

# Multi-day Neuron Tracking in High Density Electrophysiology Recordings using EMD

Augustine(Xiaoran) Yuan<sup>1,2</sup>, Jennifer Colonell<sup>1</sup>, Anna Lebedeva<sup>3</sup>, Michael Okun<sup>4</sup>, Adam S. Charles<sup>2\*</sup>, Timothy D. Harris<sup>1,2\*</sup>

\*For correspondence:

[adamsc@jhu.edu](mailto:adamsc@jhu.edu) (ASC);

[harrist@janelia.hhmi.org](mailto:harrist@janelia.hhmi.org) (TDH)

<sup>1</sup>Janelia Research Campus, Howard Hughes Medical Institute, USA; <sup>2</sup>Department of Biomedical Engineering, Center for Imaging Science Institute, Kavli Neuroscience Discovery Institute, Johns Hopkins University, USA; <sup>3</sup>Sainsbury Wellcome Centre, University College London, UK; <sup>4</sup>Department of Psychology and Neuroscience Institute, University of Sheffield, UK

**Abstract** Accurate tracking of the same neurons across multiple days is crucial for studying changes in neuronal activity during learning and adaptation. Advances in high density extracellular electrophysiology recording probes, such as Neuropixels, provide a promising avenue to accomplish this goal. Identifying the same neurons in multiple recordings is, however, complicated by non-rigid movement of the tissue relative to the recording sites (drift) and loss of signal from some neurons. Here we propose a neuron tracking method that can identify the same cells independent of firing statistics, that are used by most existing methods. Our method is based on between-day non-rigid alignment of spike sorted clusters. We verified the same cell identity in mice using measured visual receptive fields. This method succeeds on datasets separated from one to 47 days, with an 84% average recovery rate.

## 1 Introduction

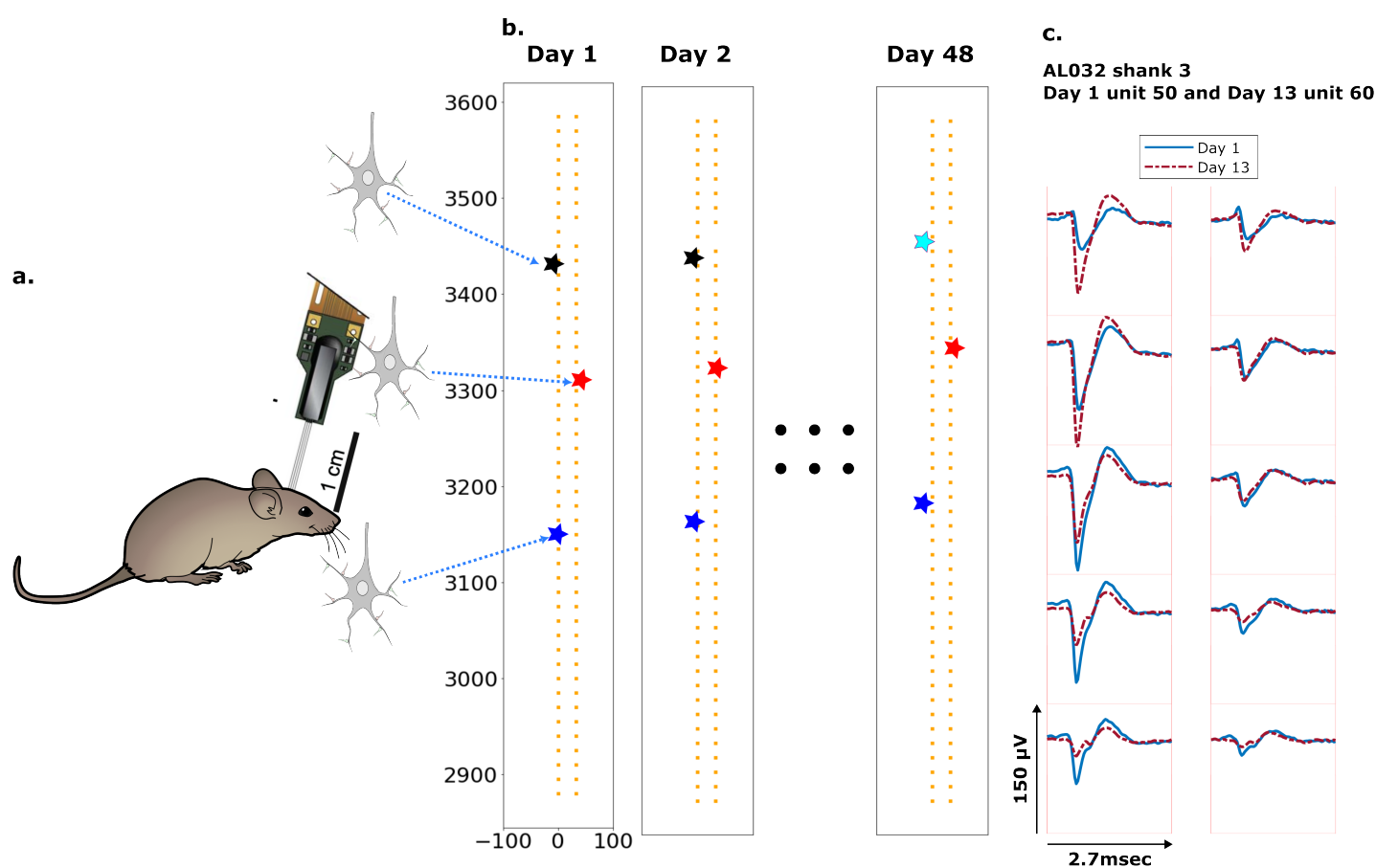
The ability to longitudinally track neural activity is crucial to understanding central capabilities and changes of neural circuits that operate on long time-scales, such as learning and plasticity,<sup>1-4</sup> motor stability,<sup>1,5,6</sup> etc. We seek to develop a method capable of tracking single units regardless of changes in functional responses for the duration of an experiment spanning one to two months.

High-density multi-channel extracellular electrophysiology (ephys) recording devices enable chronic recordings over large areas for days-to-months.<sup>7</sup> Such chronic recordings make possible experiments targeted at improving our understanding of neural computation and underlying mechanisms. Examples include perceptual decision making, exploration and navigation.<sup>8-13</sup> Electrode arrays with hundreds to thousands of sites, for example Neuropixels, are now used extensively to record the neural activity of large populations stably and with high spatio-temporal resolution, capturing hundreds of neurons with single neuron resolution.<sup>9,10</sup> Moreover, ephys retains the higher time resolution needed for single spike identification, as compared with calcium imaging that provides more spatial cues with which to track neurons over days.

The first step in analyzing ephys data is to extract single neuron signals from the recorded voltage traces, i.e., spike sorting. Spike sorting identifies individual neurons by grouping detected action potentials using waveform profiles and amplitudes. Specific algorithms include principal

40 components based methods,<sup>14, 15</sup> and template matching methods, for example, Kilosort.<sup>9, 11, 16, 17</sup>  
 41 Due to the high dimensional nature of the data, spike sorting is often computationally intensive  
 42 on large data sets (10's to 100's of GB) and optimized to run on single sessions. Thus processing  
 43 multiple sessions has received minimal attention, and the challenges therein remain largely unad-  
 44 dressed.

45 One major challenge in reliably tracking neurons is the potential for changes in the neuron  
 46 population recorded (**Figure 1a** and **Figure 1b**). In particular, since the probe is attached to the  
 47 skull, brain tissue can move relative to the probe, e.g. during licking, and drift can accumulate over  
 48 time.<sup>18</sup> Kilosort 2.5 corrects drift within a single recording by inferring tissue motion from con-  
 49 tinuous changes in spiking activity and interpolating the data to account for that motion.<sup>7</sup> Larger  
 50 between-recording drift occurs for sessions on different days, and can 1) change the size and loca-  
 51 tion of spike waveforms along the probe,<sup>19</sup> 2) lose neurons that move out of range, and 3) gain new  
 52 neurons that move into recording range. Thus clusters can change firing pattern characteristics or  
 53 completely appear/disappear. As a result the specific firing patterns classified as unit clusters may  
 54 appear and disappear in different recordings.<sup>9, 20-22</sup> Another challenge is that popular template-  
 55 matching-based spike sorting methods usually involve some randomness in template initializa-  
 56 tion.<sup>16, 23, 24</sup> As a result, action potentials can be assigned into clusters differently, and clusters can  
 57 be merged or separated differently across runs.



**Fig. 1: Schematic depiction of drift:** a. Mice were implanted with a 4-shank Neuropixels 2.0 probe in visual cortex area V1. b. Each colored star represents the location of a unit recorded on the probe. In this hypothetical case, the same color indicates unit correspondence across days. The black unit is missing on day 48, while the turquoise star is an example of a new unit. Tracking aims to correctly match the red and blue units across all datasets and determine that the black unit is undetected on day 48. c. Two example spatial-temporal waveforms of units recorded in two datasets that likely represent the same neuron, based on similar visual responses. Each trace is the average waveform on one channel across 2.7 milliseconds. The blue traces are waveforms on the peak channel and 9 nearby channels (two rows above, two rows below, and one in the same row) from the first dataset (Day 1). The red traces, similarly selected, are from the second dataset. Waveforms are aligned at the electrodes with peak amplitude, different on the two days.

58 Previous neuron tracking methods are frequently based on waveform and firing statistics, e.g.,  
 59 firing rate similarity,<sup>25</sup> action potential shape correlation and inter-spike interval histogram (ISI)  
 60 shape.<sup>26</sup> When neuronal representations change, e.g., during learning<sup>1-3</sup> or representational drift,<sup>27</sup>  
 61 neural activity statistics became less reliable. In this work, we take advantage of the rich spatial-  
 62 temporal information in the multi-channel recordings, matching units based on the estimated neu-  
 63 ron locations and unit waveforms,<sup>28</sup> instead of firing patterns.

64 As an alternative method, Steinmetz et al.<sup>7</sup> concatenated pairs of datasets after low resolution  
 65 alignment, awkward for more than 2 datasets. We report here a more flexible, expandable and  
 66 robust tracking method that can track neurons effectively and efficiently across any number of  
 67 sessions.

## 2 Results

### 2.1 Procedure

Our datasets consist of multiple recordings taken from three mice (**Figure 7a**) over 2 months. The time gap between two recordings ranges from two to 25 days. Each dataset is spike-sorted individually with a standard Kilosort 2.5 pipeline. The sorting results, including unit assignment, spike times, etc. are used as input for our method (post-processed using ecephys spike sorting pipeline<sup>29</sup>) (Sec. 4.3). To ensure the sorting results are unbiased, we performed no manual curation. As the clusters returned by Kilosort can vary in quality, we only considered the subset of units labeled as ‘good’ by Kilosort, here referred to as KSgood units (Sec. 4.4). KSgood units are mainly determined by the amount of inter-spike-interval violations and are believed to represent a single unit.<sup>16</sup>

Our overall strategy is to run spike-sorting once per session, and then to generate a unit-by-unit assignment between pairs of datasets. When tracking units across more than two sessions, two strategies are possible: match all ensuing sessions to a single session (e.g., the first session) (Sec. 2.2 and Sec. 4.2), or match consecutive pairs of sessions and then trace matched units through all sessions (Sec. 2.4).

We refer to the subset of KSgood units with strong and distinguishable visual responses in both datasets of a comparison as reference units (See Sec. 4.4 for details). Similar to Steinmetz et al.<sup>7</sup> we validated our unit matching of reference units using visual receptive field similarity. Finally, we showed that trackable units with strong visual responses are qualitatively similar to those without (**Figure 5-supplement Figure 1 to Figure 5**).

To provide registration between pairs of recordings, we used the Earth Mover’s Distance (EMD).<sup>30,31</sup> We use a feature space consisting of a geometric distance space and a waveform similarity space, to address both rigid and non-rigid neuron motion. The EMD finds matches between objects in the two distributions by minimizing the overall distances between the established matches (Sec. 4.1.1).

We use EMD in two stages: rigid drift correction and unit assignment. Importantly, the EMD distance incorporates two parameters crucial for matching units: location-based physical distance and a waveform distance metric that characterizes similarity of waveforms (Sec. 4.1.2). The EMD distance matrix is constructed with a weighted combination of the two (details in Sec. 4), i.e. a distance between two units  $d_{ik}$  is given by  $d_{ik} = d_{location_{ik}} + \omega * d_{waveform_{ik}}$  (**Figure 2a**). The first EMD stage estimates the homogeneous vertical movement of the entire population of KSgood units (**Figure 2b**). This movement estimate is used to correct the between-session rigid drift in unit locations. The rigid drift estimation procedure is illustrated in figure 2b. Post drift correction, a unit’s true match will be close in both physical distance and waveform distance. Drift-corrected units were then matched at the second EMD stage. The EMD distance between assigned units can be thought of as the local non-rigid drift combined with the waveform distortion resulting from drift. We test the accuracy of the matching by comparing with reference unit assignments based on visual receptive fields (Sec. 4.4).

For each unit, the location is determined by fitting the peak to peak amplitudes on the 10 sites nearest the site with peak signal, based on the triangulation method in Boussard, et al.<sup>32</sup> (Sec. 4.1.2). The waveform distance is an L2 norm between two spatial-temporal waveforms that spans 22 channels and 2.7 msec (Sec. 4.1.2). Physical unit distances provide a way to maintain the internal structure and relations between units in the EMD. Waveform similarity metrics will distinguish units in the local neighborhood and likely reduce the effect of new and missing units.

We analyzed the match assignment results in two ways. First, we compared all subsequent datasets to dataset 1 using recovery rate and accuracy. We define recovery rate  $R_{rec}$  as the fraction of unit assignments by our method that are the same as reference unit assignments established using visual responses (Sec. 4.4).



$$P(EMD | ref) = \frac{P(EMD \cap ref)}{P(ref)} = \frac{N_{EMD \cap ref}}{N_{ref}} \quad (1)$$

Since the EMD forces all units from the dataset with fewer neurons to have an assigned match, we use vertical z-distance to threshold out the biologically-impossible unit assignments. We then calculated the accuracy  $R_{acc}$ , i.e. the fraction of EMD unit assignments within the z-distance threshold which agree with the reference assignments.

$$P((EMD | ref) \cap threshold) = \frac{P((EMD \cap ref) | threshold)}{P(ref | threshold)} \quad (2)$$

We also retrieved non-reference units, i.e. matched units without receptive field information but whose z-distance is smaller than the threshold.

Second, we tracked units between consecutive datasets and summarized and analyzed the waveforms, unit locations, firing rates and visual responses (see **Figure 5**-supplement **Figure 1** to **Figure 5** for details) of all tracked chains, i.e. units which can be tracked across at least three consecutive datasets.

