

1 Multi-day Neuron Tracking in High 2 Density Electrophysiology 3 Recordings using EMD

4 **Augustine(Xiaoran) Yuan^{1,2}, Jennifer Colonell¹, Anna Lebedeva³, Michael Okun⁴,**
5 **Adam S. Charles^{2*}, Timothy D. Harris^{1,2*}**

*For correspondence:

adamsc@jhu.edu (ASC);

harrist@janelia.hhmi.org (TDH)

6 ¹Janelia Research Campus, Howard Hughes Medical Institute, USA; ²Department of
7 Biomedical Engineering, Center for Imaging Science Institute, Kavli Neuroscience
8 Discovery Institute, Johns Hopkins University, USA; ³Sainsbury Wellcome Centre,
9 University of Sheffield, UK; ⁴Department of Psychology and Neuroscience Institute,
10 Howard Hughes Medical Institute, USA

11
12 **Abstract** Accurate tracking of the same neurons across multiple days is crucial for studying
13 changes in neuronal activity during learning and adaptation. New advances in high density
14 extracellular electrophysiology recording probes, such as Neuropixels, provide a promising
15 avenue to accomplish this goal. Identifying the same neurons in multiple recordings is, however,
16 complicated by non-rigid movement of the tissue relative to the recording sites (drift) and loss of
17 signal from some neurons. Here we propose a neuron tracking method that can identify the
18 same cells independent of firing statistics, which are used by most existing methods. Our method
19 is based on between-day non-rigid alignment of spike sorted clusters. We verified the same cell
20 identify using measured visual receptive fields. This method succeeds on datasets separated
21 from one to 47 days, with an 84% average recovery rate.

23 1 Introduction

24 The ability to longitudinally track neural activity is crucial to understanding central capabilities and
25 changes of neural circuits that operate on long time-scales, such as learning and plasticity,¹⁻⁴ mo-
26 tor stability,^{1,5,6} etc. We seek to develop a method capable of tracking single units regardless of
27 changes in functional responses for the duration of an experiment spanning one to two months.

28 High-density multi-channel extracellular electrophysiology (ephys) recording devices enable
29 chronic recordings over large areas over days-to-months.⁷ Such chronic recordings make possi-
30 ble experiments targeted at improving our understanding of neural computation and underlying
31 mechanisms. Examples include perceptual decision making, exploration and navigation.⁸⁻¹³ Elec-
32 trode arrays with hundreds to thousands of sites, for example Neuropixels, are now used exten-
33 sively to record the neural activity of large populations stably and with high spatio-temporal reso-
34 lution, capturing hundreds of neurons with single neuron resolution.^{9,10} Moreover, ephys retains
35 the higher time resolution needed for single spike identification, as compared with calcium imaging
36 that provides more spatial cues with which to track neurons over days.

37 The first step in analyzing ephys data is to extract single neuron signals from the recorded volt-
38 age traces, i.e., spike sorting. Spike sorting identifies individual neurons by grouping detected ac-
39 tion potentials using waveform profiles and amplitudes. Specific algorithms include principal com-

40 ponents based methods¹⁴ and,¹⁵ and template matching methods, for example, Kilosort.^{9,11,16,17}
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.¹⁸ 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.⁷ 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,¹⁹ 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.^{9,20-22} Another challenge is that popular template-
55 matching-based spike sorting methods usually involve some randomness in template initializa-
56 tion.^{16,23,24} As a result, action potentials can be assigned into clusters differently, and clusters can
57 be merged or separated differently across runs.

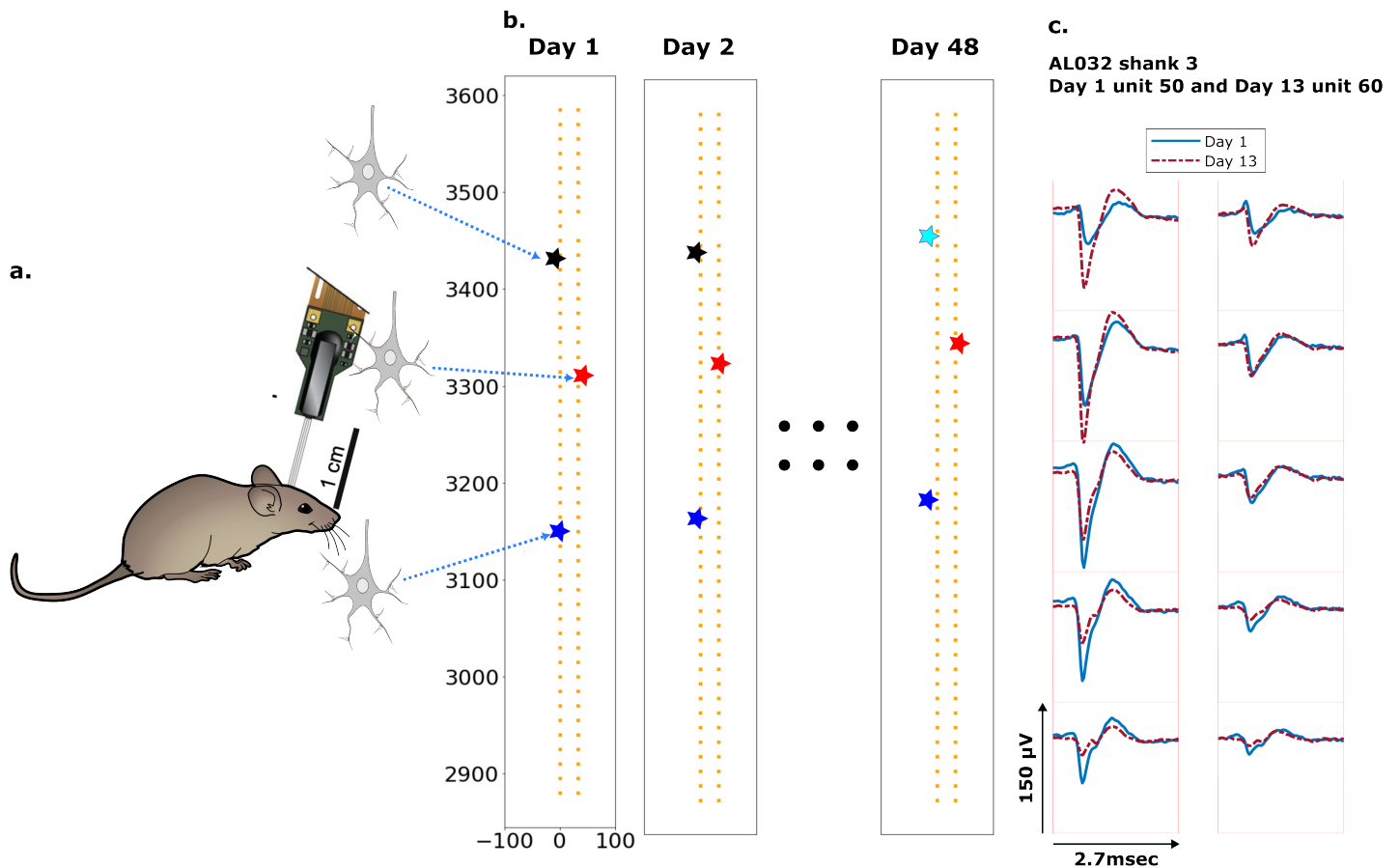


Fig. 1: Schematic depiction of drift: a. Mice were implanted with a 4-shank Neuropixels 2.0 probe in visual cortex area V1. b. Each colored star represents the location of a unit recorded on the probe. In this hypothetical case, the same color indicates unit correspondence across days. The black unit 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,²⁵ action potential shape correlation and inter-spike interval histogram (ISI)
 60 shape.²⁶ When neuronal representations change, e.g., during learning¹⁻³ or representational drift,²⁷
 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,²⁸ instead of firing patterns.

64 As an alternative method, Steinmetz et al.⁷ 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.

68 2 Results

69 2.1 Procedure

70 Our datasets consist of multiple recordings taken from three mice (**Figure 7a**) over 2 months. The
71 time gap between two recordings ranges from two to 25 days. Each dataset is spike-sorted individu-
72 ally with a standard Kilosort 2.5 pipeline. The sorting results, including unit assignment, spike times,
73 etc. are used as input for our method (post-processed using ecephys spike sorting pipeline²⁹) (Sec.
74 4.3). To ensure the sorting results are unbiased, we performed no manual curation. As the clusters
75 returned by Kilosort can vary in quality, we only considered the subset of units labeled as 'good' by
76 Kilosort, here referred to as KSgood units (Sec. 4.4). KSgood units are mainly determined by the
77 amount of inter-spike-interval violations and are believed to represent a single unit.¹⁶

78 Our overall strategy is to run spike-sorting once per session, and then to generate a unit-by-unit
79 assignment between pairs of datasets. When tracking units across more than two sessions, two
80 strategies are possible: match all ensuing sessions to a single session (e.g., the first session) (Sec.
81 2.2 and Sec. 4.2), or match consecutive pairs of sessions and then trace matched units through all
82 sessions (Sec. 2.4).

83 We refer to the subset of KSgood units with strong and distinguishable visual responses in
84 both datasets of a comparison as reference units (See Sec. 4.4 for details). Similar to Steinmetz et
85 al.⁷ we validated our unit matching of those reference units using visual receptive field similarity.
86 Finally, we showed that trackable units with strong visual responses are qualitatively similar to
87 those without (**Figure S1 to Figure S5**).

88 To provide registration between pairs of recordings, we used the Earth Mover's Distance (EMD).^{30,31}
89 We use a feature space consisting of a geometric distance space and a waveform similarity space,
90 to address both rigid and non-rigid neuron motion. The EMD finds matches between objects in
91 the two distributions by minimizing the overall distances between the established matches (Sec.
92 4.1.1).

