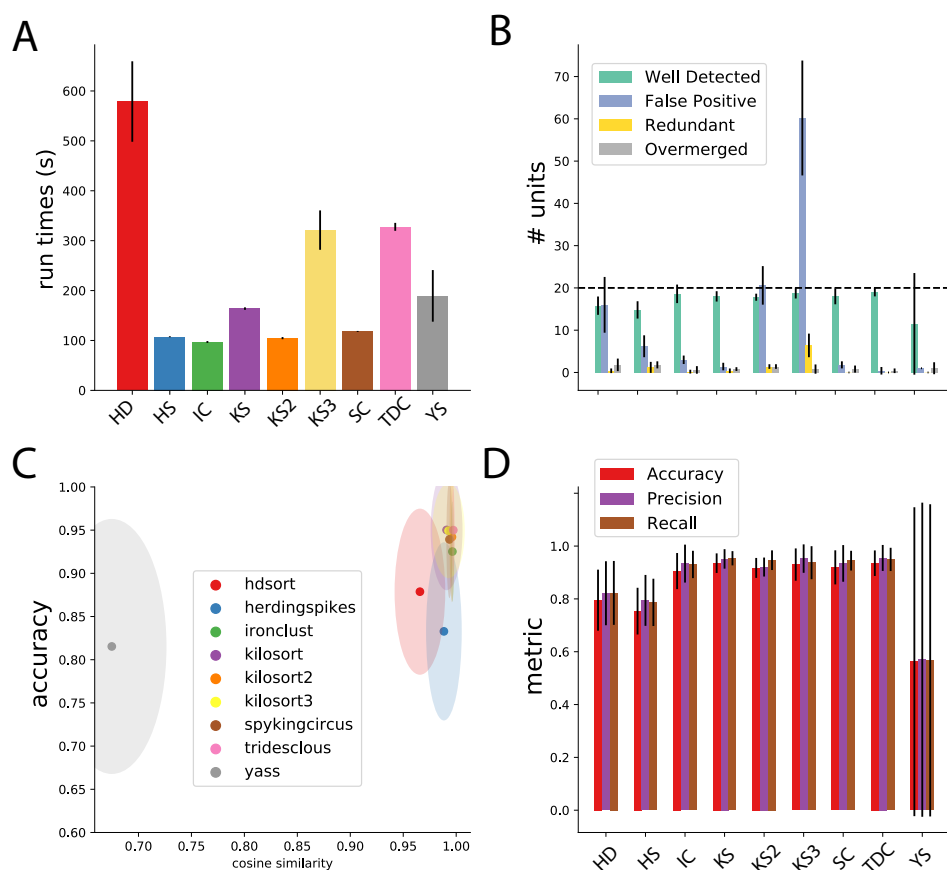


Figure 2: **Spike sorting performance.** **A)** Average run times over 5 different recordings (see Methods) for all the tested sorters. Errors bars indicate the standard deviation over multiple recordings. **B)** Average number of cells found by the sorters that are either well detected, redundant, overmerged 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.