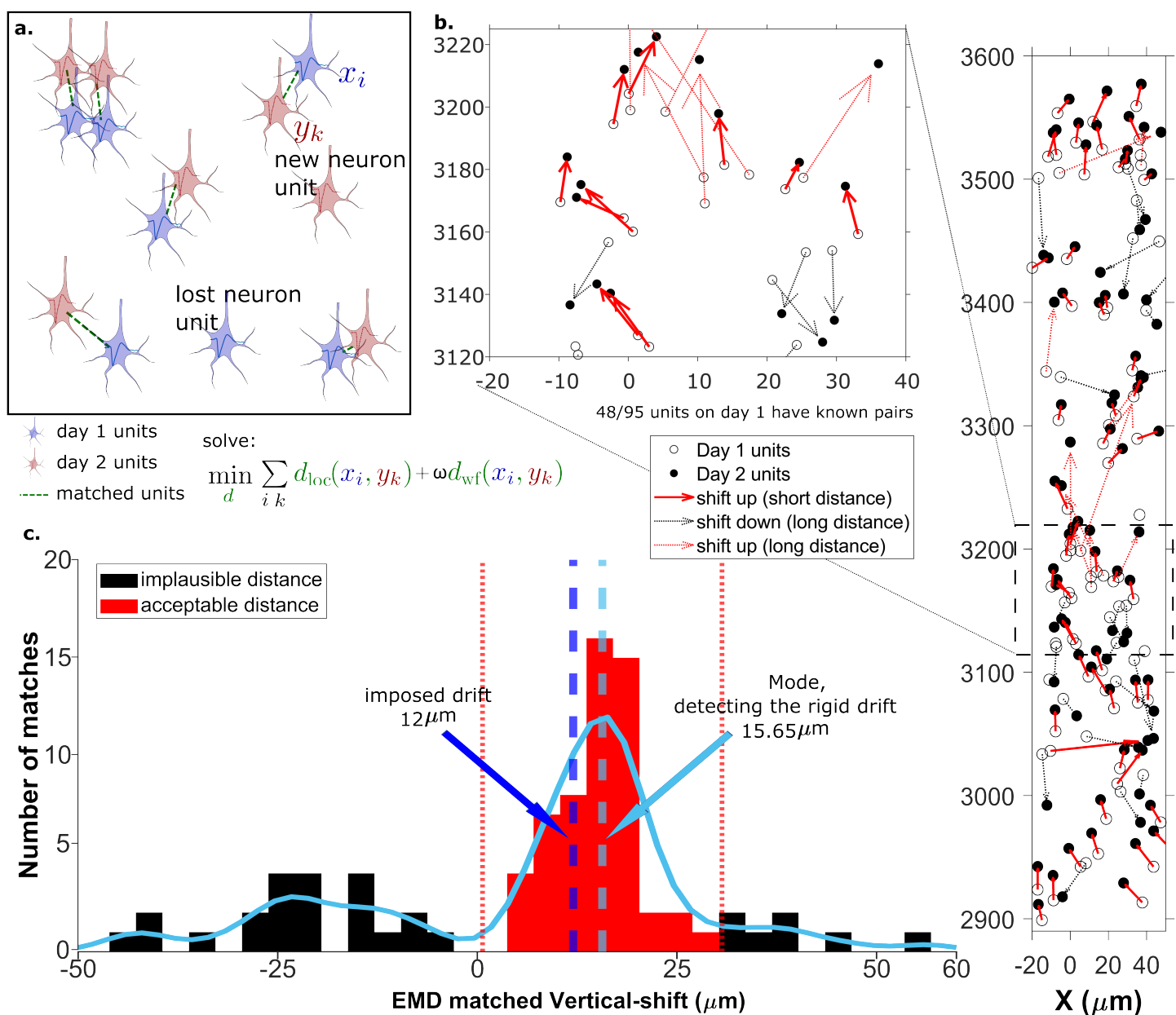
## 2.2 Measuring rigid drift using the EMD

Drift happens mostly along the direction of probe insertion (vertical or z direction). We want to estimate the amount of vertical drift under the assumption that part of the drift is rigid; this is likely a good assumption given the small ( $\approx 720\mu m$ ) z-range of these recordings. The EMD allows us to extract the homogeneous (rigid) movement of matched units. For ideal datasets with a few units consistently detected across days, this problem is relatively simple (**Figure 2a**). In the real data analyzed here, we find that only  $\approx 60\%$  of units are detected across pairs of days, so the rigid motion of the real pairs must be detected against a background of units with no true match. These units with no real match will have z-shifts far from the consensus z-shift of the paired units (**Figure 2c**).

In **Figure 2** the EMD match of units from the first dataset (**Figure 2b**, open circles) to the dataset recorded the next day (**Figure 2b**, closed circles) is indicated by the arrows between them. To demonstrate detection of significant drift, we added a 12 micron upward drift to the z-coordinate of the units from the second day. The first stage of the EMD is used to find matches using the combined distance metric as described in section 4.1.2. We used a kernel fit to the distribution of z-distances of all matched units to find the mode (Mode =  $15.65\mu m$ ); this most probable distance is the estimate of the drift (**Figure 2c**). It is close to the actual imposed drift ( $d_i = 12\mu m$ ).

As the EMD is an optimization algorithm with no biological constraints, it assigns matches to all units in the smaller dataset regardless of biophysical plausibility. As a result, some of the assigned matches may have unrealistically long distances. A distance threshold is therefore required to select correct pairs. For the illustration in **Figure 2**, the threshold is set to  $15\mu m$ , which is chosen to be larger than most of the z-shifts observed in our experimental data. The threshold value will be refined later by distribution fitting (**Figure 4**). In **Figure 2** all of the sub-threshold (short) distances belong to upward pairs (**Figure 2b** and **c**, red solid arrows), showing that the EMD can detect the homogeneous movement direction and the amount of imposed drift.

When determining matched reference units from visual response data, we require that units be spatially nearby (within  $30\mu m$ ) as well as having similar visual responses. After correcting for drift, we find that we recover more reference units (**Figure 2**-supplement **Figure 1**), indicating improved spatial match of the two ensembles. This improved recovery provides further evidence of the success of the drift correction.

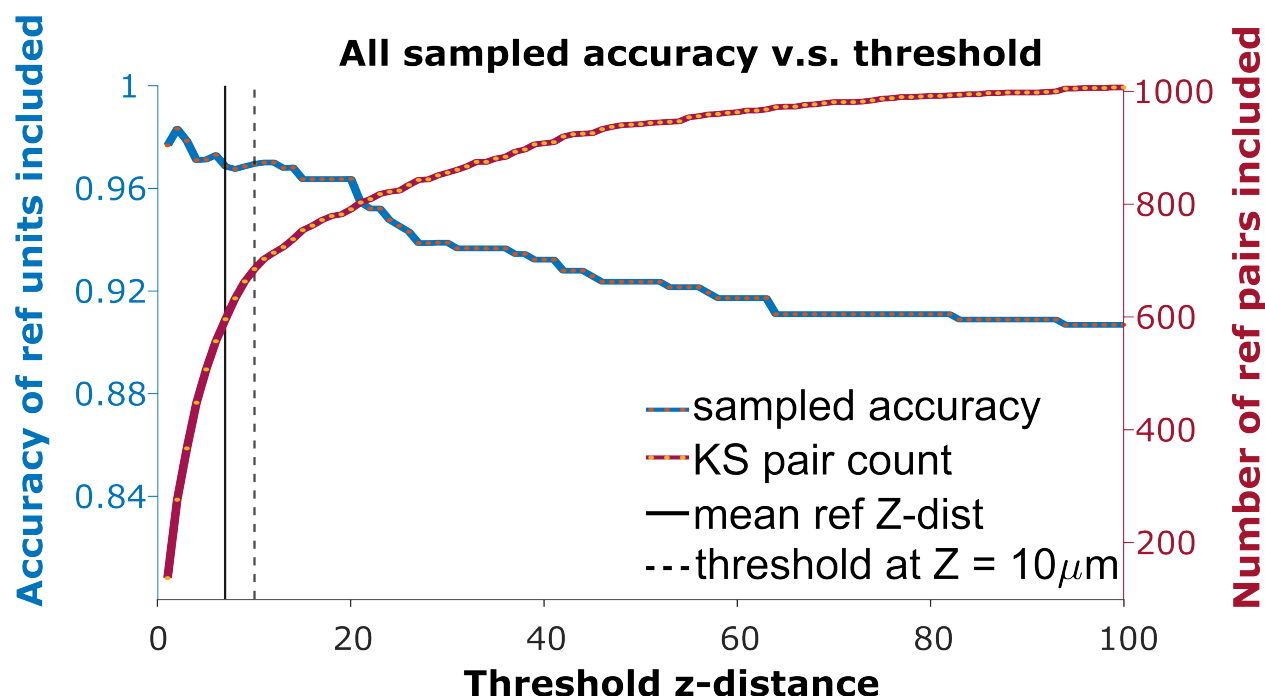


**Fig. 2: The EMD can detect the displacement of single units:** a. Schematic of EMD unit matching. Each blue unit in day 1 is matched to a red unit in day 2. Dashed lines indicate the matches to be found by minimizing the weighted sum of physical and waveform distances. b. Open and filled circles show positions of units in days 1 and 2, respectively. Arrows indicate matching using EMD. The arrow color represents the match direction; upward matches found with the EMD are in red and downward in black. Solid lines indicate a z-match distance within  $15\mu\text{m}$ , while a dashed line indicates a z distance  $> 15\mu\text{m}$ . Expanded view shows probe area from 3120 to 3220  $\mu\text{m}$ . c. Histogram of z-distances of matches (black and red bars) and kernel fit (light blue solid curve). The light blue dashed line shows the mode ( $d_m = 15.65\mu\text{m}$ ). The dark blue dashed line shows the imposed drift ( $d_i = 12\mu\text{m}$ ). The red region shows the matches within  $15\mu\text{m}$  of the mode. The EMD needs to detect the homogeneous movement against the background, i.e. units in the black region that are unlikely to be the real matches due to biological constraints.

### 2.3 A vertical distance threshold is necessary for accurate tracking

To detect the homogeneous z-shift of correct matches against the background of units without true matches, it is necessary to apply a threshold on the z-shift. When tracking units after shift cor-

158 rection, a vertical distance threshold is again required to determine which matches are reasonable  
159 in consideration of biological plausibility. The Receiver Operator Characteristic (ROC) curve in **Fig-**  
160 **ure 3** shows the fraction of reference units matched correctly and the number of reference pairs  
161 retained as a function of z-distance threshold. We want to determine the threshold that maximizes  
162 the overall accuracy in the reference units (**Figure 3**, blue curve) while including as many reference  
163 units as possible (**Figure 3**, red curve).



**Fig. 3: The ROC curve of matching accuracy vs. distance.** The blue curve shows the accuracy for reference units. The red line indicates the number of reference units included. The solid vertical line indicates the average z distance across all reference pairs in all animals ( $z = 6.96\mu\text{m}$ ). The dashed vertical black line indicates a z-distance threshold at  $z = 10\mu\text{m}$ .

164 Since reference units only account for 29% of KSgood units (units with few inter-spike-interval  
165 violations that are believed to represent a single unit), and the majority of KSgood units did not  
166 show a distinguishable visual response, we need to understand how representative the reference  
167 units are of all KSgood units.

168 We found the distribution of z-distances of reference pairs is different from the distribution  
169 of all KSgood units (**Figure 4a**, top and middle panel). While both distributions may be fit to an  
170 exponential decay, the best fit decay constant is significantly different (Kolmogorov-Smirnov test,  
171 reject  $H_0$ ,  $p = 5.5 \times 10^{-31}$ ). Therefore, the accuracy predicted by the ROC of reference pairs in Figure  
172 3 will not apply to the set of all KSgood pairs. The difference in distribution is likely due to the  
173 reference units being a special subset of KSgood units in which units are guaranteed to be found  
174 in both datasets, whereas the remaining units may not have a real match in the second dataset. To  
175 estimate the ROC curve for the set of all KSgood units, we must estimate the z-distance distribution

176 for a mixture of correct and incorrect pairs.

177 We assume that the distribution of z-distances  $P(\Delta)$  for reference units is the conditional prob-  
178 ability  $P(\Delta | H)$ ; that is, we assume all reference units are true hits. The distribution of z-distances  
179 for all KSgood units  $P(\Delta)$  includes both hits and false positives. The distance distribution of false  
180 positives is the difference between the two.

181 A Monte Carlo simulation determined that the best model for fitting the z-distance distribution  
182 of reference units  $P(\Delta | H)$  is a folded Gaussian distribution (**Figure 4a**, middle panel) and an  
183 exponential distribution for false positive units (see **Figure 4-supplement Figure 1**). The KSgood  
184 distribution is a weighted combination of the folded Gaussian and an exponential:

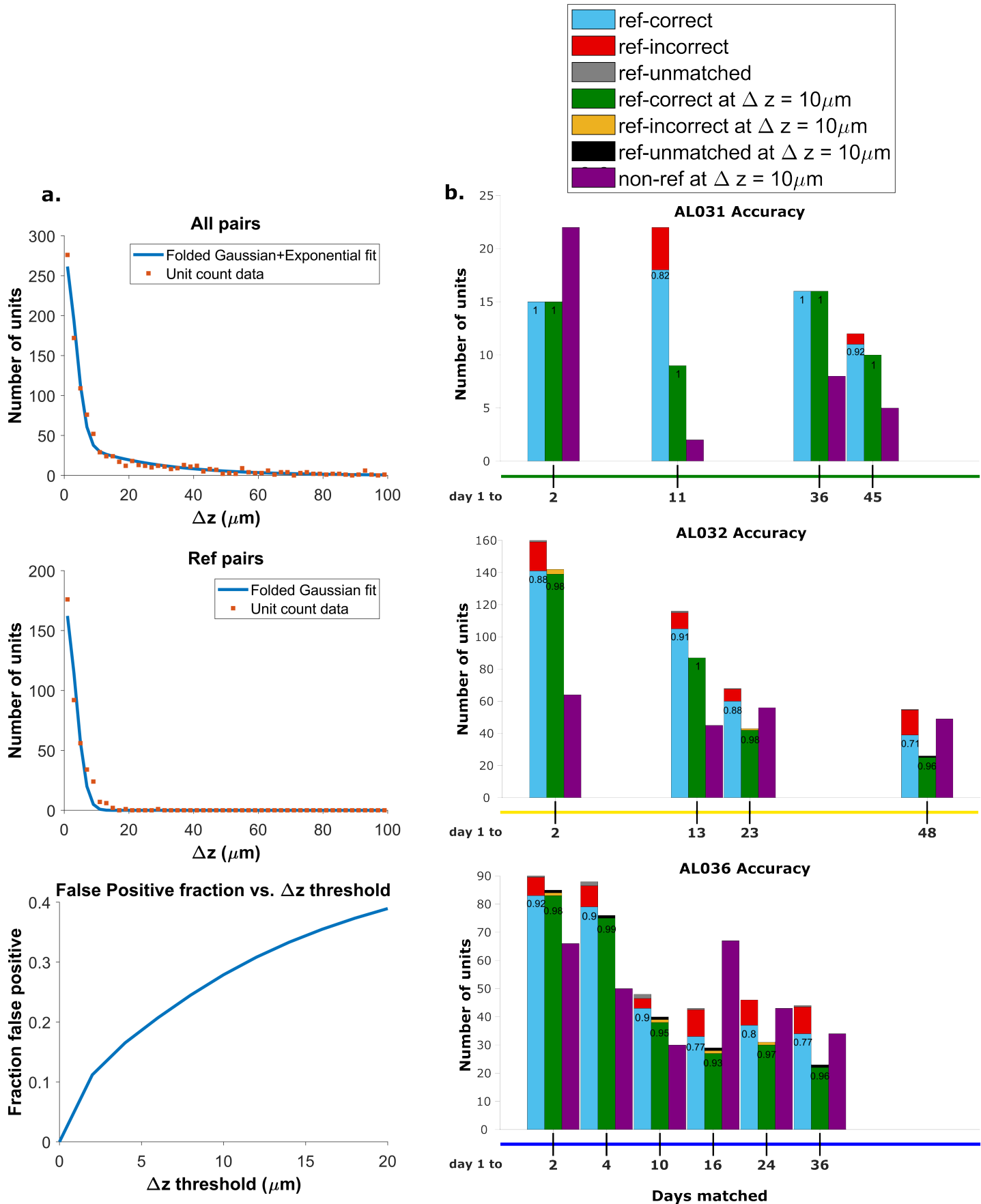
$$P(AllUnits) = f * P(FoldedGaussian) + (1 - f) * P(Exponential) \quad (3)$$

185 We fit the KSgood distribution to **Equation 3** to extract the individual distribution parameters and  
186 the fraction of true hits ( $f$ ). The full distribution can then be integrated up to any given z-threshold  
187 value to calculate the false positive rate. (**Figure 4a**, bottom panel, see **Figure 4-supplement Figure 2**  
188 for details).