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

106 For each unit, the location is determined by fitting the peak to peak amplitudes on the 10 sites
107 nearest the site with peak signal, based on the triangulation method in³² (Sec. 4.1.2). The waveform
108 distance is an L2 norm between two spatial-temporal waveforms that spans 22 channels and 2.7
109 msec (Sec. 4.1.2). Physical unit distances provide a way to maintain the internal structure and
110 relations between units in the EMD. Waveform similarity metrics will distinguish units in the local
111 neighborhood and likely reduce the effect of new and missing units (**Figure S6**).

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

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

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

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

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

122 Second, we tracked units between consecutive datasets and summarized and analyzed the
123 waveforms, unit locations, firing rates and visual responses (see **Figure S1** to **Figure S5** for details)
124 of all tracked chains, i.e. units which can be tracked across at least three consecutive datasets.

125 **2.2 Measuring rigid drift using the EMD**

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

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

141 As the EMD is an optimization algorithm with no biological constraints, it assigns matches to all
142 units in the smaller dataset regardless of biophysical plausibility. As a result, some of the assigned
143 matches may have unrealistically long distances. A distance threshold is therefore required to
144 select correct pairs. For the illustration in **Figure 2**, the threshold is set to $15\mu m$, which is chosen to
145 be larger than most of the z-shifts observed in our experimental data. The threshold value will be
146 refined later by distribution fitting (**Figure S2**). In **Figure 2** all of the sub-threshold (short) distances
147 belong to upward pairs (**Figure 2b** and **c**, red solid arrows), showing that the EMD can detect the
148 homogeneous movement direction and the amount of imposed drift.

149 When determining matched reference units from visual response data, we require that units
150 be spatially nearby (within $30\mu m$) as well as having similar visual responses. After correcting for
151 drift, we find that we recover more reference units (**Figure S7**), indicating improved spatial match
152 of the two ensembles. This improved recovery provides further evidence of the success of the drift
153 correction.

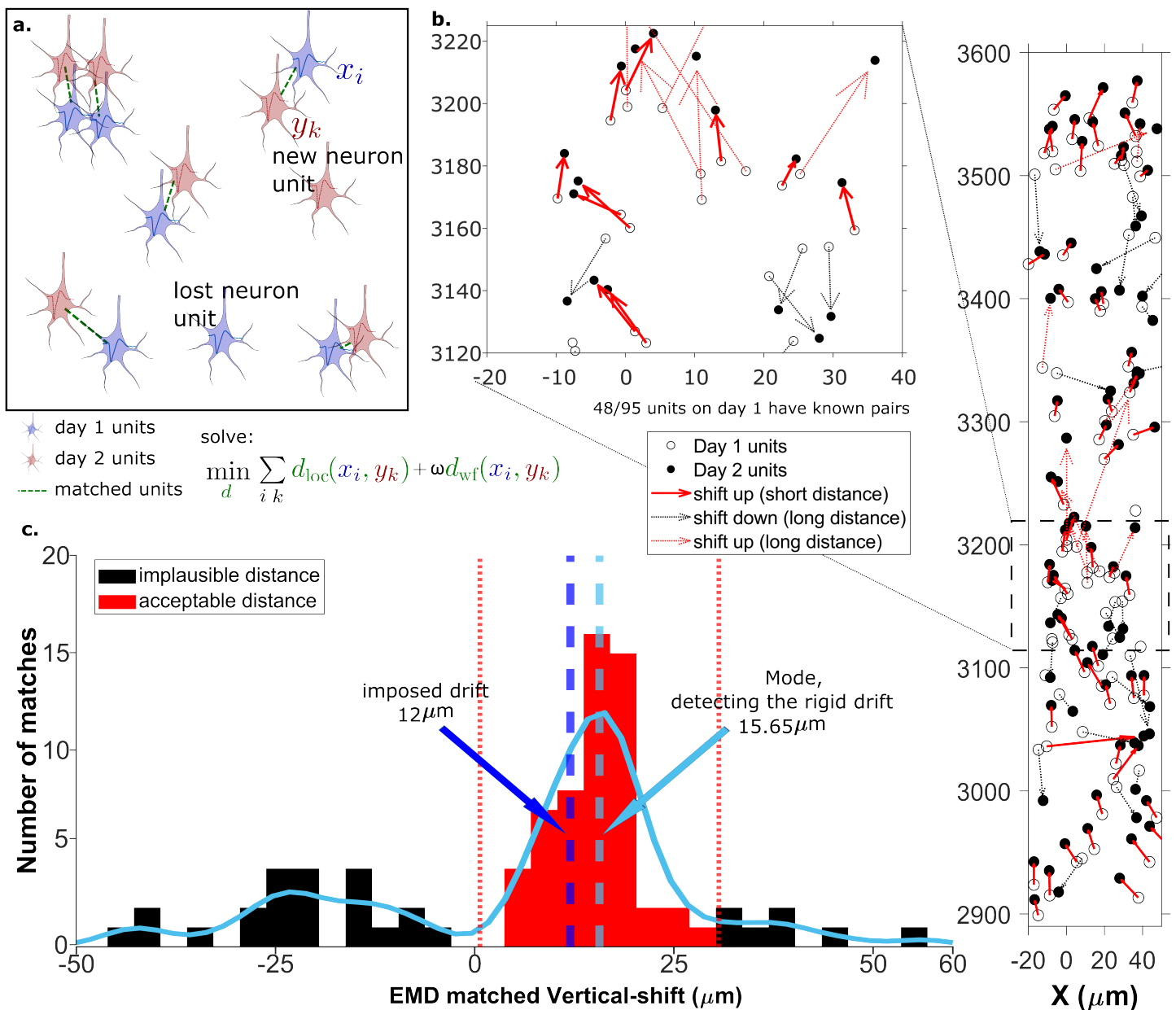


Fig. 2: The EMD can detect the displacement of single units: a. Schematic of EMD unit matching. Each blue unit in day 1 is matched to a red unit in day 2. Dashed lines indicate the matches to be found by minimizing the weighted sum of physical and waveform distances. b. Open and filled circles show positions 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 μm . 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.

154 2.3 A vertical distance threshold is necessary for accurate tracking

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

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

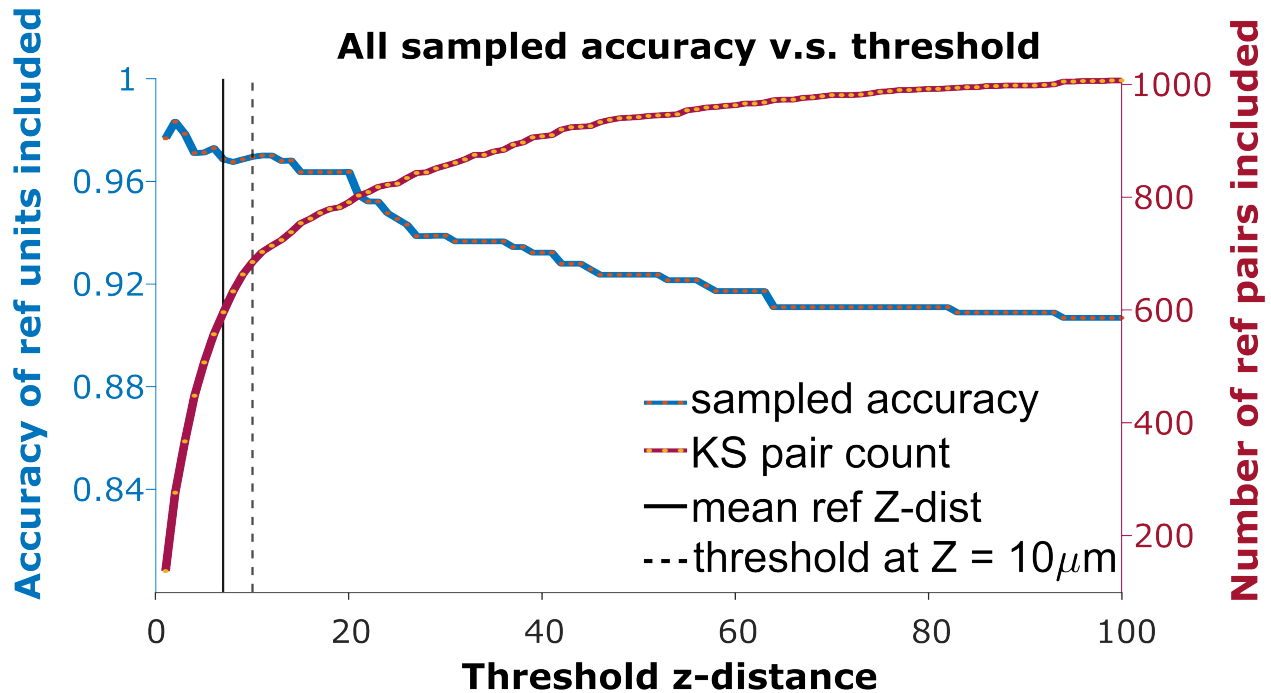


Fig. 3: The ROC curve of matching accuracy vs. distance. The blue curve shows the accuracy for reference units. The red line indicates the number of reference units included. The solid vertical line indicates the average z distance across all reference pairs in all animals ($z = 6.96\mu\text{m}$). The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

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

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

175 for a mixture of correct and incorrect pairs.

176 We assume that the distribution of z-distances $P(\Delta)$ for reference units is the conditional prob-
177 ability $P(\Delta | H)$; that is, we assume all reference units are true hits. The distribution of z-distances
178 for all KSgood units $P(\Delta)$ includes both hits and false positives. The distance distribution of false
179 positives is the difference between the two (Sec. 8.4, **Equation 6**).