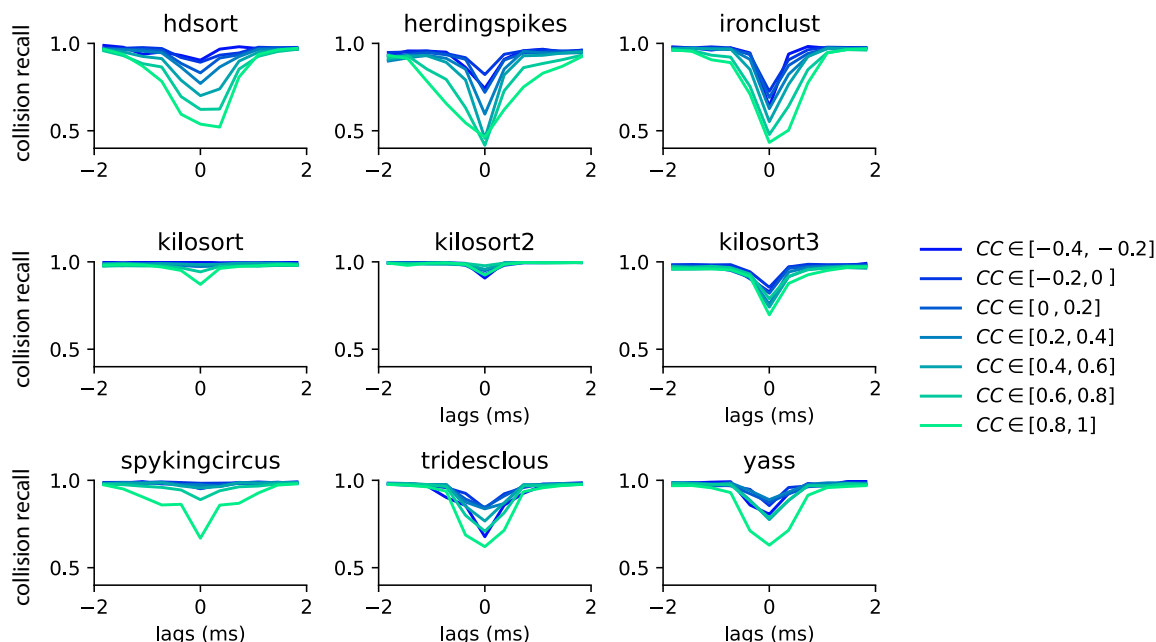


Figure 3: **Collision recall per sorter.** Error (quantified as the collision recall, see Methods) for various sorters and for all possible lags (between -2 and 2 ms), as function of the similarity between the pairs of templates (color code). All curves are averaged over multiple pairs and multiple recordings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

108 behavior. However, while performances are still good for Kilosort 1 and 2 even when the average cosine
 109 similarity between pairs is increased, they slightly degrade for Spyking-circus. Density-based sorters
 110 (HerdingSpikes and Ironclust), on the other hand, do not have a spike collision resolution strategy
 111 and this is reflected by their overall poorer performance. It is interesting to notice that Tridesclous,
 112 HDSort, YASS, and Kilsort 3, also using a template-matching based procedure to resolve the spikes,
 113 are not properly resolving the temporal correlations even for dissimilar templates. Different template-
 114 matching strategies are probably the cause of the differences among sorters. For example, HDSort
 115 and HerdingSpikes do not implement any strategy for spike collision resolution [9] and that is reflected
 116 in the quick degradation of performance as template similarity increases. KiloSort 1 and 2 used a
 117 GPU-based implementation of the k-SVD algorithm [1], used in matching learning as a dictionary
 118 learning algorithm for creating a dictionary for sparse representations. By doing so, it performs
 119 a reconstruction of the extra-cellular traces via orthogonal template matching pursuit, which is an
 120 enhancement of the greedy template matching pursuit (used in Spyking-circus and Tridesclous) more
 121 robust when templates are non-orthogonal. This might explain the boost in performance especially
 122 striking for templates with high similarity (*similarity* > 0.8).

123 2.4 Generation of controlled spike collision simulated data

124 The results shown in the previous section have been obtained only in a particular regime of activity,
 125 with all neurons firing independently as Poisson sources with an average firing rate of 5 Hz. However,
 126 neurons usually do not fire independently of each other, but rather have intrinsic correlations, also
 127 depending on different brain areas, brain states, and species. In addition, the average firing rates
 128 can also largely vary depending on brain areas. As an example, it is well known that Purkinje cells
 129 in the cerebellum have a very high firing rate [24] that networks tends to synchronize their activity