189 Based on the the estimated false positive rate (**Figure 4a**, bottom panel), we used a threshold  
190 of  $10\mu m$  (**Figure 3**, black dotted line) to obtain at least 70% accuracy in the KSgood units. We used  
191 the same threshold to calculate the number of matched reference units and the corresponding  
192 reference unit accuracy (**Figure 4b**, green bars).

193 Note that this threshold eliminates most of the known false positive matches of reference pairs  
194 (**Figure 4b**, red fraction) at the cost of recovering fewer correct pairs (**Figure 4b**, green bars). The re-  
195 covery rate varies from day to day; datasets separated by longer times tend to have higher tracking  
196 uncertainty (**Figure 4-supplement Figure 3**).

197 In addition to the units with visual response data, we can track units which have no significant  
198 visual response (**Figure 4b**, purple bars). All comparisons are between subsequent datasets and  
199 the day 1 dataset.

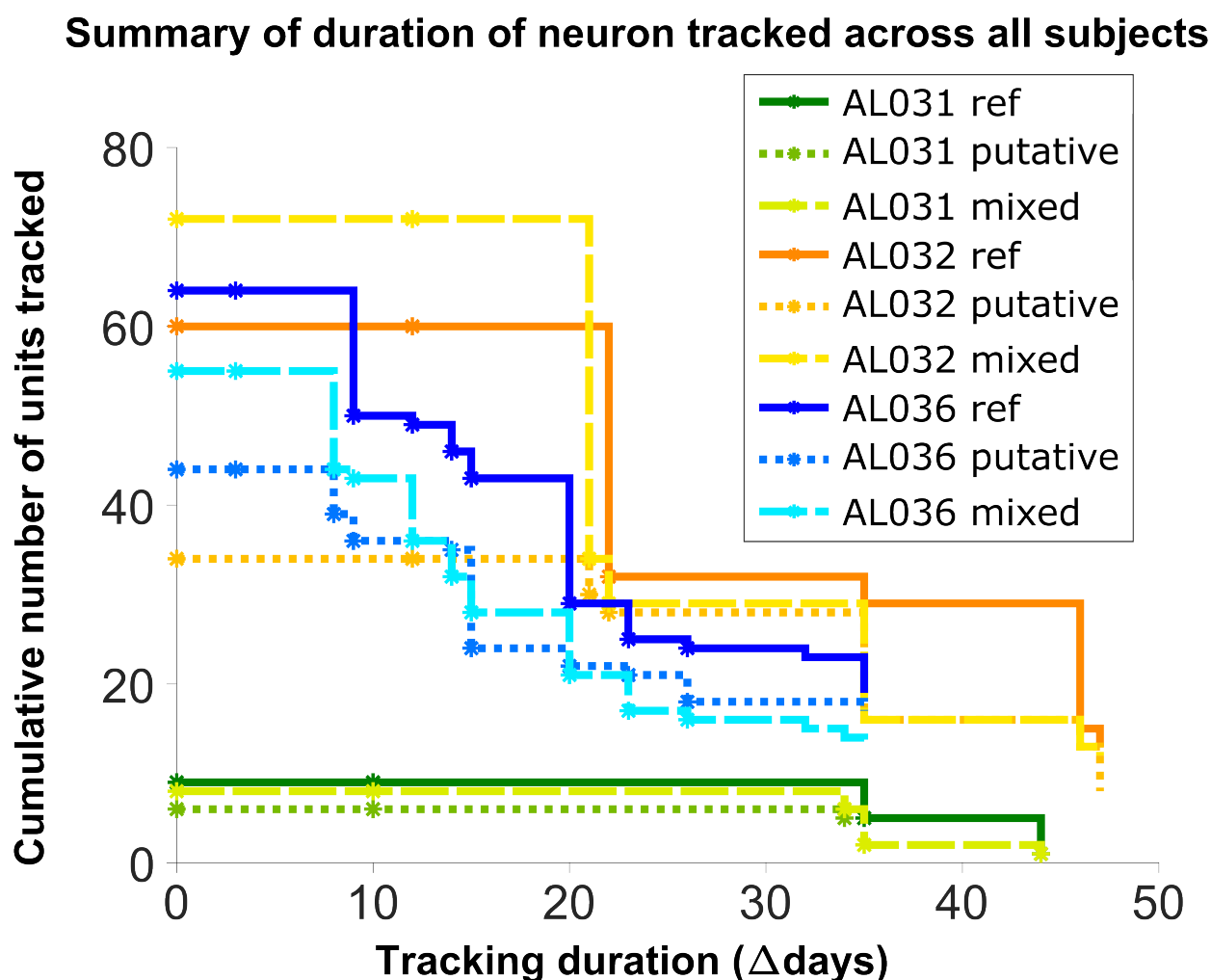


**Fig. 4: Recovery rate, accuracy and putative pairs:** a. The histogram distribution fit for all KS-good units (top) and reference units alone (middle). False positives for reference units are defined as units matched by EMD but not matched when using receptive fields. The false positive fraction for the set of all KSgood units is obtained by integration.  $z = 10\mu m$  threshold has a false positive rate = 27% for KSgood units. b. Light blue bars represent the number of reference units successfully recovered using only unit location and waveform. The numbers on the bars are the recovery rate of each dataset, and the red portion indicates incorrect matches. Incorrect matches are cases where units with a known match from receptive field data are paired with a different unit by EMD; these errors are false positives. The green bars show matching accuracy for the set of pairs with  $z$ -distance less than the  $10\mu m$  threshold. The orange portion indicates incorrect matches after thresholding. The false positives are mostly eliminated by adding the threshold. Purple bars are the number of putative units (unit with no reference information) inferred with  $z$ -threshold =  $10\mu m$ .

200

## 201 2.4 Units can be tracked in discontinuous recordings for 48 days

202 To assess long-term tracking capabilities, we tracked neurons across all datasets for each mouse.  
 203 **Figure 5** shows a survival plot of the number of unit chains successfully tracked over all durations.  
 204 All units in the plot can be tracked across at least three consecutive datasets, a chain as the term  
 205 is used here. We categorized all trackable unit chains into three types: reference chains, mixed  
 206 chains and putative chains. Reference chains have receptive field information in all datasets. Pu-  
 207 tative chains have no reference information in any of the datasets. Mixed units have at least one  
 208 dataset with no receptive field information. There are 133 reference chains, 135 mixed chains and  
 209 84 putative chains across all the subjects. Among them, 46 reference, 51 mixed, and 9 putative  
 210 units can be followed across all datasets. We refer to them as fully trackable units. One example  
 211 trackable unit in each group is shown in **Figure 6**, **Figure 6-supplement Figure 1**, and **Figure 6-**  
 212 **supplement Figure 2**.



**Fig. 5:** Number of reference units (deep blue, dark orange and green for different subjects), putative (medium green, medium orange and blue) units, and mixed units (light green, yellow, and light blue) tracked for different durations. The loss rate is similar for different chain types in the same subject. Note that chains can start on any day in the full set of recordings, so the different sets of neurons have chains with different spans between measurements.

213 We hypothesize that the three groups of units are not qualitatively different from each other,  
 214 that is, all units are equally trackable. In order to check for differences among the three groups,  
 215 we analyzed the locations, firing rates, waveforms, and receptive fields of the fully trackable units  
 216 in the three groups: reference, putative, and mixed.

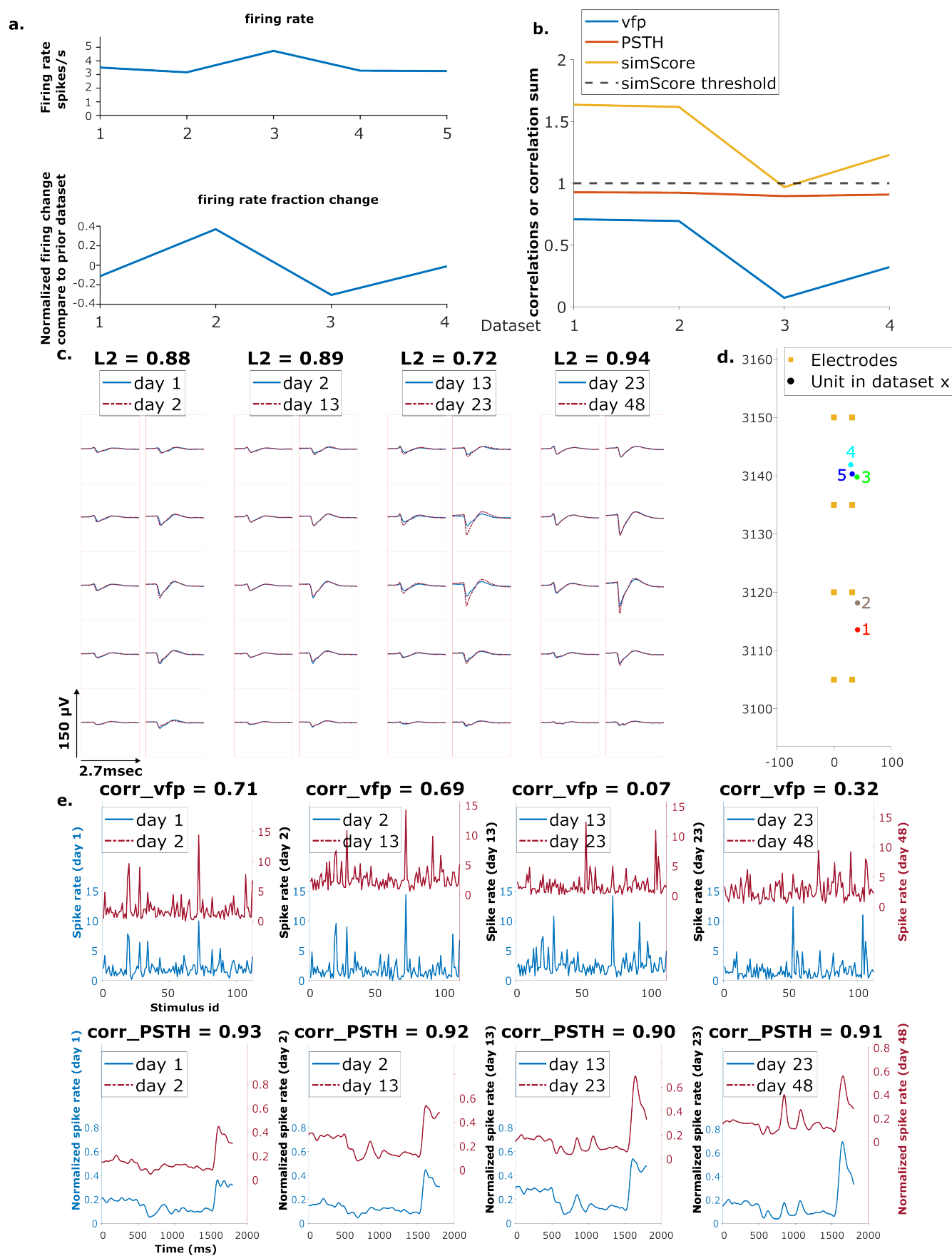
217 The spatial-temporal waveform similarity is measured by the L2 distance between waveforms  
 218 (Sec. 4.1.2). A Kruskal-Wallis test is performed on the magnitude of L2 change between all pairs  
 219 of matched waveforms among the three groups. There is no statistical difference in the waveform  
 220 similarity in reference, putative, and mixed units ( $H = 0.59$ ,  $p = 0.75$ ) (*Figure 5-supplement Figure 1*).  
 221 There is no significant difference in the physical distances of units per dataset ( $H = 1.31$ ,  $p = 0.52$ )  
 222 (*Figure 5-supplement Figure 2*, bottom panel), nor in the location change of units ( $H = 0.23$ ,  $p =$



223 0.89) (**Figure 5-supplement Figure 2**, top panel).

224 Firing rate is characterized as the average firing rate fold change of each unit chain, with firing  
225 rate of each unit in each dataset normalized by the average firing rate of that dataset. There is  
226 no difference in the firing rate fold change in the three groups of units ( $H = 1$ ,  $p = 0.6$ ) (**Figure 5-**  
227 **supplement Figure 3**).

228 The receptive field similarity between units in different datasets is described by visual finger-  
229 print (vfp) correlation and Peristimulus Time Histogram (PSTH) correlation between units, and the  
230 similarity score, the sum of the two correlations (Sec. 4.4). The change in vfp between matched  
231 units is similar among the three groups ( $H = 2.23$ ,  $p = 0.33$ ). Similarly, the change in PSTH is not  
232 different among the three groups ( $H = 1.61$ ,  $p = 0.45$ ) (**Figure 5-supplement Figure 4**).