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

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

184 We fit the KSgood distribution to **Equation 3** to extract the individual distribution parameters and
185 the fraction of true hits (f). The full distribution can then be integrated up to any given z-threshold
186 value to calculate the false positive rate. (**Figure 4a**, top panel, see Sec. 8.4 for details).

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

191 Note that this threshold eliminates most of the known false positive matches of reference pairs
192 (**Figure 4b**, red fraction) at the cost of recovering fewer correct pairs (**Figure 4b**, green bars). The re-
193 covery rate varies from day to day; datasets separated by longer times tend to have higher tracking
194 uncertainty (**Figure S10**).

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

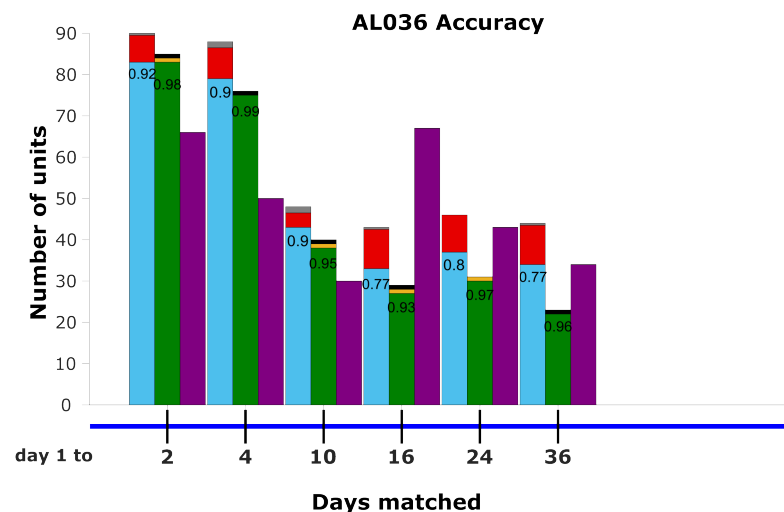
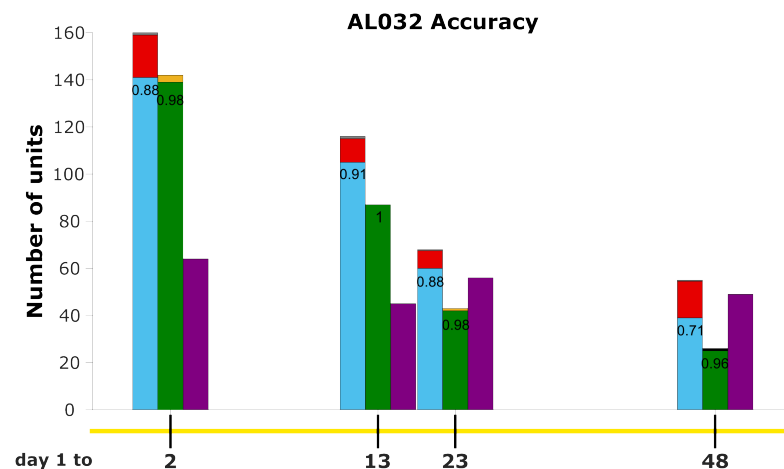
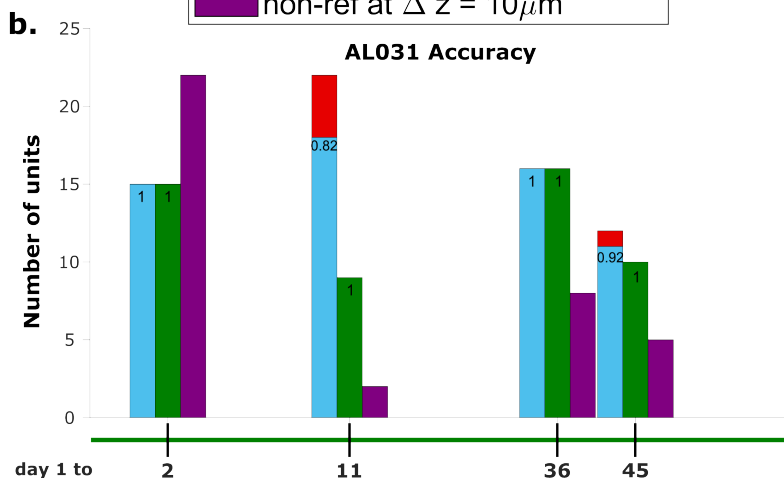
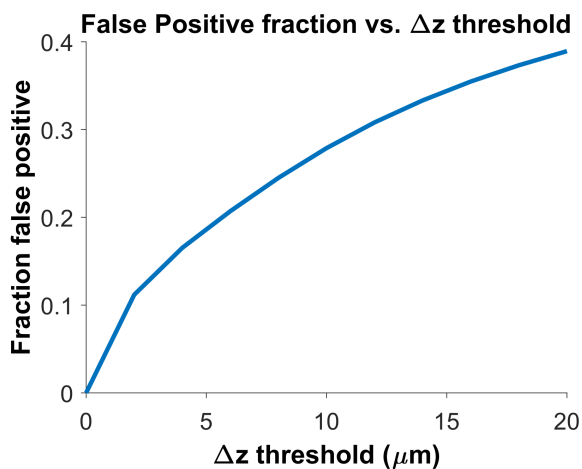
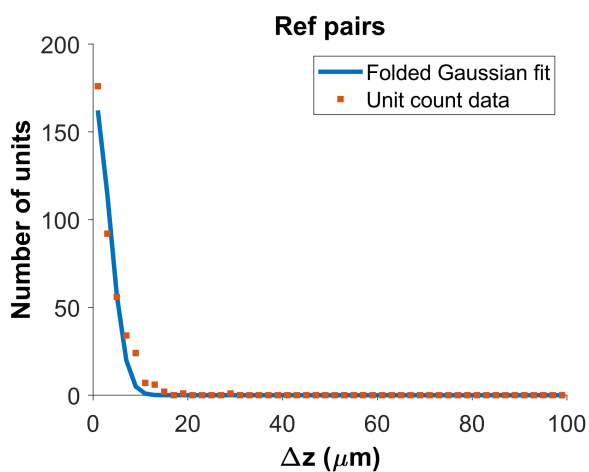
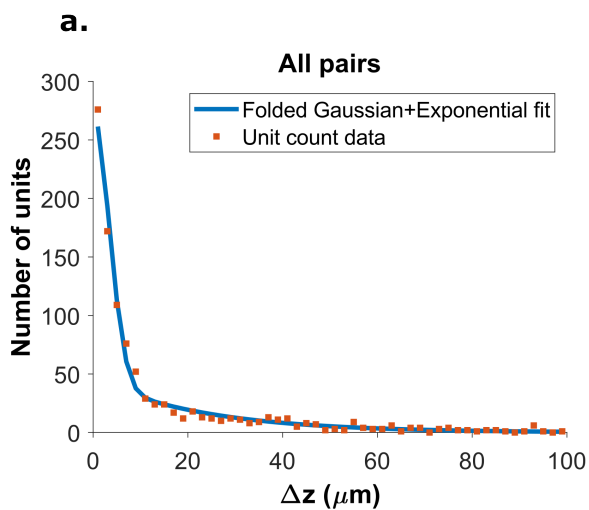
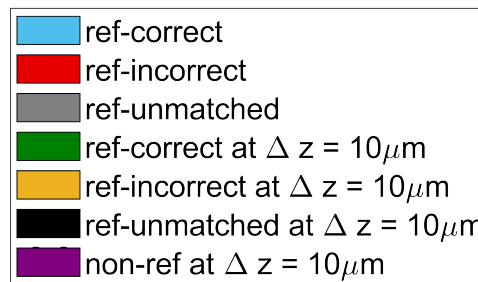


Fig. 4: Recovery rate, accuracy and putative pairs: a. The histogram distribution fit for all KS-good units (top) and reference units alone (middle). False positives for reference units are defined as units matched by EMD but not matched when using receptive fields. The false positive fraction for the 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 = $10\mu m$.

198

199 **2.4 Units can be tracked in discontinuous recordings for 48 days**

200 To assess long-term tracking capabilities, we tracked neurons across all datasets for each mouse.
201 **Figure 5** shows a survival plot of the number of unit chains successfully tracked over all durations.
202 All units in the plot can be tracked across at least three consecutive datasets, a chain as the term
203 is used here. We categorized all trackable unit chains into three types: reference chains, mixed
204 chains and putative chains. Reference chains have receptive field information in all datasets. Pu-
205 tative chains have no reference information in any of the datasets. Mixed units have at least one
206 dataset with no receptive field information. There are 133 reference chains, 135 mixed chains and
207 84 putative chains across all the subjects. Among them, 46 reference, 51 mixed, and 9 putative
208 units can be followed across all datasets. We refer to them as fully trackable units. One example
209 trackable unit in each group is shown in **Figure 6**, **Figure S16**, and **Figure S17**.