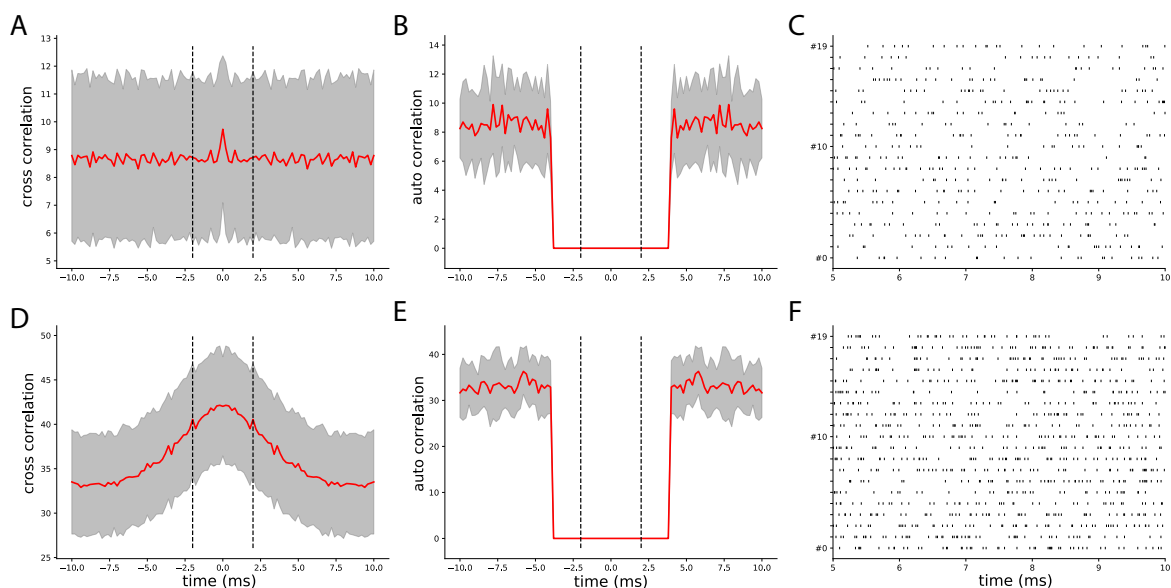


Figure 4: **Controlling spike trains correlations and firing rates.** **A)** Average cross-correlograms between all pairs of distinct neurons firing as independent Poisson sources at 5Hz (red curve, gray area represents the standard deviation) **B)** Same as **A**, but for auto-correlograms. **C)** Raster 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.

130 either in slow waves during sleep [25], or during pathological activity (such as epileptic seizures [26]).
 131 Therefore, assessing how performances may vary during different conditions is important to generalize
 132 our observations.

133 In order to study how spike sorting is affected by correlations and firing rates, we used a mixture
 134 procedure [5] that allowed us to control precisely the shape of the auto- and cross-correlograms for the
 135 injected spike trains. More precisely, we decided to explore in a systematic manner three rate levels
 136 (5, 10 and 15 Hz), and three correlation levels (0, 10, and 20 %). Note that the 5 Hz firing rate with
 137 0 % correlation corresponds to the scenario displayed in Figures 2-3.

138 Figure 4 shows the average of cross- and auto-correlograms and the spike trains of a recording where
 139 cells are firing as independent Poisson sources at 5 Hz in panels A-C (and thus with 0 % correlation,
 140 as shown by the flat average cross-correlograms in Figure 4A) and at 15 Hz with 20 % correlation
 141 (Figure 4D-F). Even though experimental recordings would contain a broader spectrum of firing rates
 142 and correlations, here we focus on assessing how different firing regimes affect spike sorting performance
 143 in a controlled setting. One would expect that the increased density of spikes (both in terms of firing
 144 rates and of synchrony) should degrade the performance of the spike sorters by affecting both the
 145 clustering step and the template-matching step, which in turn would degrade the resolution of spike
 146 collisions.

147 2.5 Do correlations and firing rates affect spike sorting of spike collisions?

148 To assess whether firing rate and spike train correlation affect spike sorting performance, we selected
 149 all unit pairs with a similarity greater than 0.5. We first averaged the recall curves for all template
 150 similarities (i.e. we averaged the curves with similarity greater than 0.5 shown in Figure 3). In
 151 Figure 5A we show the recall with respect to the spike lags averaged over all 9 configurations (3 firing
 152 rates x 3 correlations) for each sorter. The thick line represents the mean recall and the shaded area is
 153 the standard deviation over different rate-correlation configuration. All sorters, except YASS, appear

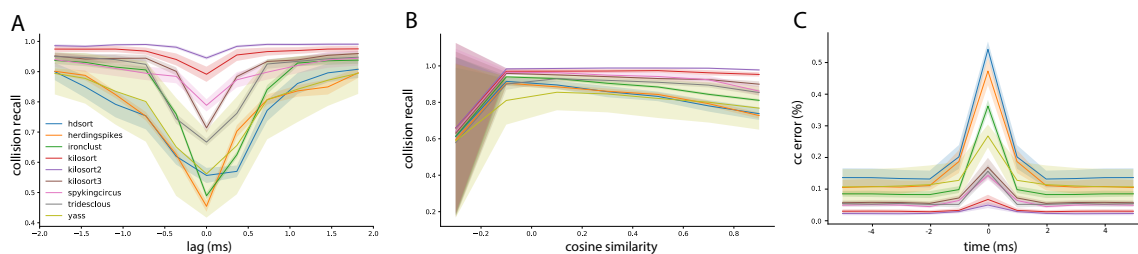


Figure 5: **Spike sorting performance for different conditions.** **A)** Average collision recall over the 9 conditions shown in Figure 5 - figure supplement 1 (3 firing rate levels and 3 correlation levels) as function of the lag between spikes, for pairs of cells with cosine similarity higher than 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

154 to have a very consistent curve (low standard deviation) over different configurations and do not seem
155 affected by changes in average firing rates and correlations in the spike trains. YASS' large
156 standard deviation can be explained by looking at individual recall curves at different rate-correlation regimes
157 (Figure 5 - figure supplement 1 - yellow lines): the spike sorting performance degrades with increasing
158 firing rates, but it does not seem to be strongly affected by increased correlation rates. However, we
159 should stress that since the collision recall is a relative measure, the same value for a larger number of
160 spikes (when firing rate is increased) means that overall, there are more misses for all sorters.