**Fig. 6: Example mixed chain:** a. Above: Firing rates of this neuron on each day (Day 1, 2, 13, 23, 48). Below: Firing rate fractional change compared to the previous day. b. Visual response similarity (yellow line), PSTH correlation (orange line), and visual fingerprint correlation (blue line). The similarity score is the sum of vfp and PSTH. The dashed black line shows the threshold to be considered a reference unit. c. Spatial-temporal waveform of a trackable unit. Each pair of traces represents the waveform on a single channel. d. Estimated location of this unit on different days. Each colored dot represents a unit on one day. The orange squares represent the electrodes. e. The pairwise vfp and PSTH traces of this unit.

233

### 234 3 Discussion

235 We present here an EMD-based neuron tracking algorithm that provides a new, automated way  
236 to track neurons over long-term experiments to enable the study of learning and adaptation with  
237 state-of-the-art high density electrophysiology probes. We demonstrate our method by tracking  
238 neurons up to 48 days without using receptive field information. Our method achieves 90% recovery  
239 rate on average for neurons separated up to one week apart and 78% on average for neurons  
240 five to seven weeks apart (**Figure 4b**, blue bars). We also achieved 99% accuracy up to one week  
241 apart and 95% five to seven weeks apart, when applying a threshold of 10  $\mu\text{m}$  (**Figure 4b**, green  
242 bars). It also retrieved a total of 552 tracked neurons with partial or no receptive field information,  
243 12 per pair of datasets on average. All the fully trackable unit chains were evaluated by wave-  
244 forms and estimated locations. Our method is simple and robust; it only requires spike sorting be  
245 performed once, independently, per dataset. In order to be more compatible and generalizable  
246 with existing sorting methods, we chose Kilosort, one of the most widely used spike sorting meth-  
247 ods.<sup>33,34</sup> We show the capability of our method to track neurons with no specific tuning preference  
248 (**Figure 6-supplement Figure 2**).

249 The method includes means to identify dataset pairs with very large drift. In our data, we can  
250 detect large drift because such datasets have very few reference units, and significantly different  
251 EMD cost. For example, datasets 1 and 2 in animal AL036 have very few reference units compared  
252 to other datasets (see **Figure 2-supplement Figure 2**, AL036). This observation is consistent with  
253 the overall relationship between the EMD cost and recovery rate (**Figure 2-supplement Figure 3**).  
254 Datasets with higher cost tend to have lower unit recovery rate and higher variation in recovery  
255 rates. Therefore, these two datasets were excluded in the tracking analysis.

256 Our validation relies on identifying reference units. The reference unit definition has limita-  
257 tions. The similarity score is largely driven by PSTHs (**Figure 7-supplement Figure 1**), the timing of  
258 stimulus triggered response, rather than vfp, the response selectivity. As a result, a single neuron  
259 can be highly correlated, i.e. similarity score greater than 1, with more than 20 other neurons. For  
260 example, in subject AL032 shank 2, one neuron on day 1 has 22 highly correlated neurons on day  
261 2, 4 of which are also within the distance of 30  $\mu\text{m}$ . Non-reference units may also have very similar  
262 visual responses: we note that 33 (5 putative neurons and 28 mixed neurons) out of 106 trackable  
263 neurons have a similarity score greater than 1 even for days with no reference unit assignment.  
264 Coincidentally similar visual responses could potentially contribute to inaccurate assignment of  
265 reference units and irregularity in trackable unit analysis. These errors would reduce the mea-  
266 sured accuracy of the EMD matching method; since the accuracy is very high (**Figure 4**), the impact  
267 of mismatches is low.

268 We note that the ratio of reference units over KSGood units decreases as recordings are further  
269 separated in time (**Figure 7-Figure 3**). This reduction in fraction of reference units might be partially  
270 due to representational drift as well as the fact that the set of active neurons are slightly different  
271 in each recording. The visual fingerprint similarity of matched neurons decreased to 60% after 40  
272 days (see reference 7 supplement).

273 We developed the new tracking algorithm based on an available visual cortex dataset, and used

a prominent sorting algorithm (Kilosort 2.5) to spikesort the data. We had reference data to assess the success of the matching and tune parameters. Applying our algorithm in other brain areas and with other sorters may require parameter adjustment. Evaluation of the results in the absence of reference data requires a change to the fitting procedure.

The algorithm has only two parameters: the weighting factor  $\omega$  that sets the relative weight of waveform distance vs. physical distance, and the z-distance threshold that selects matches that are likely correct. We found that recovery rate, and therefore accuracy, is insensitive to the value of  $\omega$  for values larger than 1500 (Figure 2-supplement (Figure 4)), so this parameter does not require precise tuning. However, the false positive rate is strongly dependent on the choice of z-distance threshold.

When reference information (unit matches known from receptive fields or other data) is available, the procedure outlined in Figure 4 can be followed. In that case, the distribution of z-distances of known pairs is fit to find the width of the distribution for correct matches. That parameter is then used in the fit of the z-distance distribution of all pairs to Equation 3. Integrating the distributions of correct and incorrect pairs yields the false positive rate vs. z-distance, allowing selection of a z-distance threshold for a target false positive rate.

In most cases, reference information is not available. However, the z-distance distributions for correct and incorrect pairs can still be estimated by fitting the distribution of all pairs. In Figure 4-supplement Figure 2 we show the results of fitting the z-distribution of all pairs without fixing the width of the distribution of correct matches. The result slightly underestimates this width, and the estimated false positive rate increases. This result is important because it suggests the accuracy estimate from this analysis will be conservative. We detail the procedure for fitting the z-distance distribution Methods section (Alg. 2).

As suggested in Dhawale et al.,<sup>5</sup> discontinuous recordings will have more false positives. Improving spike sorting and restricting the analysis to reliably sorted units will help decrease the false positive rate. Current spike sorting methods involve fitting many parameters. Due to the stochastic nature of template initialization, only around 60% to 70% units are found repeatedly in independently executed analysis passes. This leads to unpaired units which decreases EMD matching accuracy. Future users may consider limiting their analysis to the most reliably detected units for tracking; requiring consensus across analysis passes or sorters is a possible strategy. Finally, more frequent data acquisition during experiments will provide more intermediate stages for tracking and involves smaller drift between consecutive recordings.

## 4 Methods

Our neuron tracking algorithm uses the Earth Mover's Distance (EMD) optimization algorithm. The minimized distance is a weighted combination of physical distance and 'waveform distance': the algorithm seeks to form pairs that are closest in space and have the most similar waveforms. We test the performance of the algorithm by comparing EMD matches to reference pairs determined from visual receptive fields (Sec. 4.4). We calculate two performance metrics. The 'recovery rate' is the percentage of reference units that are correctly matched by the EMD procedure. The 'accuracy' is the percentage of correctly matched reference units that pass the z-distance threshold (Figure 4a). 'Putative units' are units matched by the procedure which do not have reference receptive field information. 'Chains' are units that can be tracked across at least three consecutive datasets. The full procedure is summarized in Algorithm 1.

---

### Algorithm 1 Neuron Matching Procedure

---

**Input:** channel map, unit cluster label, cluster mean waveforms (with  $K_{loc} = 2$  and  $K_{wf} = 5$  rows and  $K_{col} = 2$  columns of channels), and spike times

#### Step 1 Estimate unit locations

Estimate background amplitude for each unit

**for** all KSgood units  $u_n \in U$  **do**

**if** peak-top-peak voltage  $V_{ptp} > 60\mu V$  **then**

        Get  $u_n$ 's waveform on channels  $C_m$

        Get the peak-to-peak amplitudes  $V_{ptp_c}$  of  $u_n$  background-subtracted waveforms on channels

$C_{u_n} = \{mc_{u_n} - k_{loc}, \dots, mc_{u_n} + k_{loc}\}$  where  $mc_{u_n}$  is the peak channel

        Estimate the neuron's 3D location as in:<sup>32</sup>

$f(x, y, z) = \sum_{c \in C_{u_n}} (V_{ptp_c} - \frac{1}{\sqrt{(x-x_c)^2 + (z-z_c)^2 + y^2}})^2$  where  $x$ ,  $z$ , and  $y$  are the horizontal location, vertical location, and distance of the unit from the probe, respectively.

        Find an estimate of the global minimizer of  $f$ ,  $x_{u_n}$ ,  $y_{u_n}$ ,  $z_{u_n}$  using least-squares optimization

**end**

**end**

#### Step 2 Compute waveform similarity metrics

**for** waveforms  $wf_{xi} \in U_{N1}$  and  $wf_{yk} \in U_{N2}$  where  $U_{N1}, U_{N2}$  are the set of all units in the two datasets **do**

    Centered at peak channel  $mc_{xi}$  and  $mc_{yk}$ , respectively

    Get the sets of channels for each unit:  $C_{u_n} = \{mc_{u_n} - k_{wf}, \dots, mc_{u_n} + k_{wf}\}$

    There are  $K_{wf} * 2 * K_{col} + 2 = 22$  channels for each unit

    Compute the waveform similarity metric as  $(1/22) * \sum_{c \in C_{u_{xi}}, C_{u_{yk}}} L2(wf_{xi} - wf_{yk}) / \max(L2(wf_{xi}), L2(wf_{yk}))$  for each of the 22 channels

**end**

#### Step 3 Between-session drift correction

    Run the EMD with distances in physical and waveform space

    Estimate z-distance mode of all matched pairs with Gaussian kernel fit

    Apply correction on physical distances of all units  $\in U_2$ :  $z_{corr} = z - z_{mode}$

#### Step 4 Unit matching

    Run the EMD with corrected physical distance and waveform metrics

    Set z-distance threshold to select unit pairs likely to be the same neuron

**Output:** cost  $\sum d_{EMD}$ , unit assignments

---

## 4.1 Algorithm

### 4.1.1 Earth Mover's Distance

The EMD is an optimization-based metric developed in the context of optimal transport and measuring distances between probability distributions. It frames the question as moving dirt, in our case, units from the first dataset, into holes, which here are the neural units in the second dataset. The distance between the "dirt" and the "holes" determines how the optimization program will prioritize a given match. Specifically, the EMD seeks to minimize the total work needed to move the dirt to the holes, i.e., neurons in day 1 to day 2, by solving for a minimum overall effort, the sum of distances.<sup>30,31</sup>

$$\begin{aligned}
 & \min_{d_F} \sum_{i,k} D(x_i, y_k), \text{ where } D = d_{loc} + \omega d_{wf} \\
 & \text{subject to } f_{ik} \in [0, 1] \forall i, k \\
 & \sum_k (f_k) \leq \text{length}(Y) \\
 & \sum_i (f_i) \leq \text{length}(X) \\
 & \sum(F) = \min(\sum X, \sum Y)
 \end{aligned} \tag{4}$$

in which  $d_{loc} \in \mathcal{D}^3$  is the three-dimensional physical distance between a unit from the first dataset  $x_i$ , and a unit from the second dataset  $y_k$ .  $d_{wf} \in \mathcal{D}^1$  is a scalar representing the similarity between waveforms of units  $x_i$  and  $y_k$ .  $\omega$  is a weight parameter that was tuned to maximize the recovery rate of correctly matched reference units.  $F$  is the vector of matched objects between the two datasets (See **Figure 2**-supplement **Figure 4** for details about selecting weight).

The EMD has three benefits:

- It allows combining different types of information into the 'distance matrix' to characterize the features of units.
- The EMD can detect homogeneous movement of units (**Figure 2c**), thus providing a way for rigid drift correction, as described in section 4.1.3.
- By minimizing overall distances, the EMD has tolerance for imperfect drift correction, error in the determination of unit positions, and possible non-rigid motion of the units.

However, since the EMD is an optimization method with no assumptions about the biological properties of the data, it makes all possible matches. We therefore added a threshold on the permissible z-distance to select physically plausible matches.

#### 4.1.2 Calculating the EMD distance metric

The unit locations are estimated by fitting 10 peak-to-peak (PTP) amplitudes from adjacent electrodes and the corresponding channel positions with a 1/R distance model.<sup>32</sup> Unlike Boussard, et al.,<sup>32</sup> we operate on the mean waveforms for each unit rather than individual spikes. We found using the mean waveform yields comparable results and saves significant computation time. Unit locations are three-dimensional coordinates estimated relative to the probe, where the location of the first electrode on the left column at the tip is considered the origin. The mean waveform is computed by averaging all the spike snippets assigned to the cluster by KS 2.5.

For 10 channels  $c \in C_{u_n}$ , find the location coordinates  $x_{u_n}, y_{u_n}, z_{u_n}$  that minimizes the difference between measured amplitudes  $V_{PTP}$  and amplitudes estimated with locations  $\frac{\alpha}{\sqrt{(x-x_c)^2+(z-z_c)^2+y^2}}$ :

$$\min \sum_{c \in C_{u_n}} \left( V_{PTP_c} - \frac{1}{\sqrt{(x-x_c)^2+(z-z_c)^2+y^2}} \right)^2 \tag{5}$$

The locations are used to calculate the physical distance portion of the EMD distance.

For the waveform similarity metric, we want to describe the waveform characteristics of each unit with its spatial-temporal waveform at the channels capturing the largest signal. The waveform similarity metric between any two waveforms  $u_{n1}$  and  $u_{n2}$  in the two datasets is a scalar calculated as a normalized L2 metric (see Alg.1 Step 2) on the peak channels, namely the channel row with the highest amplitude and 5 rows above and below (a total of 22 channels). The resulting scalar reflects the 'distance' between the two units in the waveform space and is used to provide information about the waveform similarity of the units. It is used for between-session drift correction and neuron matching. **Figure 1c** shows an example waveform of a reference unit.

### 360 4.1.3 Between-session Drift Correction

361 Based on previous understanding of the drift in chronic implants, we assumed that the majority  
362 of drift occurs along the direction of the probe insertion, i.e. vertical z-direction. This rigid drift  
363 amount is estimated by the mode of the z-distance distribution of the EMD assigned units using a  
364 normal kernel density estimation implemented in MATLAB. We only included KSgood units.<sup>16</sup> The  
365 estimated drift is then applied back to correct both the reference units and the EMD distance matrix  
366 by adjusting the z coordinates of the units. For validation, the post-drift-correction reference set is  
367 compared with the post-drift-correction matching results (from step 4 in 1).

## 368 4.2 Determining Z Distance Threshold

369 Determining the z-distance threshold to achieve a target false positive rate requires estimating  
370 the widths of the z-distance distributions of correct and incorrect pairs. If reference data is avail-  
371 able, the z-distance distribution of the known correct pairs should be fit to a folded Gaussian as  
372 described in **Figure 4**. The width of the folded Gaussian, which is the error in determination of the  
373 z-positions of units, is then fixed in the fit of the z-distribution of all pairs found by the algorithm  
374 outlined in Algorithm 4.1.1. If no reference data is available, the width of the distribution of correct  
375 pairs is determined by fitting the z-distance distribution of all pairs to **Equation 3** with the folded  
376 Gaussian width as one of the parameters. This procedure is detailed in Algorithm 2. We show two  
377 examples of model fitting without reference information in **Figure 4**-supplement **Figure 2**.