Summary of duration of neuron tracked across all subjects

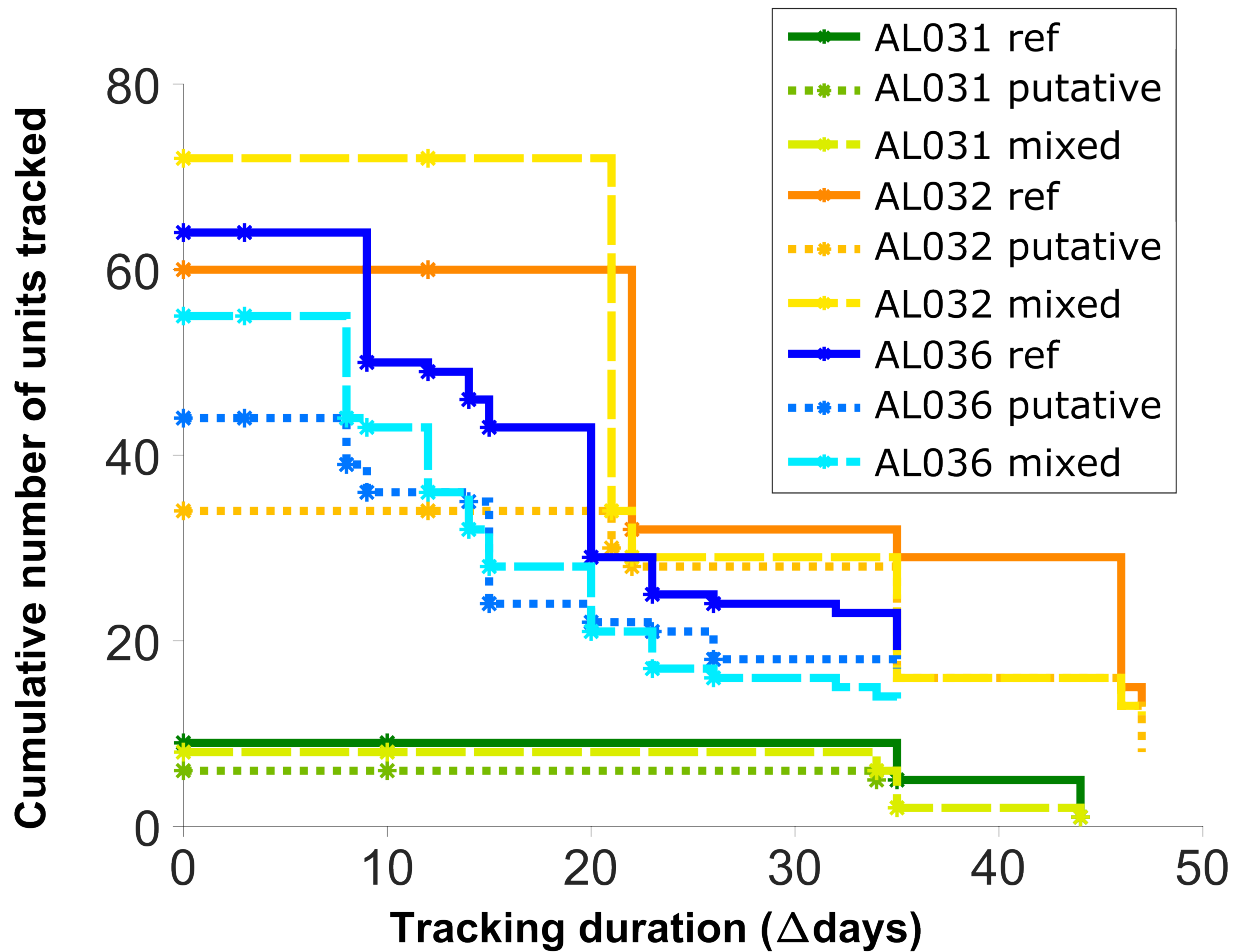


Fig. 5: Number of reference units (deep blue, dark orange and green for different subjects), putative (medium green, medium orange and blue) units, and mixed units (light green, yellow, and light blue) tracked for different durations. The loss rate is similar for different chain types in the same subject. 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.

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

214 The spatial-temporal waveform similarity is measured by the L2 distance between waveforms
215 (Sec. 4.1.2). A Kruskal-Wallis test is performed on the magnitude of L2 change between all pairs
216 of matched waveforms among the three groups. There is no statistical difference in the waveform
217 similarity in reference, putative, and mixed units ($H = 0.59$, $p = 0.75$) (Figure S1). There is no signifi-
218 cant difference in the physical distances of units per dataset ($H = 1.31$, $p = 0.52$) (Figure S2, bottom
219 panel), nor in the location change of units ($H = 0.23$, $p = 0.89$) (Figure S2, top panel).

220 Firing rate is characterized as the average firing rate fold change of each unit chain, with firing
221 rate of each unit in each dataset normalized by the average firing rate of that dataset. There is no
222 difference in the firing rate fold change in the three groups of units ($H = 1$, $p = 0.6$) (**Figure S3**).
223 The receptive field similarity between units in different datasets is described by visual finger-
224 print (vfp) correlation and Peristimulus Time Histogram (PSTH) correlation between units, and the
225 similarity score, the sum of the two correlations (Sec. 4.4). The change in vfp between matched
226 units is similar among the three groups ($H = 2.23$, $p = 0.33$). Similarly, the change in PSTH is not
227 different among the three groups ($H = 1.61$, $p = 0.45$) (**Figure S4**).

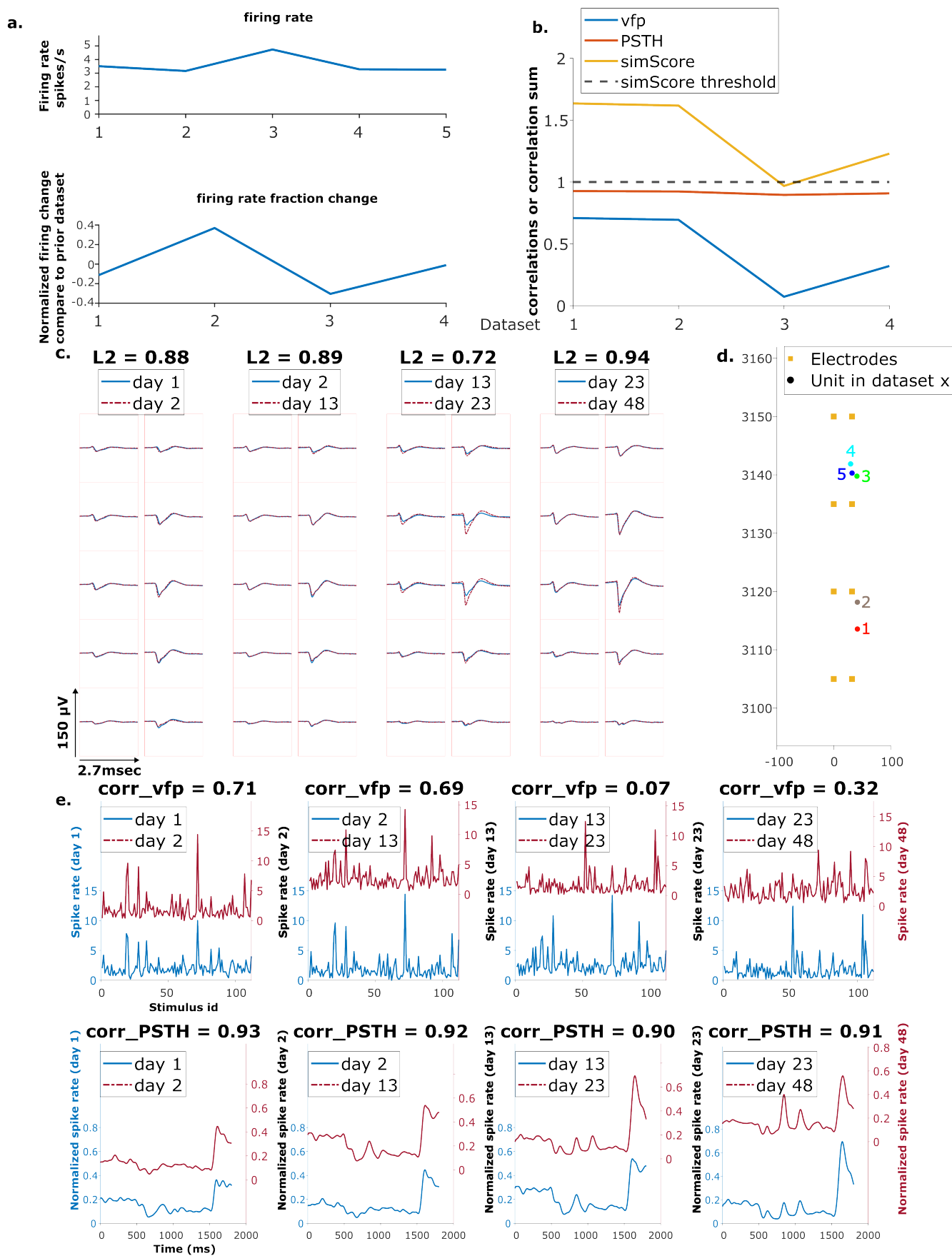


Fig. 6: Example mixed chain: a. Above: Firing rates of this neuron on each day (Day 1, 2, 13, 23, 48). Below: Firing rate fractional change compared to the previous day. b. Visual response similarity (yellow line), PSTH correlation (orange line), and visual fingerprint correlation (blue line). The similarity 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.

228

229 3 Discussion

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

242 show the capability of our method to track neurons with no specific tuning preference (**Figure S16**).
243 The method includes means to identify dataset pairs with very large drift. In our data, we can
244 detect large drift because such datasets have very few reference units, and significantly different
245 EMD cost (Sec. 8.6). For example, datasets 1 and 2 in animal AL036 have very few reference units
246 compared to other datasets (see **Figure S11**, AL036). This observation is consistent with the overall
247 relationship between the EMD cost and recovery rate (**Figure S12**). Datasets with higher cost tend to
248 have lower unit recovery rate and higher variation in recovery rates. Therefore, these two datasets
249 were excluded in the tracking analysis.