161 Similar considerations can be done by looking at the average recall with respect to template simi-
162 larity (Figure 5B). To construct this plots, we integrated the curves in Figure 3 over lags for different
163 cosine similarities. Also in this case, the curves appear consistent (low standard deviation) with the
164 exception of YASS, for which recall is reduced with increased firing rate regimes (Figure 5 - figure
165 supplement 2 - yellow lines). It is worth noticing that when the cosine similarity becomes negative, all
166 the sorters perform very poorly in properly resolving the overlaps. This could be explained by the fact
167 that when a pair of templates is anti-parallel (for example in the left panel of Figure 1A), a subset of
168 electrodes might show a negative signal for one template and a positive signal from the other (due to
169 return currents in the dendritic signals [11]). Effectively, when a spike collision between the two spikes
170 occur, this would lower the amplitude of the negative peak, which could reduce the detectability of
171 the spike.

172 The collision recall metric is mostly useful to obtain a quantitative insight on the behavior of the
173 spike sorting algorithms, but how do these errors transpose in practical situations? To assess this, we
174 measure the relative error (in percentage) between the ground-truth cross-correlograms and the ones
175 computed from the spike sorting outputs. We then averaged these error curves among all recordings
176 and experimental conditions (firing rates and synchrony levels). As shown in Figure 5, the error in
177 the estimated cross-correlogram can be as large as more than 50% for small lags, and for some spike
178 sorting algorithms such as HDsort, HerdingSpikes or IronClust. Moreover, it is also worth noticing
179 the baseline error rate is not the uniform across sorters. From this metric, we can again conclude that
180 template-matching based spike sorting algorithms such as KiloSort (1, 2, and 3), Spyking-circus or
181 Tridesclous are much better to resolve fine temporal correlations among neurons.

182 3 Discussion

183 In this study, we showed in a systematic and quantitative manner how spatio-temporal correlations
184 can be underestimated during spike sorting. Using synthetic datasets, we compared a large diversity of
185 modern spike sorters and showed how they behaved as function of the similarity between the templates