---

### Algorithm 2 Determining an appropriate z distance threshold

---

**Input:** Z distances of all matched units, target false positive rate, width  $\sigma$  of the z-distance distribu-  
tion of correct pairs, if available

**Step 1** Fit z distance distribution of all pairs to decompose into distributions of correct and incor-  
rect pairs

Fit the z-distance distribution of all pairs to the sum of a folded Gaussian (for correct pairs) and  
an exponential (for incorrect pairs). If the width  $\sigma$  of the distribution of correct pairs is known  
from reference data, fix at that value. Otherwise, include in the fit parameters. The functional  
form is:  $P(z) = d(fNe^{-\frac{z^2}{2\sigma^2}} + \frac{1-f}{c}e^{-\frac{z}{c}})$

Where:  $f$  = fraction of correct pairs;  $\sigma$  = width of the distribution of correct pairs;  $c$  = decay  
constant of distribution of incorrect pairs;  $d$  = amplitude normalization; and  $N = \frac{2}{\sigma\sqrt{2\pi}}$ , the  
normalization factor of the folded Gaussian.

**Step 2** Determine z threshold to achieve a target false positive rate

For Neuropixels 1.0 and 2.0 probes, the width of the z-distance distribution of correct matches  
( $\sigma$ ) should be  $<10 \mu\text{m}$ ; a larger width, or a very small value of the fraction of correct pairs  
suggests few or no correct matches. In this case, the EMD cost is likely to be large as well (See  
**Figure 2**-supplement **Figure 2** Animal AL036 first two rows).

For a range of z values, integrate the z-distance distribution of incorrect pairs from 0  
to z, and divide by the integral of the distribution of all pairs over that range. This gener-  
ates the false positive rate vs. z-distance threshold, as shown in **Figure 4**-supplement **Figure 2**.

---

**Output:**  $\sigma$  (uncertainty of position estimation), threshold at the target false positive rate

---

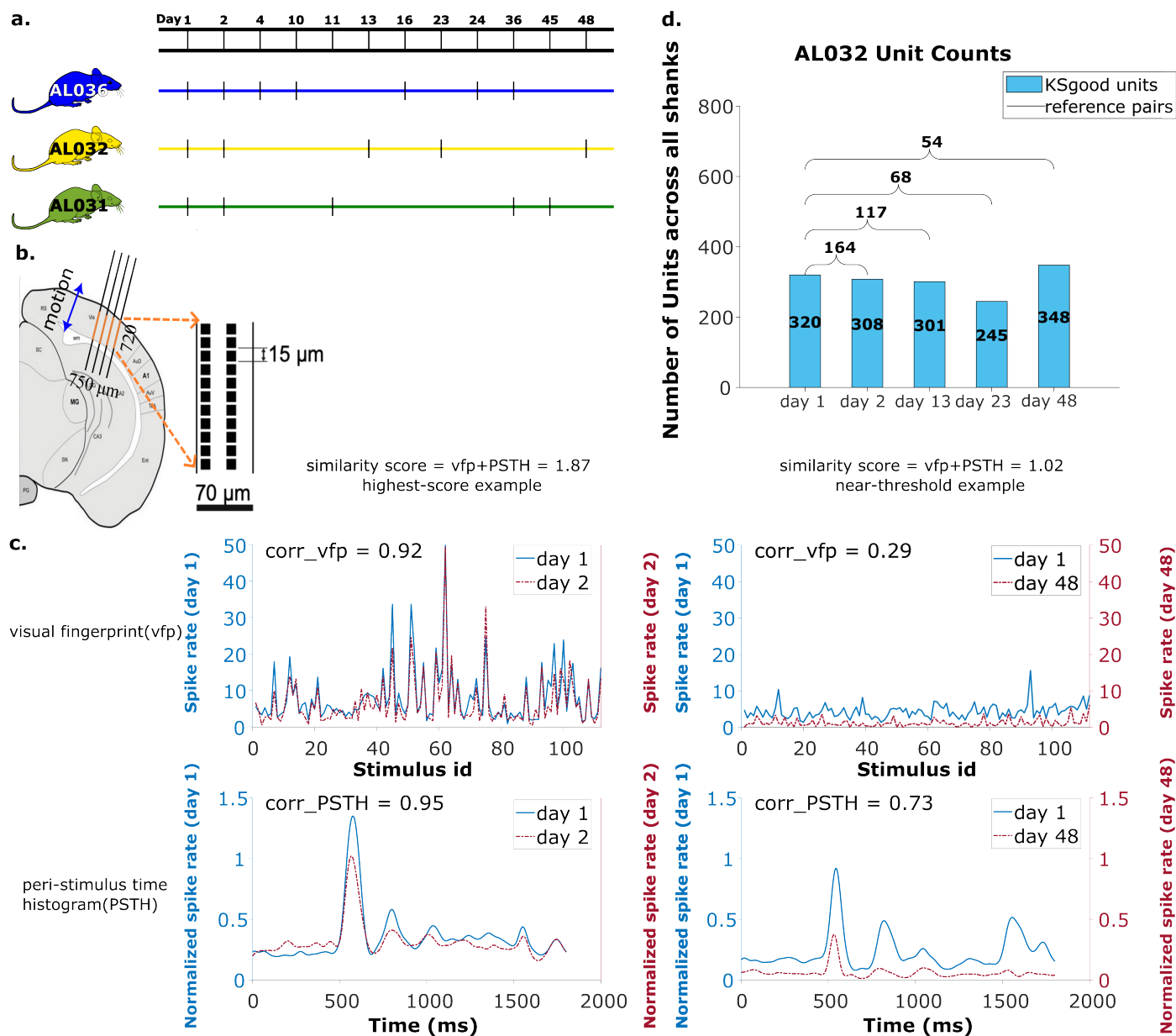
## 378 4.3 Dataset

379 The data used in this work are recordings collected from two chronically implanted NP 2.0 four-  
380 shank probes and one chronically implanted one-shank NP 2.0 probe in the visual cortex of three  
381 head fixed mice (**Figure 7b**, see Steinmetz et al.<sup>7</sup> for experiment details). The recordings were taken  
382 while 112 visual stimuli were shown from three surrounding screens (data from Steinmetz et al.<sup>7</sup>



Supplement Section 1.2). The same bank of stimuli was presented five times, with order shuffled. The 4-shank probes had the 384 recording channels mapped to 96 sites on each shank.

We analyzed 65 recordings, each from one shank, collected in 17 sessions (5 sessions for animal AL031, 5 sessions for animal AL032, and 7 sessions for animal AL036). The time gap between recordings ranges from one day to 47 days (**Figure 7a**), with recording duration ranging from 1917 to 2522 seconds. The sample rate is 30kHz for all recordings. There are a total of 2958 KSorted units analyzed across all animals and shanks, with an average of 56 units per dataset (**Figure 7d** and **Figure 7-supplement Figure 2**).



**Fig. 7: Summary of dataset:** a. The recording intervals for each animal. A black dash indicates one recording on that day. b. All recordings are from visual cortex V1 with a 720  $\mu\text{m}$  section of the probe containing 96 recording sites. The blue arrow indicates the main drift direction. c. Examples of visual fingerprint(vfp) and peri-stimulus time histogram(PSTH) from a high correlation (left column) and a just-above-threshold (right column) correlation unit. Both vfp and PSTH values vary from  $[-1,1]$ . d. Kilosort-good and reference unit counts for animal AL032, including units from all four shanks.

# 4.4 Reference set

To track clusters across days, Steinmetz et al.<sup>7</sup> concatenated two recording sessions and took advantage of the within-recording drift correction feature of Kilosort 2.0 to extract spikes from the two days with a common set of templates. They first estimated the between session drift of each recording from the pattern of firing rate and amplitude on the probe and applied a position correction of an integer number of probe rows (15  $\mu\text{m}$  for the probes used). Then two corrected recordings were concatenated and sorted as a single recording. This procedure ensured that the same templates are used to extract spikes across both recordings, so that putative matches are extracted with the same template. A unit from the first half of the recording is counted as the same neuron if its visual response is more similar to that from the same cluster in the second half of the recording than to the visual response of the physically nearest neighbor unit. Using this procedure and matching criteria, 93% of the matches were correct for recordings < 16 days apart, and 85% were correct for recordings from 3-9 weeks (See Steinmetz et al.,<sup>7</sup> Fig. 4). In addition, although mean fingerprint similarity decreases for recordings separated by more than 16 days, this decline is only 40% for the same unit recorded from 40 days apart (see Steinmetz et al.<sup>7</sup> Supplement S3). This procedure, while successful in their setting, was limited to the use of integral row adjustments of the data for between-session drift correction and relied on a customized version of Kilosort 2.0. Although up to three recordings can be sorted together, they must come from recording sessions close in time. In addition, a separate spike sorting session needs to be performed for every pair of recordings to be matched, which is time consuming and introduces extra sorting uncertainty.

To find units with matched visual responses, we examine the visual response similarity across all possible pairs. The visual response similarity score follows Steinmetz et al.,<sup>7</sup> and consists of two measurements. 1) The peristimulus time histogram (PSTH), which is the histogram of the firing of a neuron across all presentations of all images, in a 1800 msec time window starting 400 msec before and ending 400 msec after the stimulus presentation. The PSTH is calculated by histogramming spike times relative to stimulus on time for all stimuli, using 1 ms bins. This histogram is then smoothed with a Gaussian filter. 2) The visual fingerprint(vfp) is the average response of the neuron to each of the 112 images. The vfp is calculated by averaging the spike counts in response to each natural image from the stimulus onset to 1 second afterwards across 5 shuffled trials.

Following Steinmetz et al.,<sup>7</sup> the similarity score between two neurons is the sum of the correlation of the PSTH and the correlation of the vfp across two sessions. The two correlations have values in the range (-1,1), and the similarity score ranges from (-2, 2).

The pool of reference units is established with three criteria: 1) The visual response similarity score of the pair, as described above, is greater than 1 and their physical distance, both before and after drift correction, is smaller than 30  $\mu\text{m}$ . A physical distance criterion is necessary, because some units have several potential partners with high visual response similarity (*Figure 7-supplement Figure 1*). We impose the 30  $\mu\text{m}$  threshold on both pre- and post-correction data because the drift is relatively small in our case, and we can reduce false positives by constraining the reference units to be in a smaller region without losing units. In general, one could apply the threshold only on corrected data (after drift correction). 2) A Kruskal-Wallis test is applied on all trials of the vfps to ensure the triggered response to the stimulus is significantly distinguishable from a flat line. 3) Select units from each recording that meet the good criteria in Kilosort. Kilosort assigns a label of either single-unit (good) or multi-unit (MUA) to all sorted clusters based on ISI violations.<sup>16</sup> This step aims to ensure included units are well separated. If there are multiple potential partners for a unit, the pair with the highest similarity score is selected as the reference unit. The complete pool of reference units includes comparisons of all pairs of recordings for each shank in each animal. The portion of units with qualified visual response ranges from 5% to 61%, depending on the time gap between datatets (*Figure 7-supplement Figure 3*). Overall, these reference units made up 29% of all KSgood units (*Figure 7-supplement Figure 2*) across all three animals in our dataset. *Figure 7c* shows examples of visual responses from a high similarity reference unit and a

441 reference unit with similarity just above threshold.

## 442 **5 Code sharing**

443 All code used can be accessed at: [https://github.com/janelia-TDHarrisLab/Yuan-Neuron\\_Tracking](https://github.com/janelia-TDHarrisLab/Yuan-Neuron_Tracking).

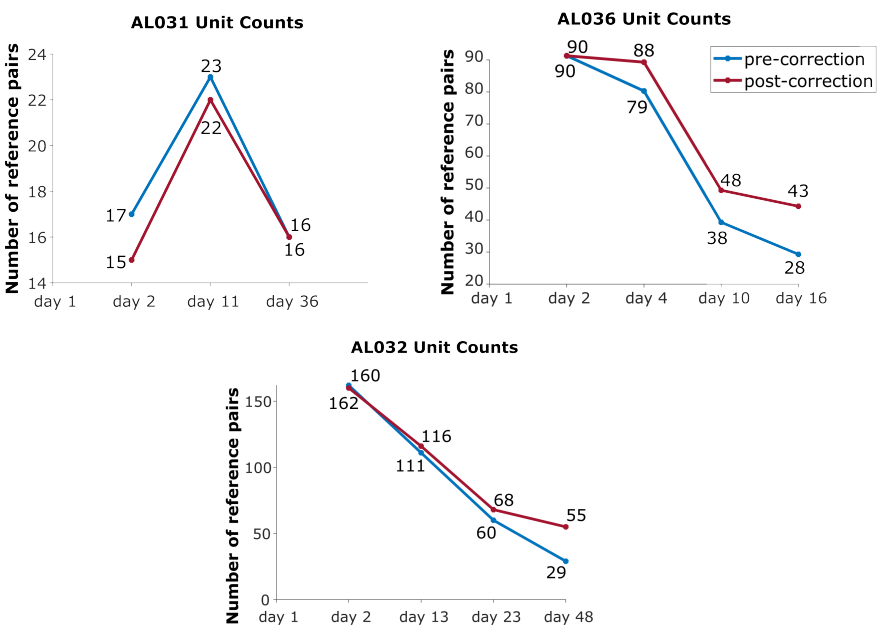
## 444 **6 Acknowledgments**

445 NIH grant U01 NS115587 in part supported TDH and AXY. We thank Claudia Böhm and Albert Lee  
446 for allowing us to use their data in *Figure 4*-supplement *Figure 2*.