250 Our validation relies on identifying reference units. The reference unit definition has limitations.
251 The similarity score is largely driven by PSTHs (**Figure 6**, **Figure S11**), the timing of stimulus triggered
252 response, rather than vfp, the response selectivity. As a result, a single neuron can be highly corre-
253 lated, i.e. similarity score greater than 1, with more than 20 other neurons. For example, in subject
254 AL032 shank 2, one neuron on day 1 has 22 highly correlated neurons on day 2, 4 of which are
255 also within the distance of $30\ \mu\text{m}$. Non-reference units may also have very similar visual responses:
256 we note that 33 (5 putative neurons and 28 mixed neurons) out of 106 trackable neurons have
257 a similarity score greater than 1 even for days with no reference unit assignment. Coincidentally
258 similar visual responses could potentially contribute to inaccurate assignment of reference units
259 and irregularity in trackable unit analysis. These errors would reduce the measured accuracy of
260 the EMD matching method; since the accuracy is very high (**Figure 4**), the impact of mismatches is
261 low.

262 We note that the ratio of reference units over KGood units decreases as recordings are further
263 separated in time (**Figure S13**). This reduction in fraction of reference units might be partially due
264 to representational drift as well as the fact that the set of active neurons are slightly different in
265 each recording. The visual fingerprint similarity of matched neurons decreased to 60% after 40
266 days (see reference 7 supplement).

267 We developed the new tracking algorithm based on an available visual cortex dataset, and used
268 a prominent sorting algorithm (Kilosort 2.5) to spikesort the data. We had reference data to assess

269 the success of the matching and tune parameters. Applying our algorithm in other brain areas and
270 with other sorters may require parameter adjustment. Evaluation of the results in the absence of
271 reference data requires a change to the fitting procedure.

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

277 When reference information (unit matches known from receptive fields or other data) is avail-
278 able, the procedure outlined in section 8.4 can be followed. In that case, the distribution of z-
279 distances of known pairs is fit to find the width of the distribution for correct matches. That pa-
280 rameter is then used in the fit of the z-distance distribution of all pairs to *Equation 3*. Integrating
281 the distributions of correct and incorrect pairs yields the false positive rate vs. z-distance, allowing
282 selection of a z-distance threshold for a target false positive rate.

283 In most cases, reference information is not available. However, the z-distance distributions
284 for correct and incorrect pairs can still be estimated by fitting the distribution of all pairs. In sec-
285 tion 8.4, *Figure S9* we show the results of fitting the z-distribution of all pairs without fixing the
286 width of the distribution of correct matches. The result slightly underestimates this width, and the
287 estimated false positive rate increases. This result is important because it suggests the accuracy
288 estimate from this analysis will be conservative. We detail the procedure for fitting the z-distance
289 distribution Methods section (Alg. 2).

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

299 4 Methods

300 Our neuron tracking algorithm uses the Earth Mover's Distance (EMD) optimization algorithm. The
301 minimized distance is a weighted combination of physical distance and 'waveform distance': the al-
302 gorithm seeks to form pairs that are closest in space and have the most similar waveforms. We test
303 the performance of the algorithm by comparing EMD matches to reference pairs determined from
304 visual receptive fields (Sec. 4.4). We calculate two performance metrics. The 'recovery rate' is the
305 percentage of reference units that are correctly matched by the EMD procedure. The 'accuracy' is
306 the percentage of correctly matched reference units that pass the z-distance threshold (*Figure 4a*).
307 'Putative units' are units matched by the procedure which do not have reference receptive field
308 information. 'Chains' are units that can be tracked across at least three consecutive datasets. The
309 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:³²

$$f(x, y, z) = \sum_{c \in C_{u_n}} (V_{ptp_c} - \frac{1}{\sqrt{(x-x_c)^2 + (z-z_c)^2 + y^2}})^2$$

 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

310 **4.1 Algorithm**

311 4.1.1 Earth Mover's Distance

312 The EMD is an optimization-based metric developed in the context of optimal transport and mea-
 313 suring distances between probability distributions. It frames the question as moving dirt, in our
 314 case, units from the first dataset, into holes, which here are the neural units in the second dataset.
 315 The distance between the "dirt" and the "holes" determines how the optimization program will pri-
 316 oritize a given match. Specifically, the EMD seeks to minimize the total work needed to move the
 317 dirt to the holes, i.e., neurons in day 1 to day 2, by solving for a minimum overall effort, the sum of
 318 distances.^{30,31}

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

319 in which $d_{loc} \in \mathcal{D}^3$ is the three-dimensional physical distance between a unit from the first
 320 dataset x_i , and a unit from the second dataset y_k . $d_{wf} \in \mathcal{D}^1$ is a scalar representing the similar-
 321 ity between waveforms of units x_i and y_k . ω is a weight parameter that was tuned to maximize the
 322 recovery rate of correctly matched reference units. F is the vector of matched objects between the
 323 two datasets (See **Figure S14** for details about selecting weight).

324 The EMD has three benefits:

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

331 However, since the EMD is an optimization method with no assumptions about the biological prop-
 332 erties of the data, it makes all possible matches. We therefore added a threshold on the permissible
 333 z-distance to select physically plausible matches. Supplement **Figure S14** shows the recovery rate
 334 change as a function of weight parameters to combine neuron location and waveform metrics into
 335 a distance matrix.

336 4.1.2 Calculating the EMD distance metric

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

344 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
 345 between measured amplitudes V_{PTP} and amplitudes estimated with locations $\frac{\alpha}{\sqrt{(x-x_c)^2+(z-z_c)^2+y^2}}$:

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

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

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

355 4.1.3 Between-session Drift Correction

356 Based on previous understanding of the drift in chronic implants, we assumed that the majority
357 of drift occurs along the direction of the probe insertion, i.e. vertical z-direction. This rigid drift
358 amount is estimated by the mode of the z-distance distribution of the EMD assigned units using
359 a normal kernel density estimation implemented in MATLAB. We only included KSGood units.¹⁶
360 The estimated drift is then applied back to correct both the reference units and the EMD distance
361 matrix by adjusting the z coordinates of the units. A post-correction reference set is compared
362 with the post-correction matching results for validation.

363 4.2 Determining Z Distance Threshold

364 Determining the z-distance threshold to achieve a target false positive rate requires estimating
365 the widths of the z-distance distributions of correct and incorrect pairs. If reference data is avail-
366 able, the z-distance distribution of the known correct pairs should be fit to a folded Gaussian as
367 described in 8.4. The width of the folded Gaussian, which is the error in determination of the z-
368 positions of units, is then fixed in the fit of the z-distribution of all pairs found by the algorithm
369 outlined in Algorithm 4.1.1. If no reference data is available, the width of the distribution of correct
370 pairs is determined by fitting the z-distance distribution of all pairs to **Equation 3** with the folded
371 Gaussian width as one of the parameters. This procedure is detailed in Algorithm 2. We show two
372 examples of model fitting without reference information in section **Figure S9**.

Algorithm 2 Determining an appropriate z distance threshold

Input: Z distances of all matched units, target false positive rate, width σ 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 σ of the distribution of correct pairs is known
from reference data, fix at that value. Otherwise, include in the fit parameters. (See section
8.4 for details). 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; σ = 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
(σ) 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 S11 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 generates
the false positive rate vs. z-distance threshold, as shown in **Figure S9**. (Code available at:
https://github.com/AugustineY07/Neuron_Tracking/tree/main/Pipeline/Plot/Fit)

Output: σ (uncertainty of position estimation), threshold at the target false positive rate

373 4.3 Dataset

374 The data used in this work are recordings collected from two chronically implanted NP 2.0 four-
375 shank probes and one chronically implanted one-shank NP 2.0 probe in the visual cortex of three
376 head fixed mice (**Figure 7b**, see Steinmetz et al.⁷ for experiment details). The recordings were taken

377 while 112 visual stimuli were shown from three surrounding screens (data from Steinmetz et al.⁷
378 Supplement Section 1.2). The same bank of stimuli was presented five times, with order shuffled.
379 The 4-shank probes had the 384 recording channels mapped to 96 sites on each shank.

380 We analyzed 65 recordings, each from one shank, collected in 17 sessions (5 sessions for animal
381 AL031, 5 sessions for animal AL032, and 7 sessions for animal AL036). The time gap between
382 recordings ranges from one day to 47 days (**Figure 7a**), with recording durations ranging from
383 1917 to 2522 seconds. The sample rate is 30kHz for all recordings. There are a total of 2958 KGood
384 units analyzed across all animals and shanks, with an average of 56 units per dataset (**Figure 7d**
385 and **Figure S15**).