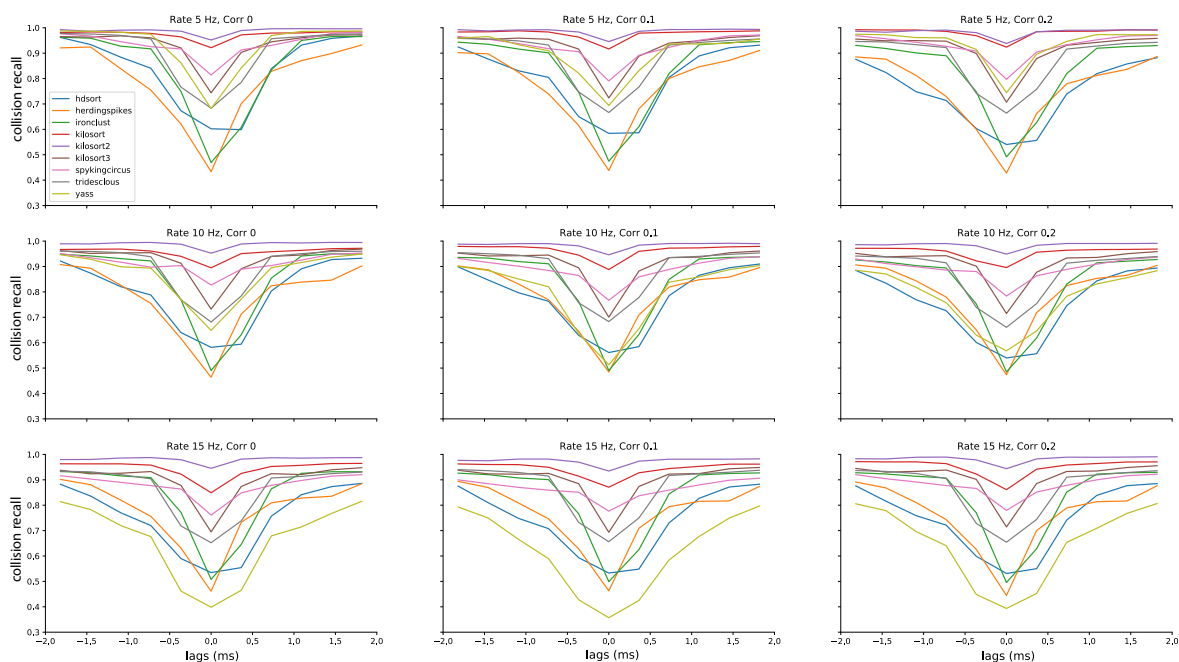


Figure 5 - figure supplement 1: Average performances of the spike sorters as function of the temporal lags. Each panel shows the average collision recall for template pairs with a similarity above 0.5 for a different condition, in terms of firing rate and correlation levels.

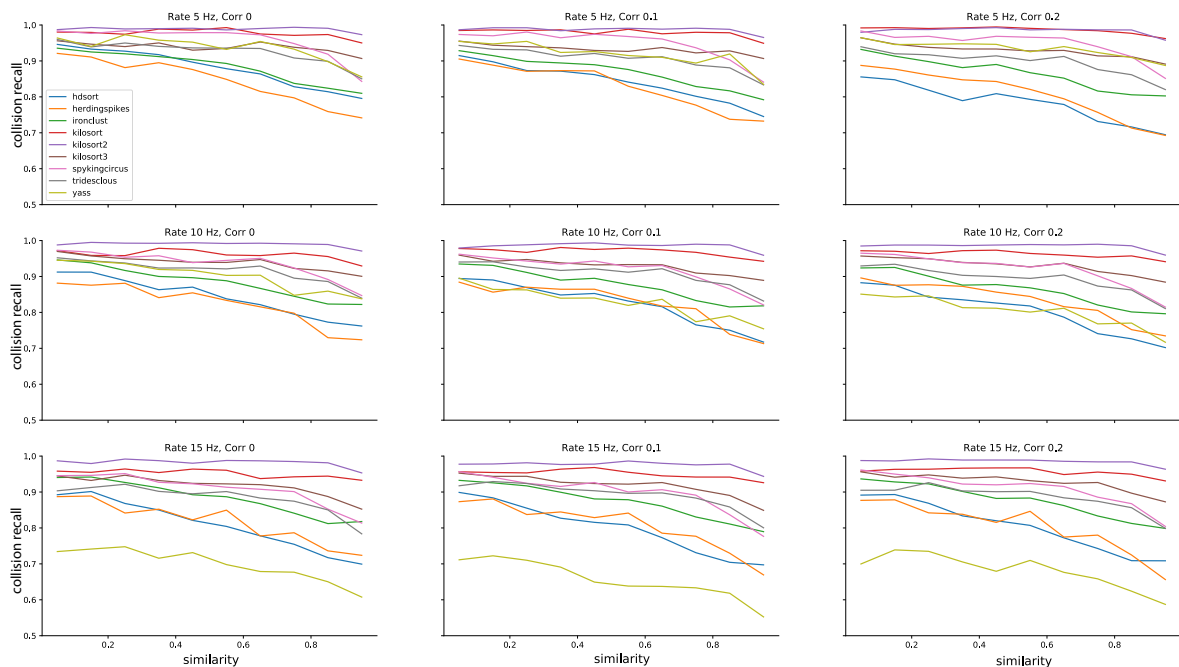


Figure 5 - figure supplement 2: Average performances of the spike sorters as function of the template similarity. Each panel shows the average collision recall over all lags in $[-2, 2]$ ms for a different condition, in terms of firing rate and correlation levels.

186 and the temporal lags between spikes. As expected, the closer the spikes are in time, the harder is
187 it, for all sorter, to properly resolve the overlaps. However, more interestingly, the more similar the
188 templates are, the higher the failures are. These failures are striking especially for spike sorters that
189 are not relying on template-matching based approaches (HerdingSpikes, Ironclust). For the ones using
190 a template-matching based approach (KiloSort, Spyking-circus, Tridesclous, HDSort), the problem
191 is less pronounced (with the exception of HDSort) but still present, and therefore this phenomenon
192 should be taken into account when making claims about the synchrony.

193 To our surprise, the global behavior of the spike sorters did not depend much on the overall
194 firing rate and/or the correlation levels. This allows us to generalize the findings and we think that
195 the quantitative results shown here could be translated to various *in vitro* or *in vivo* recordings from
196 different brain regions and species. As shown in Figure 5, while the variability over different conditions
197 is rather high for some algorithms, template-matching based algorithms tend to be rather robust and
198 overall better in resolving spike collisions. This is a very encouraging sign towards a unified and
199 reproducible automated solution for spike sorting [17, 7], agnostic of the recording conditions.