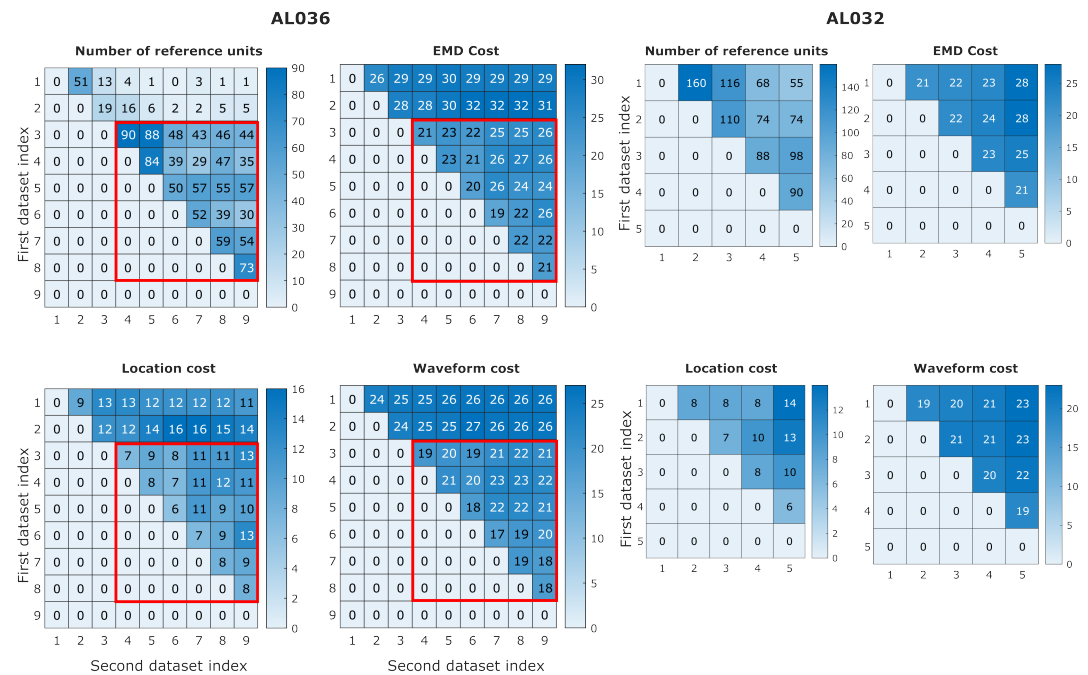
## 447 **7 Declaration of interests**

448 The authors declare no competing interests.

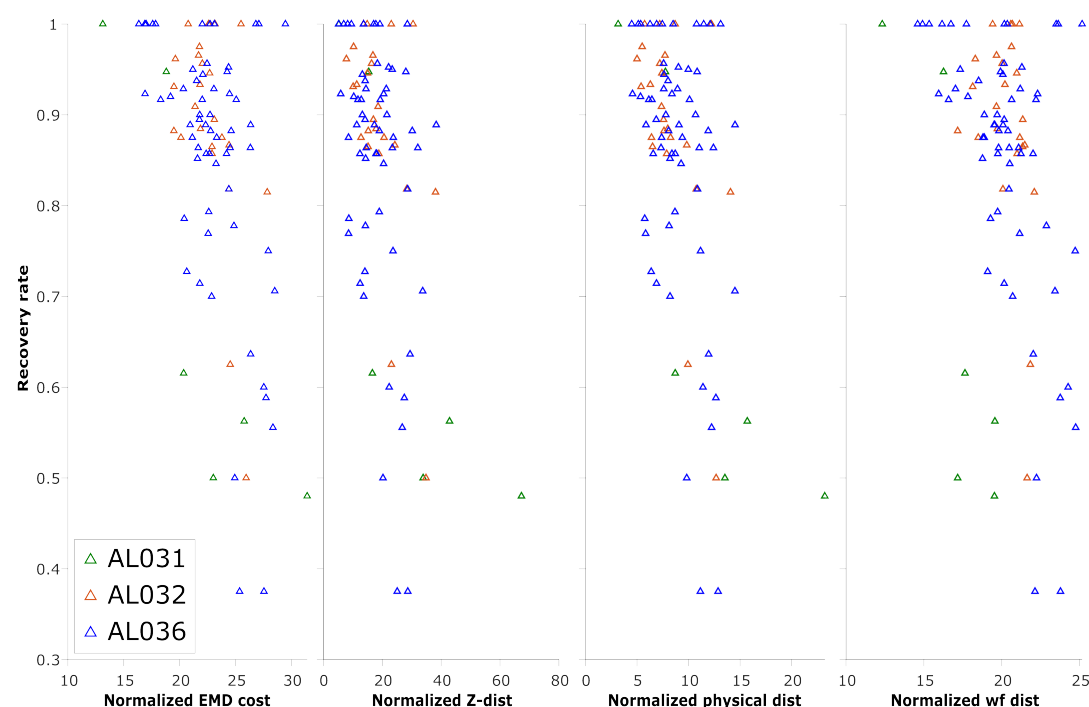
449 **8 Figure 2 supplement**



**Figure 2 - figure supplement 1:** The effect of drift correction on reference units yield for all three animals. Note that drift correction improves the recovery rate for most cases; the degree of improvement is a function of the magnitude of the drift.

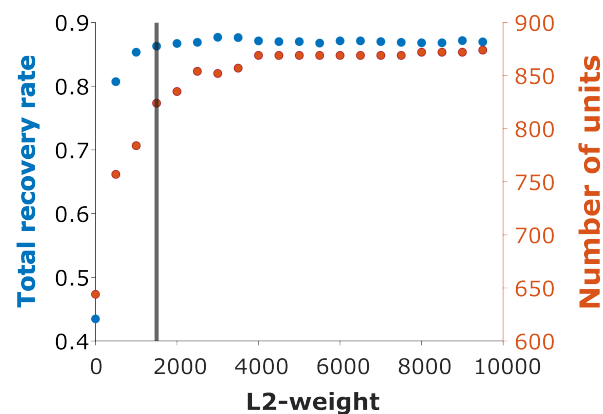


**Figure 2 - figure supplement 2:** EMD cost can be used to detect discontinuities in the data. In animal AL036, we noted a large decrease in the number of reference units (units with matched visual responses, see Sec. 4.4) after the second dataset. This likely indicates a large physical shift in the tissue relative to the probe. It is important to be able to detect such discontinuities to eliminate datasets from consideration. We find that the discontinuity can be detected in the EMD mean cost, location mean cost and waveform mean cost. The four heatmaps on the left, show reference counts and pairwise costs for units matched on one shank in animal AL036. Note that the days with few reference units also have higher EMD cost. To show that days 1-2 (first two rows) are significantly different from days 3-9, we use the Mann-Whitney U Test. All three cost values show significant differences between the groups (EMD mean cost, reject  $H_0$ ,  $p = 6 \times 10^{-7}$ ; location mean cost, reject  $H_0$ ,  $p = 6 \times 10^{-5}$ ; waveform mean cost, reject  $H_0$ ,  $p = 5 \times 10^{-7}$ ). To show that days 3-9 come from the same distribution, we compare odd and even rows using the same test. All three cost values show no significant difference between odd and even days (accept  $H_0$ ,  $p = 0.92$ ). Based on this significant difference between days 1-2 and later days (datasets in the red rectangles), we infer that the first two datasets sampled a different population of units than the later recordings. These first two datasets were eliminated from our analysis. Matrices on the right show similar information for animal AL032 for reference. To estimate the relative magnitude of EMD cost in related datasets versus unrelated datasets, we calculated the cost between unrelated datasets with similar number of units (AL032 shank 1 and AL036 shank 1, EMD cost = 78, location cost = 67, and waveform cost = 32). The EMD cost is between 70-80, much larger than observed for related datasets (between 20-30).



**Figure 2 - figure supplement 3:** The normalized EMD cost (unitless), z distance ( $\mu\text{m}$ ), physical distance ( $\mu\text{m}$ ), and waveform distance (unitless) and the corresponding recovery rate of reference unit (units with matched visual responses) in pairwise matches of all to all pairs of recordings, on each shank. Each triangle represents the recovery rate in a pair of datasets. Animal AL031 has 6 sets of matching, with one outlier removed. Animal AL032 has 24 sets of matching. Animal AL036 has 60 sets of matched units. Overall, most of the datasets with high recovery rates have per-unit EMD in the range 20-30, but datasets with lower recovery are in the same range. Therefore, while very high EMD cost reveals discontinuous data, EMD cost in the normal range is not predictive of reference unit recovery, which is a metric of match success.

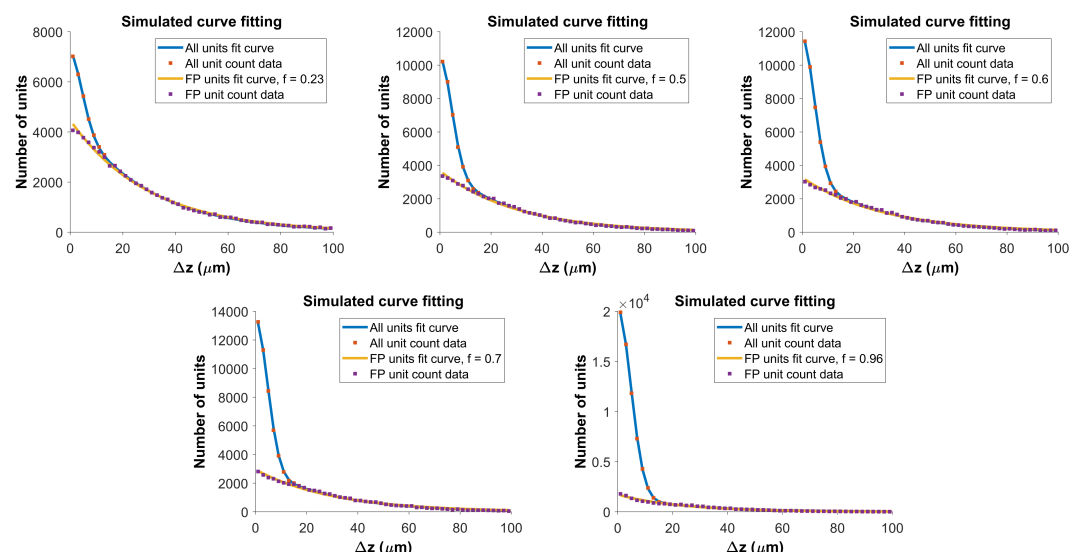
# Recovery rate across subjects v.s. waveform metrics weight



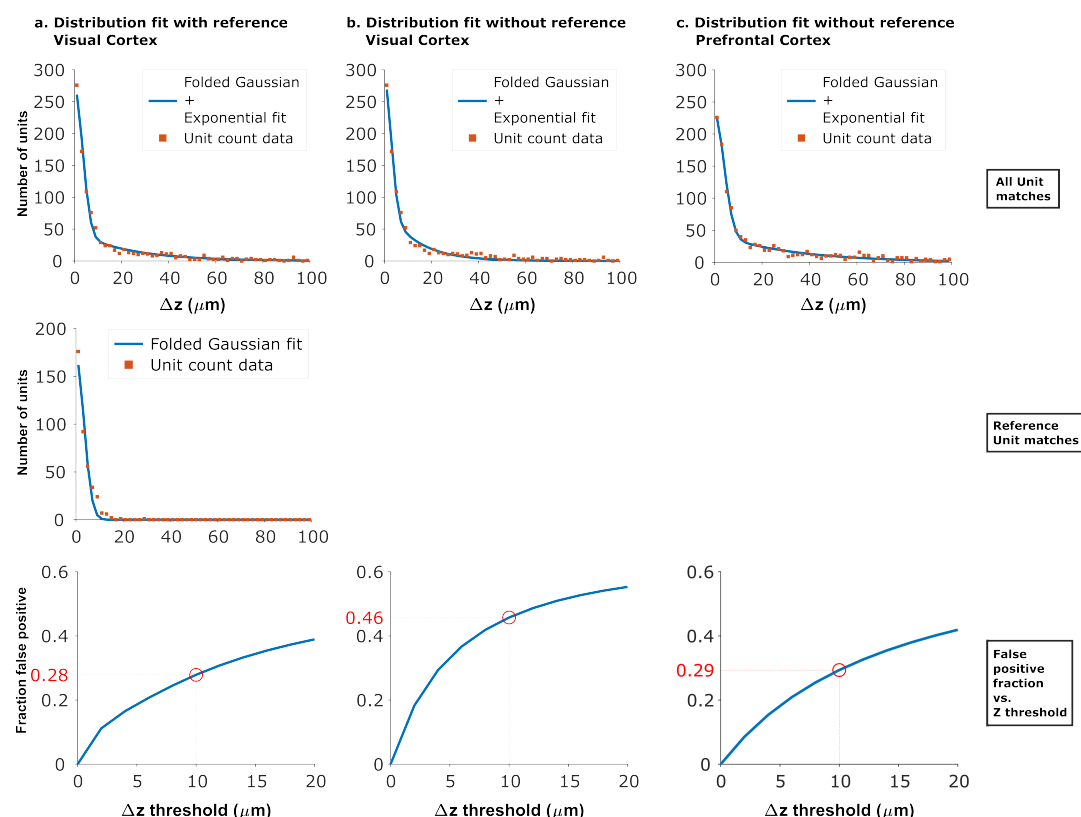
**Figure 2 - figure supplement 4:** Recovery rate vs. L2-weight. We varied the weight  $\omega$  in *Equation 4* used to combine the physical and waveform distances in increments of 500. The vertical line indicates weight = 1500, where the overall recovery rate = 86.29%. The maximum recovery rate = 87.68% occurs at weight = 3000. We chose weight = 1500 for all subsequent analysis.



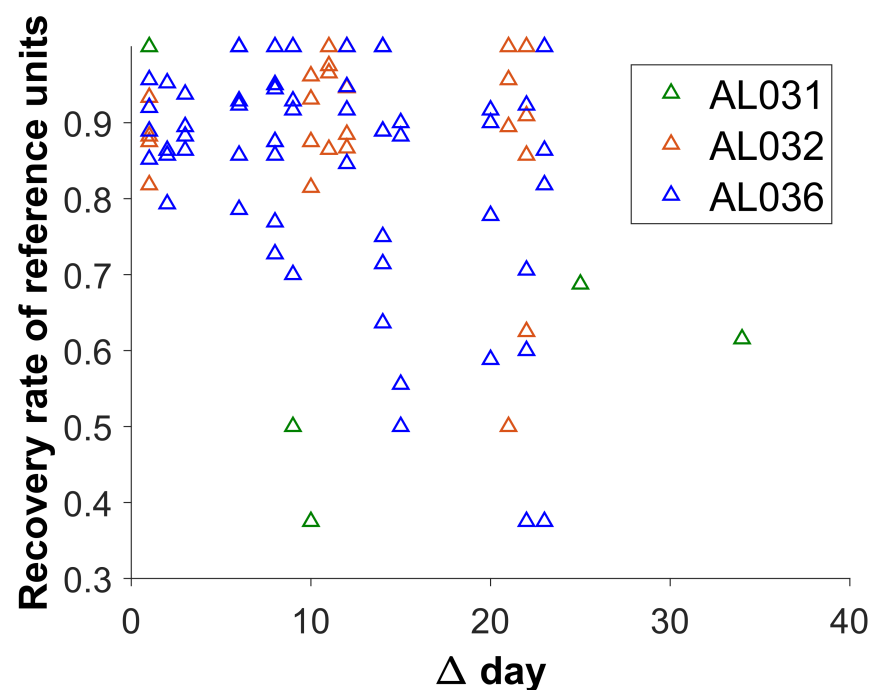
## 450 9 Figure 4 supplement



**Figure 4 - figure supplement 1:** Determining the functional form for the z-distance distribution of all pairs. As shown in **Figure 4a**, the z distance distribution of reference pairs differs significantly from that of all pairs. The z-distance distribution for all pairs is the sum of z-distance distributions for true hits ( $P(\Delta | H)$ ) and false positives ( $P(\Delta | \sim H)$ ), weighted by the fraction correct,  $f$ :  $P(\Delta) = f * P(\Delta | H) + (1 - f) * P(\Delta | \sim H)$ . We built a Monte Carlo model, with 150 units (the average density of subject AL032), normally distributed error  $\sigma = 5\mu m$  for the measured location of the units in true pairs, and random placement of false positives. For each value of fraction correct, we ran the model 500 times. The figure shows fits to model distributions with fraction correct = 0.23, 0.5, 0.6 (top row) and  $f = 0.7, 0.96$  (bottom row). The resulting z-distance distributions are well fit using a folded Gaussian for the distance distribution of true hits and an exponential for the distance distribution of false positives (see Algorithm 2). We use these functional forms to fit the experimental z-distance distribution and estimate the false positive rate.