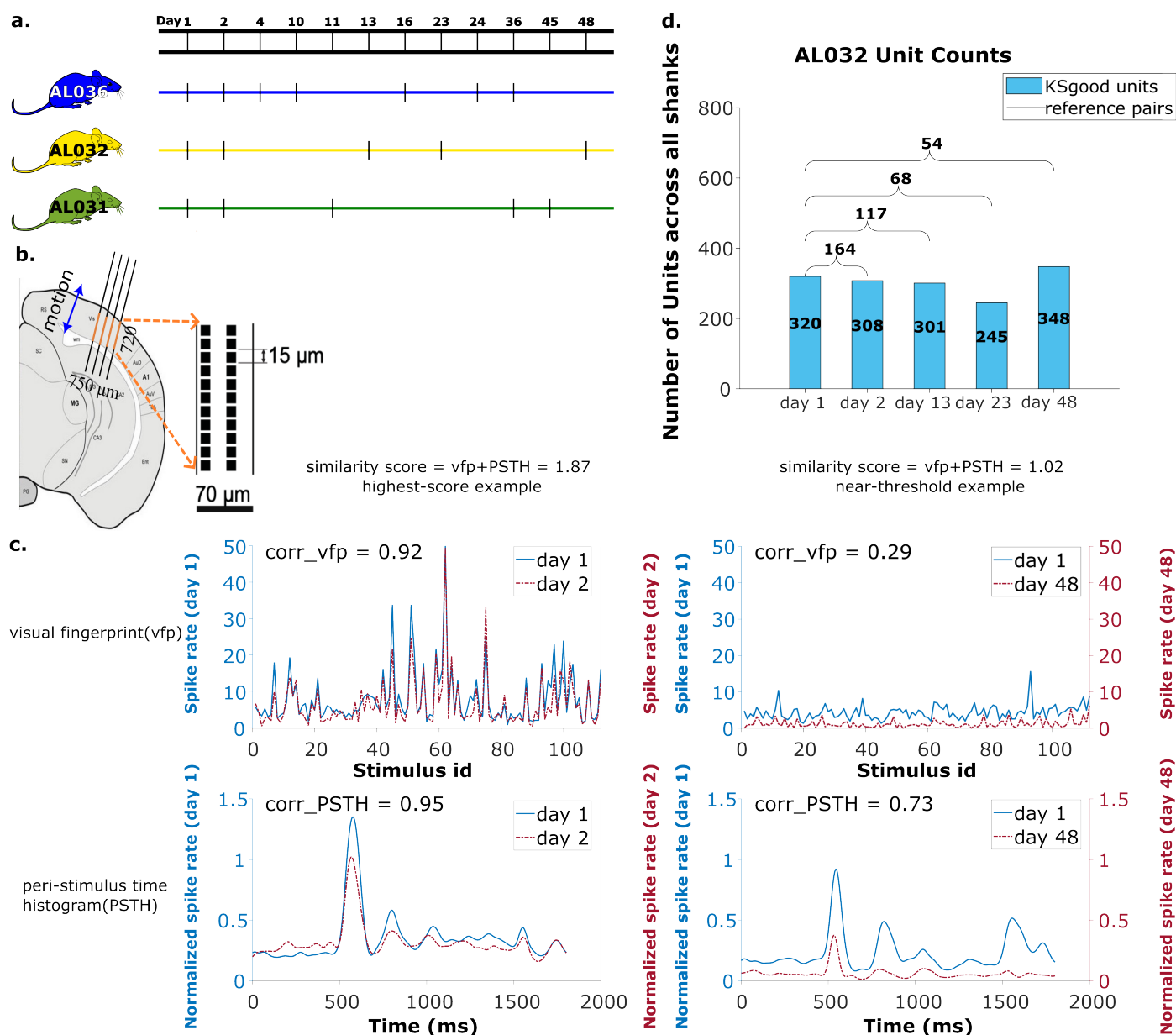


Fig. 7: Summary of dataset: a. The recording intervals for each animal. A black dash indicates one recording on that day. b. All recordings are from visual cortex V1 with a 720 μm section of the probe containing 96 recording sites. The blue arrow indicates the main drift direction. c. Examples of visual 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.

386 4.4 Reference set

387 To track clusters across days, Steinmetz et al.⁷ concatenated two recording sessions and took
 388 advantage of the within-recording drift correction feature of Kilosort 2.0 to extract spikes from
 389 the two days with a common set of templates. They first estimated the between session drift of
 390 each recording from the pattern of firing rate and amplitude on the probe and applied a position
 391 correction of an integer number of probe rows (15 μm for the probes used). Then two corrected
 392 recordings were concatenated and sorted as a single recording. This procedure ensured that the

393 same templates are used to extract spikes across both recordings, so that putative matches are
394 extracted with the same template. A unit from the first half of the recording is counted as the same
395 neuron if its visual response is more similar to that from the same cluster in the second half of the
396 recording than to the visual response of the physically nearest neighbor unit. Using this procedure
397 and matching criteria, 93% of the matches were correct for recordings < 16 days apart, and 85%
398 were correct for recordings from 3-9 weeks (See Steinmetz et al.,⁷ Fig. 4). In addition, although
399 mean fingerprint similarity decreases for recordings separated by more than 16 days, this decline
400 is only 40% for the same unit recorded from 40 days apart (see Steinmetz et al.⁷ Supplement S3).
401 This procedure, while successful in their setting, was limited to the use of integral row adjustments
402 of the data for between-session drift correction and relied on a customized version of Kilosort 2.0.
403 Although up to three recordings can be sorted together, they must come from recording sessions
404 close in time. In addition, a separate spike sorting session needs to be performed for every pair of
405 recordings to be matched, which is time consuming and introduces extra sorting uncertainty.

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

415 Following Steinmetz et al.,⁷ the similarity score between two neurons is the sum of the corre-
416 lation of the PSTH and the correlation of the vfp across two sessions. The two correlations have
417 values in the range (-1,1), and the similarity score ranges from (-2, 2).

418 The pool of reference units is established with three criteria: 1) The visual response similarity
419 score of the pair, as described above, is greater than 1 and their physical distance, both before and
420 after drift correction, is smaller than 30 μm . We impose the 30 μm threshold on both pre- and post-
421 correction data because the drift is relatively small in our case, and we can reduce false positives
422 by constraining the reference units to be in a smaller region without losing units. In general, one
423 could apply the threshold only on corrected data (after drift correction). 2) A Kruskal-Wallis test
424 is applied on all trials of the vfps to ensure the triggered response to the stimulus is significantly
425 distinguishable from a flat line. 3) Select units from each recording that meet the good criteria in
426 Kilosort. Kilosort assigns a label of either single-unit (good) or multi-unit (MUA) to all sorted clusters
427 based on ISI violations.¹⁶ This step aims to ensure included units are well separated. If there are
428 multiple potential partners for a unit, the pair with the highest similarity score is selected as the
429 reference unit. The complete pool of reference units includes comparisons of all pairs of recordings
430 for each shank in each animal. The portion of units with qualified visual response ranges from 5%
431 to 61%, depending on the time gap between datatets (**Figure S13**). Overall, these reference units
432 made up 29% of all KSGood units (**Figure S15**) across all three animals in our dataset. **Figure 7c**
433 shows examples of visual responses from a high similarity reference unit and a reference unit with
434 similarity just above threshold.

435 5 Code sharing

436 All code used can be accessed at: https://github.com/AugustineY07/Neuron_Tracking.

437 6 Acknowledgments

438 This work was supported supported in part by NIH grant U01 NS115587. We thank Claudia Böhm
439 and Albert Lee for allowing us to use their data in **Figure S9**.

440 **7 Declaration of interests**

441 The authors declare no competing interests.

442 8 Supplement

443 8.1 Trackable units statistics

444 To show that trackable reference, putative, and mixed units are qualitatively similar, we summa-
445 rized the median, maximum and minimum change of firing rate, visual receptive field, and loca-
446 tion in the box plots in *Figure S1* to *Figure S5*. A Kruskal-Wallis test performed for each feature
447 suggested no difference among the three groups (see Sec. 2.4 for details).

Waveform L2 change per dataset

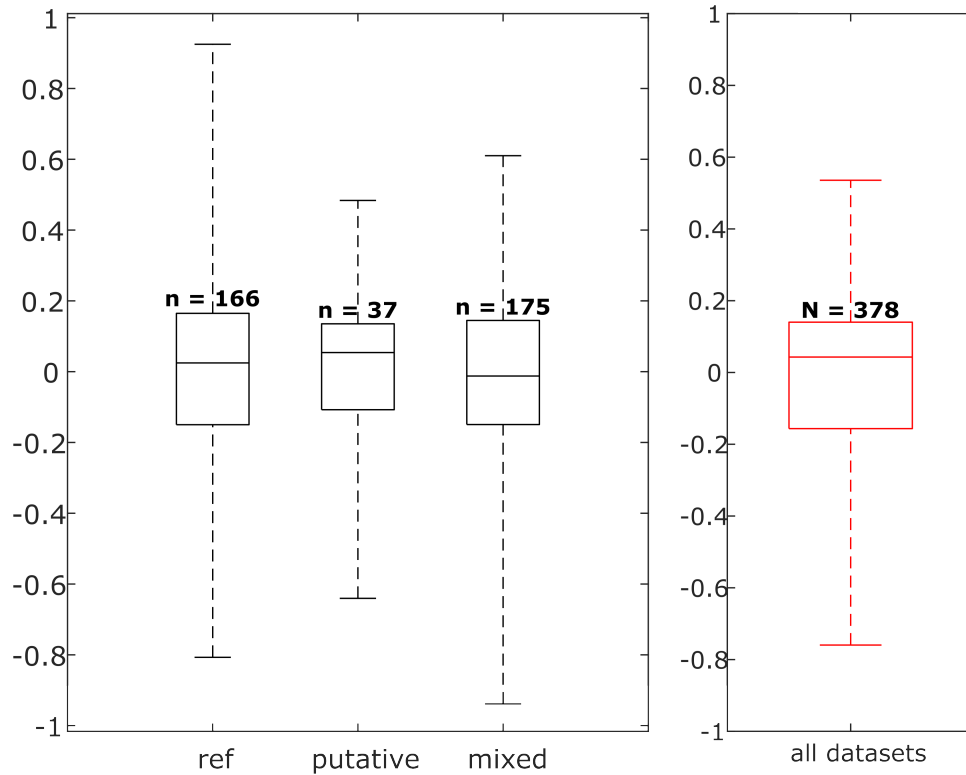
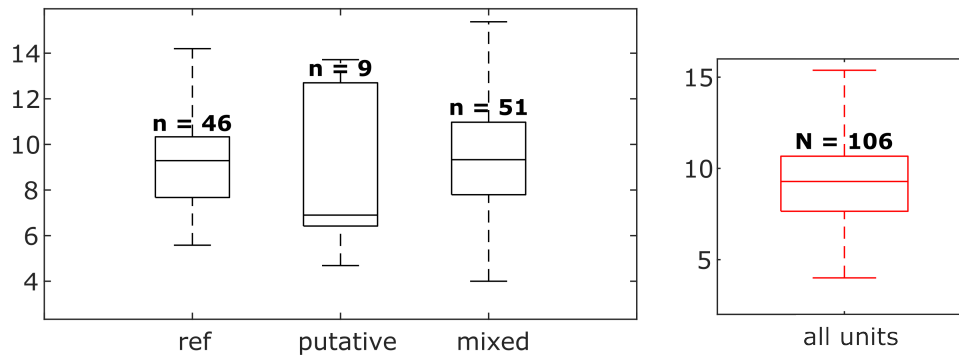


Fig. S1: Distribution of waveform L2 similarity change per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of unit comparisons, i.e. (number of units) \times (number of datasets - 1).