200 The results shown in the paper were obtained with purely artificial recordings, since we need
201 exhaustive information on the ground-truth spiking activity of all neurons to quantitatively compare
202 and benchmark different spike sorters. However, it would be interesting to generalize these observations
203 with real recordings, assuming one would have a proper ground truth at the population level. Indeed,
204 such a ground truth is needed to compute the *collision recall* and see how sorters behave as function
205 of lags and similarities between templates. To our knowledge, such a ground truth does not exist
206 [28, 19, 9]. While one could try to generate an “approximated” ground truth by combining the output
207 of several spike sorters with an *ensemble* spike sorting approach (as in [7]), the disagreements among
208 sorters are currently so high that this process is hard if not impossible, if one want to sample from a
209 large number of pairs.

210 While missing spikes for very dissimilar templates and small lags is problematic, the errors made for
211 very similar templates may be less frequent depending on the probe layout and neuronal preparation.
212 Indeed, such errors strongly depends on the distribution of template similarities between all pairs of
213 recorded cells, and this distribution might differ from recording to recording. For example, in the
214 retina [27] one would expect highly synchronous cells, of the same functional type, to be far apart from
215 each other because of an intrinsic tiling of the visual space. Such properties are unknown *in vivo* or
216 in cortical structures, but might bias the distribution of template similarities between nearby neurons,
217 and thus modify the estimation of collision recalls.

218 4 Methods

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

220 4.1 Simulated datasets

221 We used the MEArec simulator [6] to generate synthetic ground truth recordings. In brief, MEArec
222 uses biophysically detailed multicompartment models to simulate the extracellular action potentials, or
223 so called “templates”. For this study, we used 13 cell models from layer 5 of a juvenile rat somatosensory
224 cortex [18, 23]. Templates are then combined with spike trains and slightly modulated in amplitude
225 to add physiological variability. Additive uncorrelated Gaussian noise is finally added to the traces.
226 We generated simulated recordings with 20 neurons randomly positioned in front of the probe, a noise
227 level of 5 μ V and a sampling rate of 32 kHz. To obtain more robust results, we generated 5 recording
228 per conditions with various random seeds. The spike times were kept unchanged, but the positions
229 and the templates of the 20 neurons were changed in each of the individual recording. This allowed us
230 to populate the distribution of cosine similarities between pairs.

231 4.2 Generating spike trains with controlled correlations

232 To generate the recordings with various firing rates and correlations levels, we used the mixture process
233 method described in [5]. Since by default the method generates controlled cross-correlograms with a
234 decaying exponential profile, we modified it to generate cross-correlograms with a Gaussian profile, in
235 order to have more synchronous firing for small lags. The Gaussian profile can be seen in Figure 4D,
236 with a standard deviation $\sigma = 2.5\text{ms}$. By setting three different rate levels (5, 10 and 15 Hz) and three
237 different correlation levels (0, 10 and 20 %) this gave rise to 9 conditions, so to 45 recordings in total
238 (5 recordings per conditions, see above).

239 4.3 Template similarity

240 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
241 C the number of channels. After *flattening* the template by concatenating the signals from different
242 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
243 defined as follows:

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

244 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
245 is therefore bounded between -1 (templates are anti-parallel) and 1 (templates are parallel). A cosine
246 similarity of 0 means that the templates are orthogonal.

247 4.4 Spike sorters

248 All the spike sorters used in this study were run using the SpikeInterface framework [7], with default
249 parameters. The following are the exact versions that we used for the different spike sorters: Tridesclous
250 (1.6.4), Spyking-circus (1.0.9) [28], Herdingspikes (0.3.7) [13], Kilosort (v1, 2, or 3) [20], YASS (2.0)
251 [15], Ironclust (5.9.8) [8], HDSort (1.0.3) [9]. The desktop machine used has 36 Intel Xeon(R) Gold
252 5220 CPU @ 2.20GHz, 200Go of RAM and a Quadro RTX 5000 with 16Gb of RAM as a GPU.

253 4.5 Spike sorting comparison

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

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

$$\text{score}_{ij} = \frac{\#n_{\text{matches}}}{\#n_{i_{\text{gt}}} + \#n_{j_{\text{sorted}}} - \#n_{\text{matches}}}$$

258 where $\#n_{i_{\text{gt}}}$ and $\#n_{j_{\text{sorted}}}$ are the numbers of spikes in the i -th ground-truth spike train and the
259 j -th sorted spike trains, respectively. $\#n_{\text{matches}}$ is the number of spikes within 0.4 ms between the
260 two spike trains.