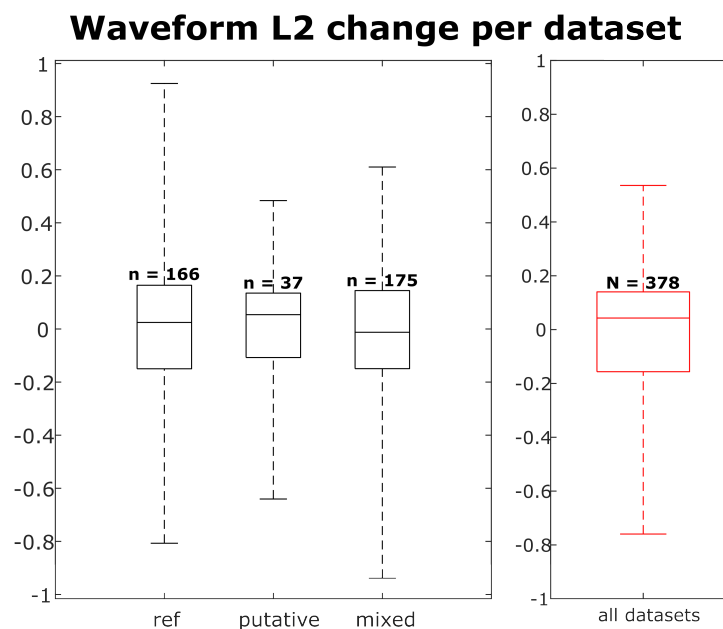


**Figure 4 - figure supplement 2:** Fits of experimental z-distance distributions to the model. When reference data is available, the z-distance distribution of these known true hits can be fit to obtain the width  $\sigma$  of the folded Gaussian.  $\sigma$  can then be fixed in the fit of the distribution of all KSGood units to *Equation 3*, which is used to estimate the false positive rate. When no reference data are available,  $\sigma$  can be estimated from fitting the distribution of all KSGood units to all four parameters in *Equation 3*. Panels a and b show the dataset from *Figure 4* fit with and without fixing the folded Gaussian distribution width. The resulting false positive rate from the no-reference fit at threshold  $z = 10\mu\text{m}$  is larger than that from the fit using reference data, so the procedure gives a conservative estimate of the accuracy. Panel c. shows the model fit of data from an unrelated dataset acquired with from mouse prefrontal cortex using Neuropixels 1.0.<sup>35</sup> The similar shape of the distribution and a 29% false positive rate suggests that the method can be generalized.

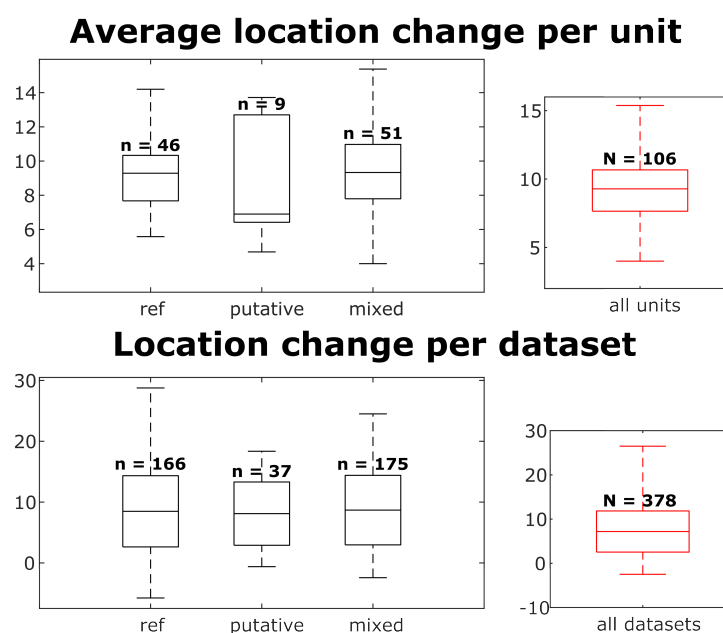


**Figure 4 - figure supplement 3:** The reference unit recovery rate vs. days between matched recordings. Each triangle represents the matching results of two datasets. Animal AL031 has 6 sets of matched units, with one outlier removed. Animal AL032 has 24 sets of matched units. Animal AL036 has 60 sets of matching. The recovery rate is lower for longer durations.

451 **10 Figure 5 supplement**

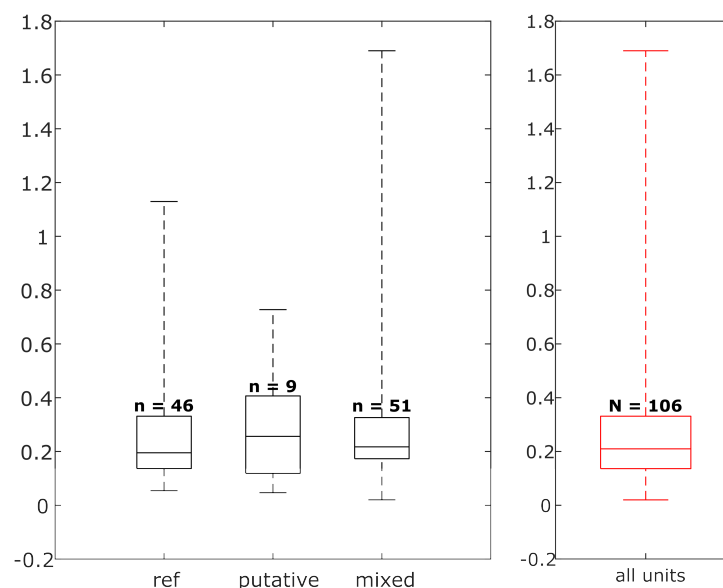


**Figure 5 - figure supplement 1:** Distribution of waveform L2 similarity change per dataset for each neuron group (reference, putative and mixed) and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of unit comparisons, i.e. (number of units)×(number of datasets - 1). A Kruskal-Wallis test indicates no difference among the three groups.

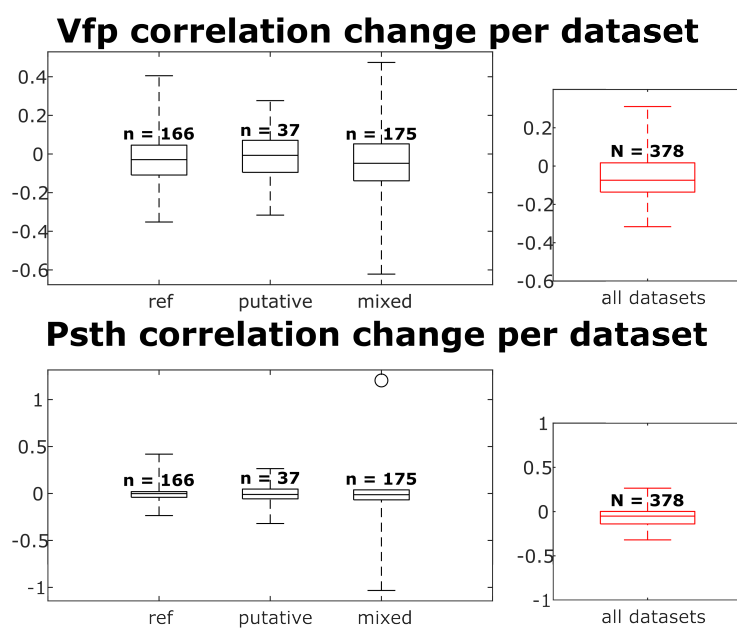


**Figure 5 - figure supplement 2:** Distributions of individual unit location changes over whole chains (top) and unit location changes between pairs of datasets (bottom), for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. In the top plot, n and N are the number of units. In the bottom plot, n and N are the number of unit comparisons, i.e. (number of units)x(number of datasets - 1). A Kruskal-Wallis test indicates no difference among the three groups.

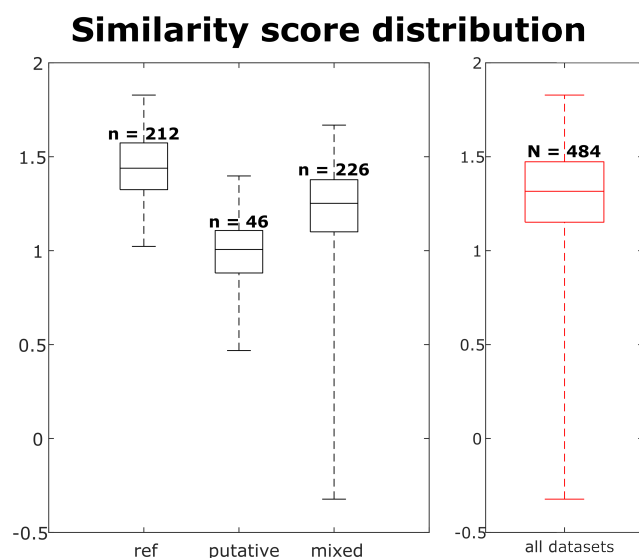
## Average firing rate change ratio per unit



**Figure 5 - figure supplement 3:** Distribution of firing rate fold change per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of units. A Kruskal-Wallis test indicates no difference among the three groups.



**Figure 5 - figure supplement 4:** The visual fingerprint and PSTH change distributions per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of unit comparisons, i.e. (number of units) × (number of datasets - 1). A Kruskal-Wallis test indicates no difference among the three groups.



**Figure 5 - figure supplement 5:** The similarity score distribution per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of observations of the units, i.e.  $\sum_{units} (\text{observations of this unit})$

**Figure 6 - figure supplement 1:** Example reference chain. a. Above: Firing rates of this neuron on each day. Below: Firing rate fractional change compared to the previous day. b. Visual response similarity (yellow line), PSTH correlation (orange line), and visual fingerprint correlation (blue line). The similarity score is the sum of vfp and PSTH. The dashed black line shows the threshold to be considered a reference unit. c. Spatial-temporal waveform of a trackable unit. Each pair of traces represent the waveform on a single channel. d. Estimated location of this unit on different days. Each colored dot represents a unit on one day. The orange squares represent the electrodes. e. The pairwise vfp and PSTH traces of this unit.



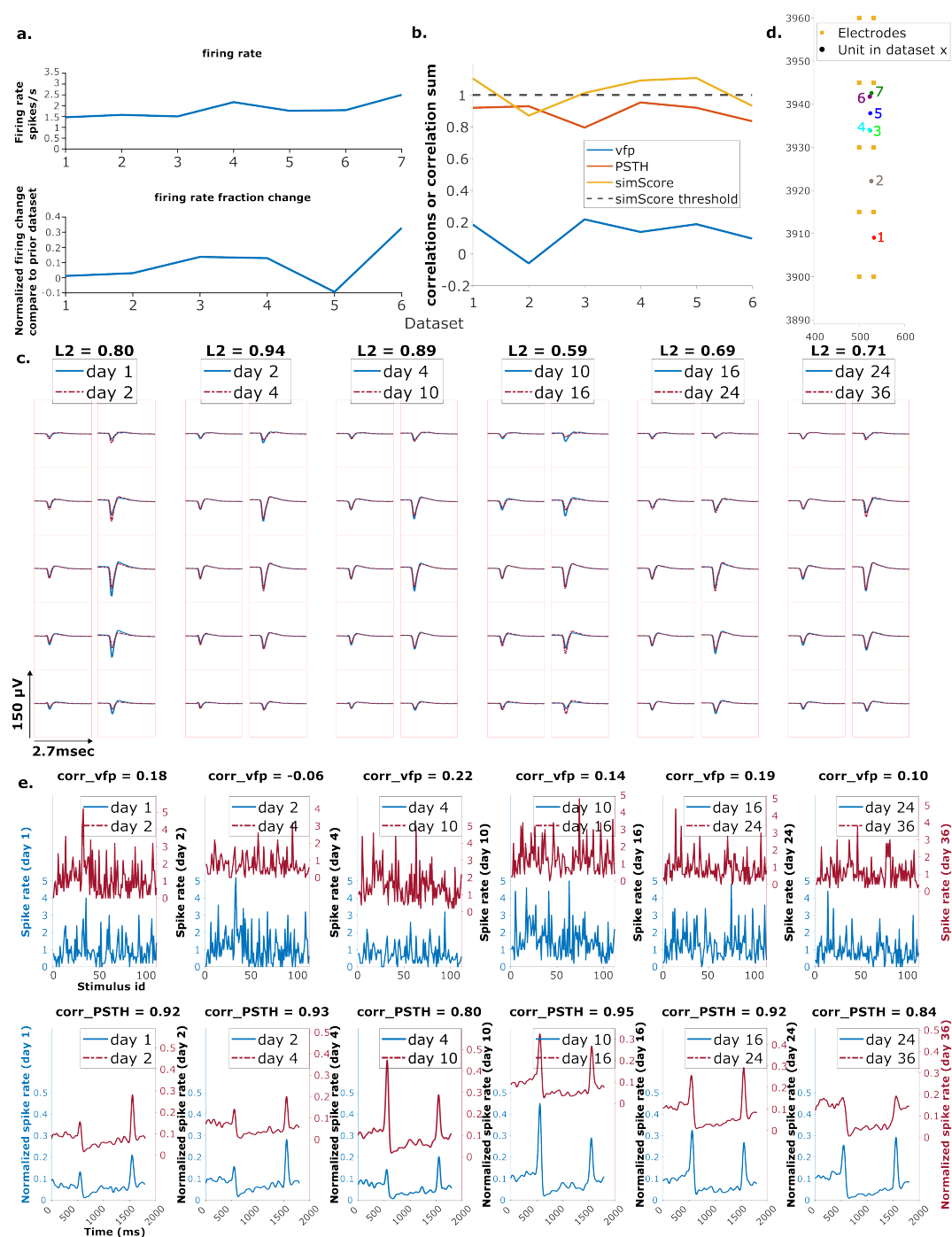
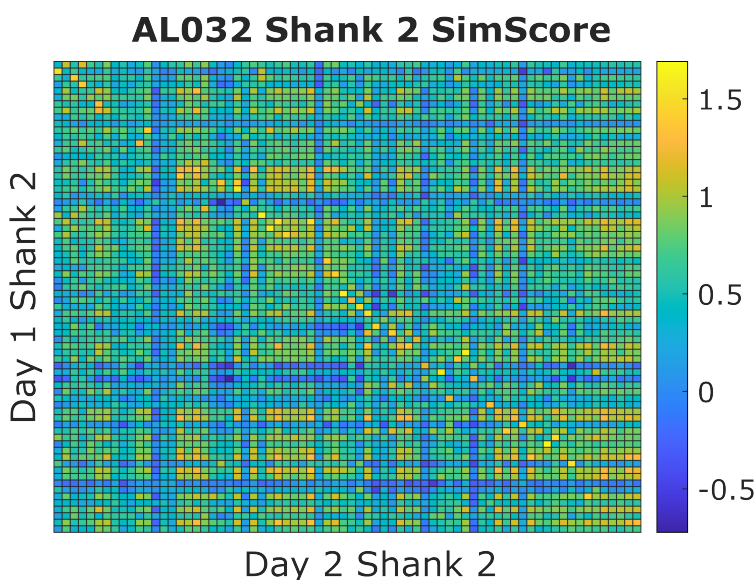
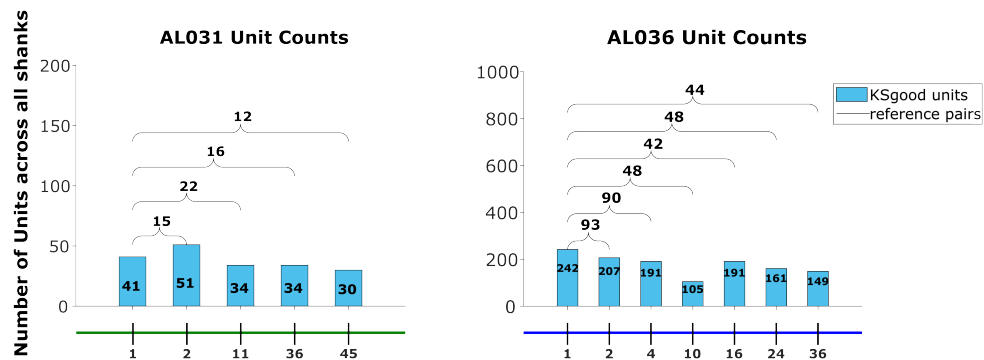


Figure 6 - figure supplement 2: Example putative chain. Order is the same as the previous figure.