Average location change per unit



Location change per dataset

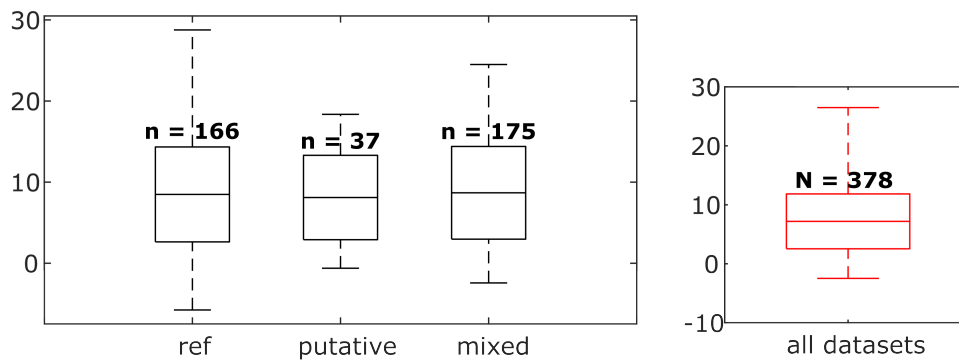


Fig. S2: Distributions of individual unit location changes over whole chains (top) and unit location changes between pairs of datasets (bottom), for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum 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) \times (number of datasets - 1).

Average firing rate change ratio per unit

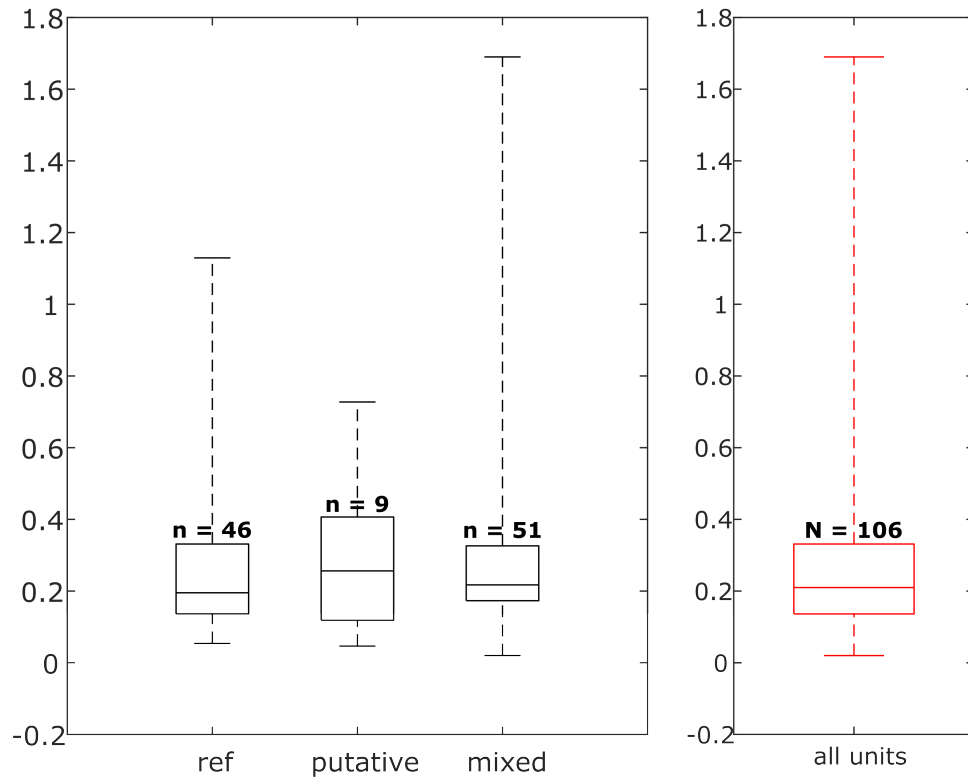


Fig. S3: Distribution of firing rate fold change per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N represent the number of units.

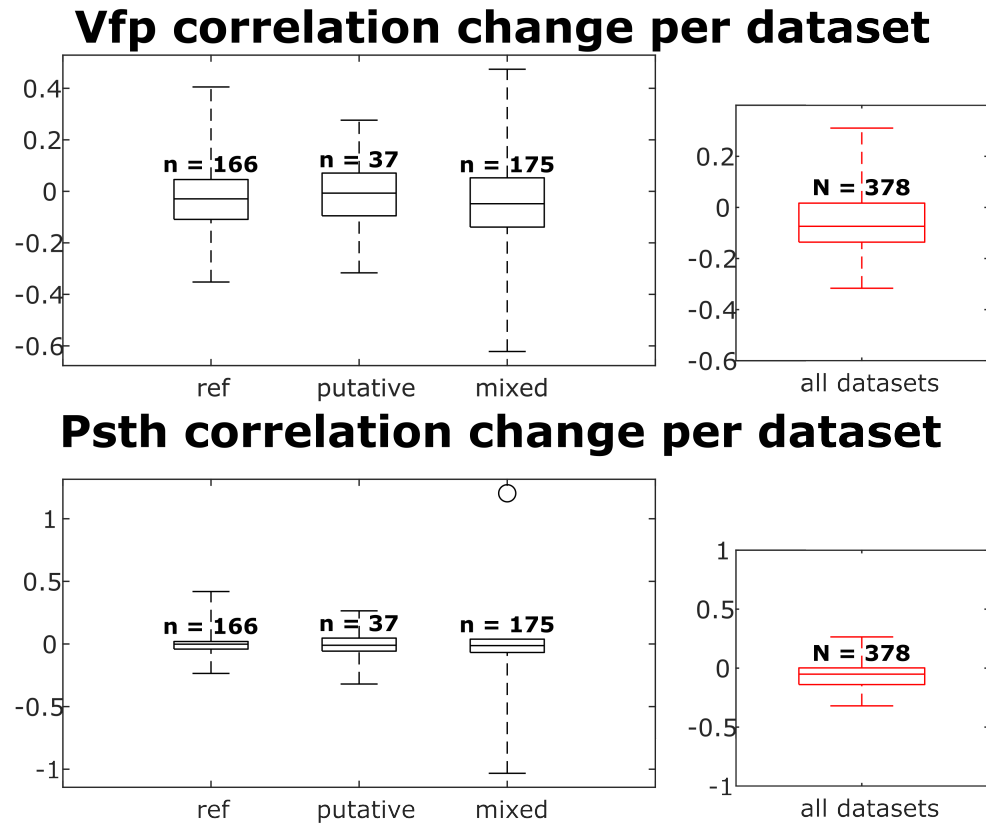


Fig. S4: The visual fingerprint and PSTH change distributions per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of unit comparisons, i.e.(number of units) \times (number of datasets - 1).

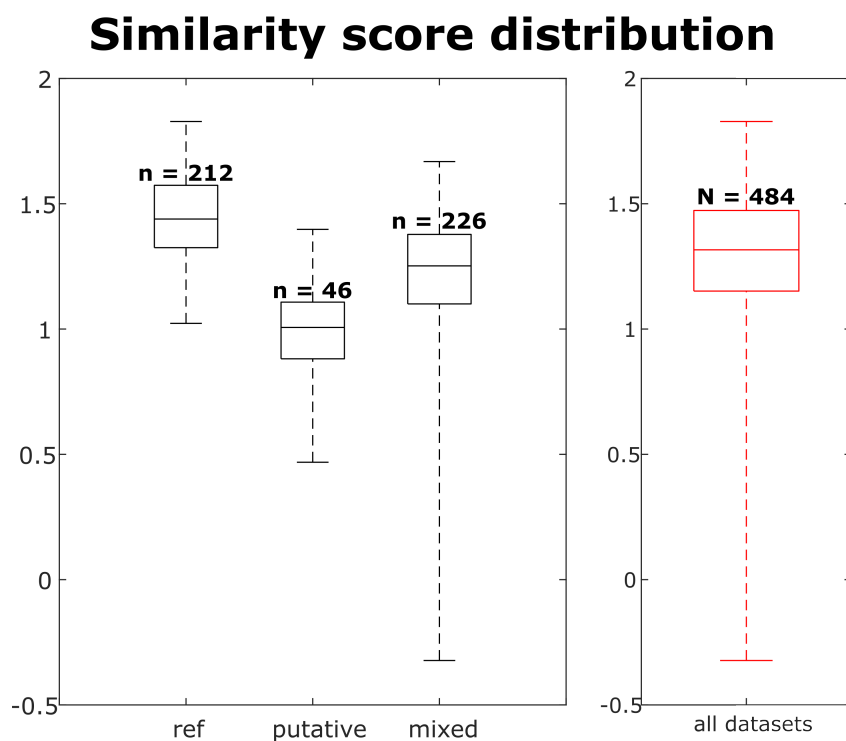


Fig. S5: The similarity score distribution per dataset for each neuron group and across all neurons. Box plots indicate 25% percentile, medians, and 75% percentile. Whiskers at the ends of the box plot show maximum and minimum values. n and N are the number of observations of the units, i.e. \sum_{units} (observations of this unit)

448 8.2 Similarity score heatmap

449 We identify reference pairs as units that are close in space (peak channels separated by $< 30\mu m$)
450 and high similarity score (> 1). Multiple partners can meet these criteria due to oversplitting – these
451 correspond to blocks of high scores in the heatmap. We only include a unit as a reference if its
452 highest similarity score counterpart in the other dataset is within the $30\mu m$ distance threshold.