261 Once scores for all pairs are computed, an hungarian assignment is used to match ground-truth
262 units to sorted units [7]. All spikes from matched spike trains are then labeled as: true positive (TP),
263 if the spike is found both in the ground-truth and the sorted spike train; false positive (FP), if the
264 spike is found in the sorted spike train, but not in the ground-truth one; and false negative (FN), if
265 the spike is only found in the ground-truth spike train.

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

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

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

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

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

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

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

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

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

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

276 In order to generate the spike lag versus recall figures (e.g. Figures 3-5 - figure supplement 1) we
277 expanded the SpikeInterface software with several novel comparison methods and visualization widgets.
278 In particular, we extended the ground-truth comparison class to the `CollisionGTComparison`, which
279 computes performance metrics by spike lag. In addition to the agreement score computation and the
280 matching described in the previous paragraphs, this method first detects and flags all “synchronous
281 spike events” in the ground-truth spike trains. Two spikes from two separate units are considered to
282 be a “synchronous spike event” if their spike times occur within a time lag of 2 ms. The synchronous
283 events are then binned in 11 bins spanning the $[-2, 2]$ ms interval and the *collision recall* is computed
284 for each bin. With a similar principle, we implemented the `CorrelogramGTComparison` to compute
285 the lag-wise relative errors in cross-correlograms between ground-truth units and spike sorted units
286 (Figure 5C).

287 References

- 288 [1] M. Aharon, M. Elad, and A. Bruckstein. rm K-SVD: An Algorithm for Designing Overcomplete
289 Dictionaries for Sparse Representation. *IEEE Transactions on Signal Processing*, 54(11):4311–
290 4322, Nov. 2006.
- 291 [2] G. N. Angotzi, F. Boi, A. Lecomte, E. Miele, M. Malerba, S. Zucca, A. Casile, and L. Berdon-
292 dini. Sinaps: An implantable active pixel sensor cmos-probe for simultaneous large-scale neural
293 recordings. *Biosensors and Bioelectronics*, 126:355–364, 2019.
- 294 [3] B. B. Averbeck, P. E. Latham, and A. Pouget. Neural correlations, population coding and com-
295 putation. *Nature reviews neuroscience*, 7(5):358–366, 2006.
- 296 [4] L. Berdondini, K. Imfeld, A. Maccione, M. Tedesco, S. Neukom, M. Koudelka-Hep, and S. Marti-
297 noia. Active pixel sensor array for high spatio-temporal resolution electrophysiological recordings
298 from single cell to large scale neuronal networks. *Lab on a Chip*, 9(18):2644–2651, 2009.
- 299 [5] R. Brette. Generation of correlated spike trains. *Neural computation*, 21(1):188–215, 2009.
- 300 [6] A. P. Buccino and G. T. Einevoll. Mearec: a fast and customizable testbench simulator for
301 ground-truth extracellular spiking activity. *Neuroinformatics*, pages 1–20, 2020.
- 302 [7] A. P. Buccino, C. L. Hurwitz, S. Garcia, J. Magland, J. H. Siegle, R. Hurwitz, and M. H. Hennig.
303 Spikeinterface, a unified framework for spike sorting. *Elife*, 9:e61834, 2020.
- 304 [8] J. E. Chung, J. F. Magland, A. H. Barnett, et al. A fully automated approach to spike sorting.
305 *Neuron*, 95(6):1381–1394, 2017.
- 306 [9] R. Diggelmann, M. Fiscella, A. Hierlemann, and F. Franke. Automatic spike sorting for high-
307 density microelectrode arrays. *Journal of neurophysiology*, 120(6):3155–3171, 2018.
- 308 [10] U. Frey, U. Egert, F. Heer, S. Hafizovic, and A. Hierlemann. Microelectronic system for high-
309 resolution mapping of extracellular electric fields applied to brain slices. *Biosensors and Bioelec-
310 tronics*, 24(7):2191–2198, 2009.
- 311 [11] C. Gold, C. C. Girardin, K. A. Martin, and C. Koch. High-amplitude positive spikes recorded
312 extracellularly in cat visual cortex. *Journal of neurophysiology*, 102(6):3340–3351, 2009.
- 313 [12] M. H. Hennig, C. Hurwitz, and M. Sorbaro. Scaling spike detection and sorting for next-generation
314 electrophysiology. *In Vitro Neuronal Networks*, pages 171–184, 2019.
- 315 [13] G. Hilgen, M. Sorbaro, S. Pirmoradian, J.-O. Muthmann, I. E. Kepiro, S. Ullo, C. J. Ramirez,
316 A. P. Encinas, A. Maccione, L. Berdondini, et al. Unsupervised spike sorting for large-scale,
317 high-density multielectrode arrays. *Cell reports*, 18(10):2521–2532, 2017.
- 318 [14] J. J. Jun, N. A. Steinmetz, J. H. Siegle, D. J. Denman, M. Bauza, B. Barbarits, A. K. Lee, C. A.
319 Anastassiou, A. Andrei, Ç. Aydın, et al. Fully integrated silicon probes for high-density recording
320 of neural activity. *Nature*, 551(7679):232, 2017.
- 321 [15] J. Lee, C. Mitelut, H. Shokri, I. Kinsella, N. Dethe, S. Wu, K. Li, E. B. Reyes, D. Turcu, E. Batty,
322 et al. Yass: Yet another spike sorter applied to large-scale multi-electrode array recordings in
323 primate retina. *bioRxiv*, 2020.
- 324 [16] B. Lefebvre, P. Yger, and O. Marre. Recent progress in multi-electrode spike sorting methods.
325 *Journal of Physiology-Paris*, 110(4):327–335, 2016.