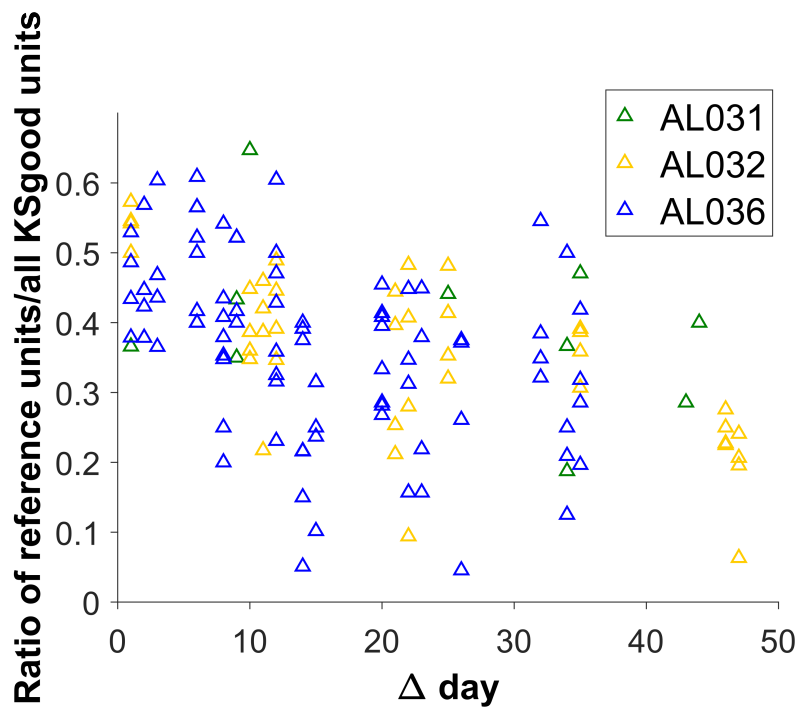
453 **12 Figure 7 supplement**



**Figure 7 - figure supplement 1:** An example similarity score (vfp + PSTH) heatmap from animal AL032, shank 2, Kilosort-good units between day 1 and 2. Each small square represents the similarity score (value range from [-2,2]) between one unit from day 1 and one unit from day 2. A warm colored square indicates a higher score. The clusters are ordered by their physical locations on the probe. There is a diagonal line with brightest color blocks, indicating that units with more similar firing responses across days tend to be physically close. This confirms our assumption that neurons are physically stable over time. Also notice that, on each column, there might be more than one bright block in the more distant clusters. We minimize the effect of distant units by constraining the feasible region during selection of reference units. There are also columns without bright yellow blocks. This happens because some units do not respond to the stimulus and those units are not included in the reference set.



**Figure 7 - figure supplement 2:** The Kilosort-good and reference unit counts for the animals AL031 and AL036, as shown for animal AL032 in Figure 7.



**Figure 7 - figure supplement 3:** The ratio of the count of reference units to KSgood units decreases for pairs of datasets with larger time intervals. However, the variability of the number of reference units is generally large for all time intervals.

# References

- [1] Carmena JM, Lebedev MA, Henriquez CS, Nicolelis MAL. Stable Ensemble Performance with Single-Neuron Variability during Reaching Movements in Primates. *J Neurosci*. 2005;25:10712–10716. <https://doi.org/10.1523/JNEUROSCI.2772-05.2005>.
- [2] Huber D, Gutnisky DA, Peron S, O'Connor DH, Wiegert JS, Tian L, et al. Multiple dynamic representations in the motor cortex during sensorimotor learning. *Nature*. 2012;484:473–478. <https://doi.org/10.1038/nature11039>.
- [3] Liberti WA, Markowitz JE, Perkins LN, Liberti DC, Leman DP, Guitchoynts G, et al. Unstable neurons underlie a stable learned behavior. *Nat Neurosci*. 2016;19:1665–1671. <https://doi.org/10.1038/nn.4405>.
- [4] Clopath C, Bonhoeffer T, Hübener M, , Rose T. Variance and invariance of neuronal long-term representations. *Phil Trans R Soc*. 2017;372. <https://doi.org/10.1098/rstb.2016.0161>.
- [5] Dhawale AK, Poddar R, Wolff SB, Normand VA, Kopelowitz E, Ölveczky BP. Automated long-term recording and analysis of neural activity in behaving animals. *eLife*. 2017;6:e27702. <https://doi.org/10.7554/eLife.27702>.
- [6] Jensen KT, Harpaz NK, Dhawale AK, Wolff SBE, , Ölveczky BP. Long-term stability of single neuron activity in the motor system. *Nat Neurosci*. 2022;25:1664–1674. <https://doi.org/10.1038/s41593-022-01194-3>.
- [7] Steinmetz NA, Aydin C, Lebedeva A, Okun M, Pachitariu M, Bauza M, et al. Neuropixels 2.0: A miniaturized high-density probe for stable, long-term brain recordings. *Science*. 2021;372:eabf4588. <https://doi.org/10.1126/science.abf4588>.
- [8] Luo TZ, Bondy AG, Gupta D, Elliott VA, Kopec CD, Brody CD. An approach for long-term, multi-probe Neuropixels recordings in unrestrained rats. *eLife*. 2020;9. <https://doi.org/10.7554/eLife.59716>.
- [9] Harris KD, Quiroga RQ, Freeman J, Smith SL. Improving data quality in neuronal population recordings. *Nature Neuroscience*. 2016;19:1165–1174. <https://doi.org/10.1038/nn.4365>.
- [10] Buzsáki G. Large-scale recording of neuronal ensembles. *Nature Neuroscience*. 2004;7:446–451. <https://doi.org/10.1038/nn1233>.
- [11] Brown EN, Kass RE, Mitra PP. Multiple neural spike train data analysis: state-of-the-art and future challenges. *Nature Neuroscience*. 2004;7:456–461. <https://doi.org/10.1038/nn1228>.
- [12] Quiroga RQ, Panzeri S. Extracting information from neuronal populations: information theory and decoding approaches. *Nature Reviews Neuroscience*. 2009;10:173–185. <https://doi.org/10.1038/nrn2578>.
- [13] Harris KD. Neural signatures of cell assembly organization. *Nature Reviews Neuroscience*. 2005;6:399–407. <https://doi.org/10.1038/nrn1669>.
- [14] Quiroga RQ, Nadasdy Z, Ben-Shaul Y. Unsupervised Spike Detection and Sorting with Wavelets and Superparamagnetic Clustering. *Neural Computation*. 2004;16:1661–1687. <https://doi.org/10.1162/089976604774201631>.
- [15] Chah E, Hok V, Della-Chiesa A, Miller JJH, O'Mara SM, , et al. Automated spike sorting algorithm based on Laplacian eigenmaps and k -means clustering. *J Neural Eng*. 2011;8:016006. <https://doi.org/10.1088/1741-2560/8/1/016006>.

- 500 [16] Pachitariu M, Steinmetz N, Kadir S, Carandini M, Harris KD.: Kilosort: realtime spike-sorting  
501 for extracellular electrophysiology with hundreds of channels. Preprint at [https://www.biorxiv.  
502 org/content/10.1101/061481v1](https://www.biorxiv.org/content/10.1101/061481v1).
- 503 [17] Carlson D, Carin L. Continuing progress of spike sorting in the era of big data. Current Opinion  
504 in Neurobiology. 2019;55:90–96. <https://doi.org/10.1016/j.conb.2019.02.007>.
- 505 [18] Jun JJ, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, et al. Fully integrated silicon  
506 probes for high-density recording of neural activity. Nature. 2017;551:232–236. [https://doi.  
507 org/10.1038/nature24636](https://doi.org/10.1038/nature24636).
- 508 [19] Hall NJ, Herzfeld DJ, Lisberger SG. Evaluation and resolution of many challenges of neural  
509 spike sorting: a new sorter. Journal of Neurophysiology. 2021;126:2065–2090. [https://doi.org/  
510 10.1152/jn.00047.2021](https://doi.org/10.1152/jn.00047.2021).
- 511 [20] Tolias AS, Ecker AS, Siapas AG, Hoenselaar A, Keliris GA, Logothetis NK. Recording Chron-  
512 ically From the Same Neurons in Awake, Behaving Primates. Journal of Neurophysiology.  
513 2007;98:3780–3790. <https://doi.org/10.1152/jn.00260.2007>.
- 514 [21] Swindale NV, Spacek MA. Spike sorting for polytrodes: a divide and conquer approach. Fron-  
515 tiers in Systems Neuroscience. 2014;8. <https://doi.org/10.3389/fnsys.2014.00006>.
- 516 [22] Bar-Hillel A, Spiro A, Stark E. Spike sorting: Bayesian clustering of non-stationary data. Journal  
517 of Neuroscience Methods. 2006;157:303–316. <https://doi.org/10.1016/j.jneumeth.2006.04.023>.
- 518 [23] Lee J, Mitelut C, Shokri H, Kinsella I, Dethe N, Wu S, et al. YASS: Yet Another Spike Sorter ap-  
519 plied to large-scale multi-electrode array recordings in primate retina. 2020;p. 10712–10716.  
520 Preprint at <https://www.biorxiv.org/content/10.1101/2020.03.18.997924v1>. [https://doi.org/10.  
521 1101/2020.03.18.997924](https://doi.org/10.1101/2020.03.18.997924).
- 522 [24] Chung JE, Magland JF, Barnett AH, Tolosa VM, Tooker AC, Lee KY, et al. A Fully Automated  
523 Approach to Spike Sorting. Neuron. 2017;95:1381–1394.e6. [https://doi.org/10.1016/j.neuron.  
524 2017.08.030](https://doi.org/10.1016/j.neuron.2017.08.030).
- 525 [25] Chung JE, Joo HR, Fan JL, Liu DF, Barnett AH, Chen S, et al. High-Density, Long-Lasting, and Multi-  
526 region Electrophysiological Recordings Using Polymer Electrode Arrays. Neuron. 2019;101:21–  
527 31.e5. <https://doi.org/10.1016/j.neuron.2018.11.002>.
- 528 [26] Vasil'eva LN, Badakva AM, Miller NV, Zobova LN, Roshchin VY, Bondar IV. Long-Term Record-  
529 ing of Single Neurons and Criteria for Assessment. Neuroscience and Behavioral Physiology.  
530 2016;46:264–269. <https://doi.org/10.1007/s11055-016-0227-8>.
- 531 [27] Rokni U, Richardson AG, Bizzi E, Seung HS. Motor Learning with Unstable Neural Representa-  
532 tions. Neuron. 2007;54:653–666. <https://doi.org/10.1016/j.neuron.2007.04.030>.
- 533 [28] Lewicki MS. A review of methods for spike sorting: the detection and classification of neural ac-  
534 tion potentials Michael S Lewicki. Network. 1998;9:R53–78. [https://doi.org/10.1088/0954-898X/  
535 9/4/001](https://doi.org/10.1088/0954-898X/9/4/001).
- 536 [29] Colonell J.: ecephys spike sorting. GitHub. [https://github.com/jenniferColonell/ecephys\\_spike\\_  
537 sorting](https://github.com/jenniferColonell/ecephys_spike_sorting).
- 538 [30] Cohen S. FINDING COLOR AND SHAPE PATTERNS IN IMAGES (Stanford University, Palo Alto,  
539 1999). 1999;.
- 540 [31] Bertrand NP, Charles AS, Lee J, Dunn PB, , Rozell CJ. Efficient Tracking of Sparse Signals via  
541 an Earth Mover's Distance Dynamics Regularizer. IEEE. 2020;27:1120–1124. [https://doi.org/10.  
542 1109/LSP.2020.3001760](https://doi.org/10.1109/LSP.2020.3001760).

- 538 [32] Boussard J, Varol E, Lee HD, Dethé N, Paninski L. Three-dimensional spike localization and  
539 improved motion correction for Neuropixels recordings. *NeurIPS Proceedings*. 2021;<https://doi.org/10.1101/2021.11.05.467503>.  
540
- 541 [33] Sauerbrei BA, Guo JZ, Cohen JD, Mischiati M, Guo W, Kabra M, et al. Cortical pattern generation  
542 during dexterous movement is input-driven. *Nature*. 2020;577:386–391. <https://doi.org/10.1038/s41586-019-1869-9>.  
543
- 544 [34] Stringer C, Pachitariu M, Steinmetz N, Carandini M, Harris KD. High-dimensional geometry  
545 of population responses in visual cortex. *Nature*. 2019;571:361–365. <https://doi.org/10.1038/s41586-019-1346-5>.  
546
- 547 [35] Böhm C, Lee AK.: Functional specialization and structured representations for space and time  
548 in prefrontal cortex. Preprint at <https://www.biorxiv.org/content/10.1101/2023.01.16.524214v1>.