- 326 [17] J. Magland, J. J. Jun, E. Lovero, A. J. Morley, C. L. Hurwitz, A. P. Buccino, S. Garcia, and A. H.
327 Barnett. Spikeforest, reproducible web-facing ground-truth validation of automated neural spike
328 sorters. *Elife*, 9:e55167, 2020.
- 329 [18] H. Markram, E. Muller, S. Ramaswamy, et al. Reconstruction and simulation of neocortical
330 microcircuitry. *Cell*, 163(2):456–492, 2015.
- 331 [19] J. P. Neto, G. Lopes, J. Frazão, J. Nogueira, P. Lacerda, P. Baião, A. Aarts, A. Andrei, S. Musa,
332 E. Fortunato, et al. Validating silicon polytrodes with paired juxtacellular recordings: method
333 and dataset. *Journal of neurophysiology*, 116(2):892–903, 2016.
- 334 [20] M. Pachitariu, N. A. Steinmetz, S. N. Kadir, et al. Fast and accurate spike sorting of high-
335 channel count probes with kilosort. In *Advances in Neural Information Processing Systems*, pages
336 4448–4456, 2016.
- 337 [21] J. W. Pillow, J. Shlens, E. Chichilnisky, and E. P. Simoncelli. A model-based spike sorting
338 algorithm for removing correlation artifacts in multi-neuron recordings. *PloS one*, 8(5):e62123,
339 2013.
- 340 [22] R. Q. Quiroga, Z. Nadasdy, and Y. Ben-Shaul. Unsupervised spike detection and sorting with
341 wavelets and superparamagnetic clustering. *Neural computation*, 16(8):1661–1687, 2004.
- 342 [23] S. Ramaswamy, J. Courcol, M. Abdellah, et al. The neocortical microcircuit collaboration portal:
343 a resource for rat somatosensory cortex. *Front Neural Circuits*, 9, 2015.
- 344 [24] E. Sedaghat-Nejad, M. A. Fakharian, J. Pi, P. Hage, Y. Kojima, R. Soetedjo, S. Ohmae, J. F. Med-
345 ina, and R. Shadmehr. P-sort: an open-source software for cerebellar neurophysiology. *bioRxiv*,
346 2021.
- 347 [25] M. Steriade. Slow-wave sleep: serotonin, neuronal plasticity, and seizures. *Arch Ital Biol*,
348 142(4):359–367, Jul 2004.
- 349 [26] W. Truccolo, J. A. Donoghue, L. R. Hochberg, E. N. Eskandar, J. R. Madsen, W. S. Anderson,
350 E. N. Brown, E. Halgren, and S. S. Cash. Single-neuron dynamics in human focal epilepsy. *Nat*
351 *Neurosci*, 14(5):635–641, May 2011.
- 352 [27] H. Wässle. Parallel processing in the mammalian retina. *Nat Rev Neurosci*, 5(10):747–757, Oct
353 2004.
- 354 [28] P. Yger, G. L. Spampinato, E. Esposito, B. Lefebvre, S. Deny, C. Gardella, M. Stimberg, F. Jetter,
355 G. Zeck, S. Picaud, et al. A spike sorting toolbox for up to thousands of electrodes validated with
356 ground truth recordings in vitro and in vivo. *Elife*, 7:e34518, 2018.