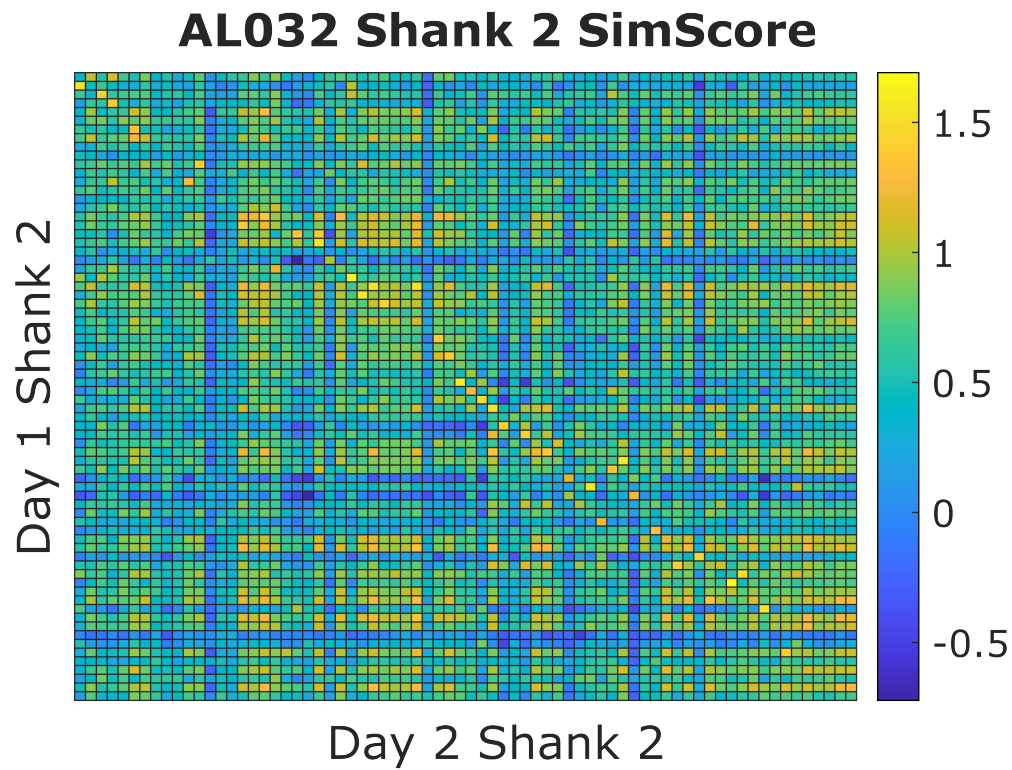


Fig. S6: An example similarity score (vfp + PSTH) heatmap from animal AL032 shank 2 Kilosort-good units between day 1 and 2. Each small square represents the similarity score (value range from [-2,2]) between one unit from day 1 and one unit from day 2. A warm colored square indicates a higher score. The clusters are ordered by their physical locations on the probe. There is a diagonal line with brighter color blocks, indicating that units with more similar visual 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; these units do not respond to the stimulus and are not included in the reference set.

453 **8.3 Pre- and post-drift correction reference unit counts**

454 We showed that between-session drift correction improved yield of reference units.

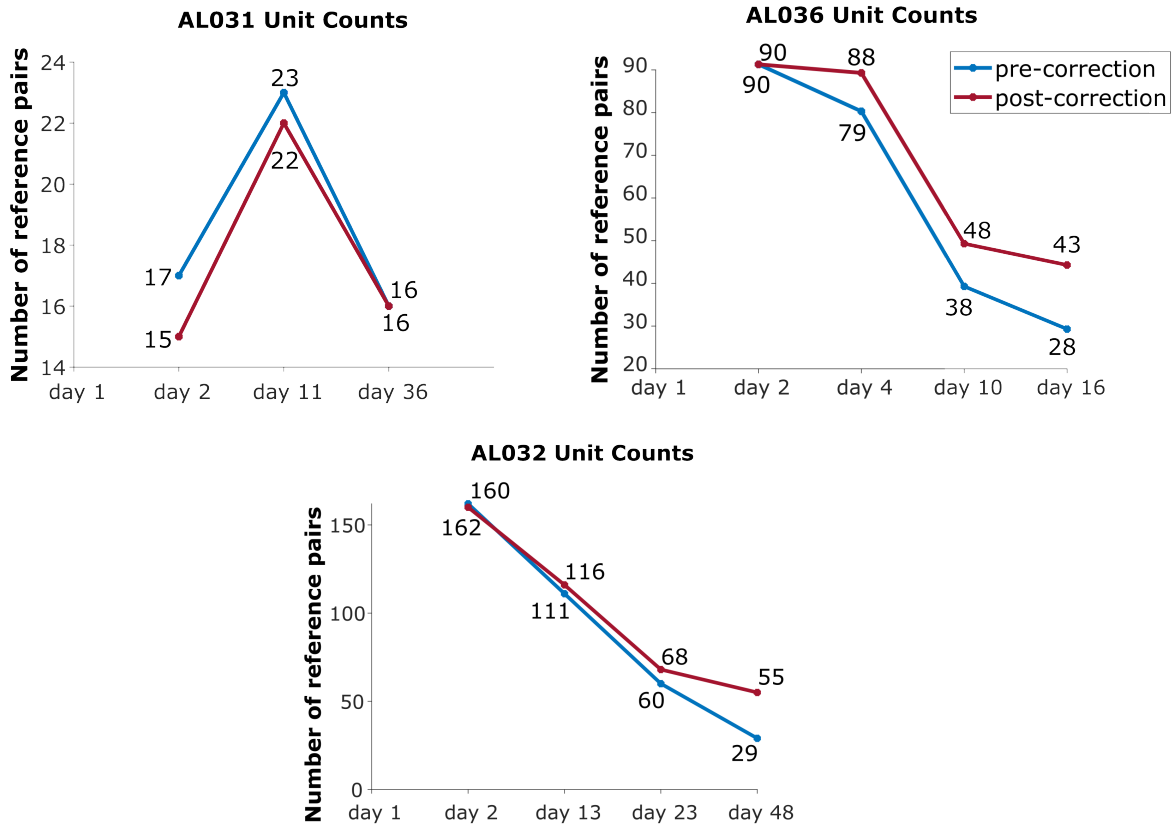


Fig. S7: The effect of drift correction on reference unit yield for all three animals. Note that drift correction improves the recovery rate for most cases; the degree of improvement is a function of the magnitude of the drift.

455 **8.4 Modeling the z-distance distribution for all units**

456 As shown in **Figure 4a**, the z-distance distribution of reference pairs differs significantly from that
 457 of all pairs. To estimate the false positive rate for all pairs, we need to account for this difference.
 458 We cannot simply extrapolate from the measured false positive rate of the reference units. The
 459 difference arises from a bias in the selection of reference units: Because reference units must be
 460 detected in two datasets, they must be easily isolated. We created a simple model to determine
 461 an appropriate functional form to fit the z-distance distribution of all pairs and estimate the false
 462 positive rate.

463 Assume the following distributions:

464 1. The z-distance distribution of all matched neurons, i.e. Ksgood unit distribution, ($\Delta > 0$) is

$$P(\Delta)$$

465 2. The z-distance distribution of matched neurons that are true hits (H : correct match/hits) is

$$P(\Delta | H)$$

466 3. The z-distance distribution of false positive matched neurons is

$$P(\Delta | \sim H)$$

467 Let f be the fraction of units with true hits, then the z-distance distribution for all units is:

$$P(\Delta) = f * P(\Delta | H) + (1 - f) * P(\Delta | \sim H) \quad (6)$$

468 To estimate the distribution of $P(\Delta | H)$, we assume that drift correction works properly. In this
 469 case, the z shift between the two units of a reference pair, or any true hit, is due to the error in
 470 measuring the position of the unit. The distribution of Δz , which is the absolute value of the z shift,
 471 is expected to be a folded Gaussian with $\mu = 0$, and $\sigma = 2 * (\text{error in measured z position})$.

472 To estimate the distribution of $P(\Delta | \sim H)$, we performed a Monte Carlo simulation. In the
 473 simulation, the number of units is 150, the average density of subject AL036. A fraction f will have
 474 real partners in the second dataset. The unit positions in each dataset have normally distributed
 475 errors with $\sigma = 5 \mu m$, matching the observed distribution of z-distance in the reference units.

476 To determine a range of values of f (fraction of true hits) that matches the real data, we can
 477 estimate probability of a hit in terms of probability of being a reference neuron $P(R)$ using Bayes
 478 rule

$$P(H) = P(H | R)P(R) + P(H | \sim R)P(\sim R)$$

479 $P(H | R)$ can be estimated from the reference units recovery rate 0.86, and $P(R)$ can be estimated
 480 from the ratio of reference units, which is 0.29. $P(\sim R) = 1 - P(R) = 0.73$. Then

$$P(H | R)P(R) \leq P(H) \leq P(H | R)P(R) + P(\sim R) \quad (7)$$

$$0.25 \leq P(H) \leq 0.96 \quad (8)$$

481 We modeled the distribution at values of $f = 0.23, 0.5, 0.6, 0.7$ and 0.96 . For each value of f ,
 482 we generate 500 datasets, and compile the z-distance distributions for H and $\sim H$, from the EMD
 483 solution. From these simulations, we learned that the false positive distribution is well fit by an
 484 exponential decay. Therefore, the z-distance distribution for all units is the sum of the two, as
 485 shown in **Equation 3** and Alg. 2.

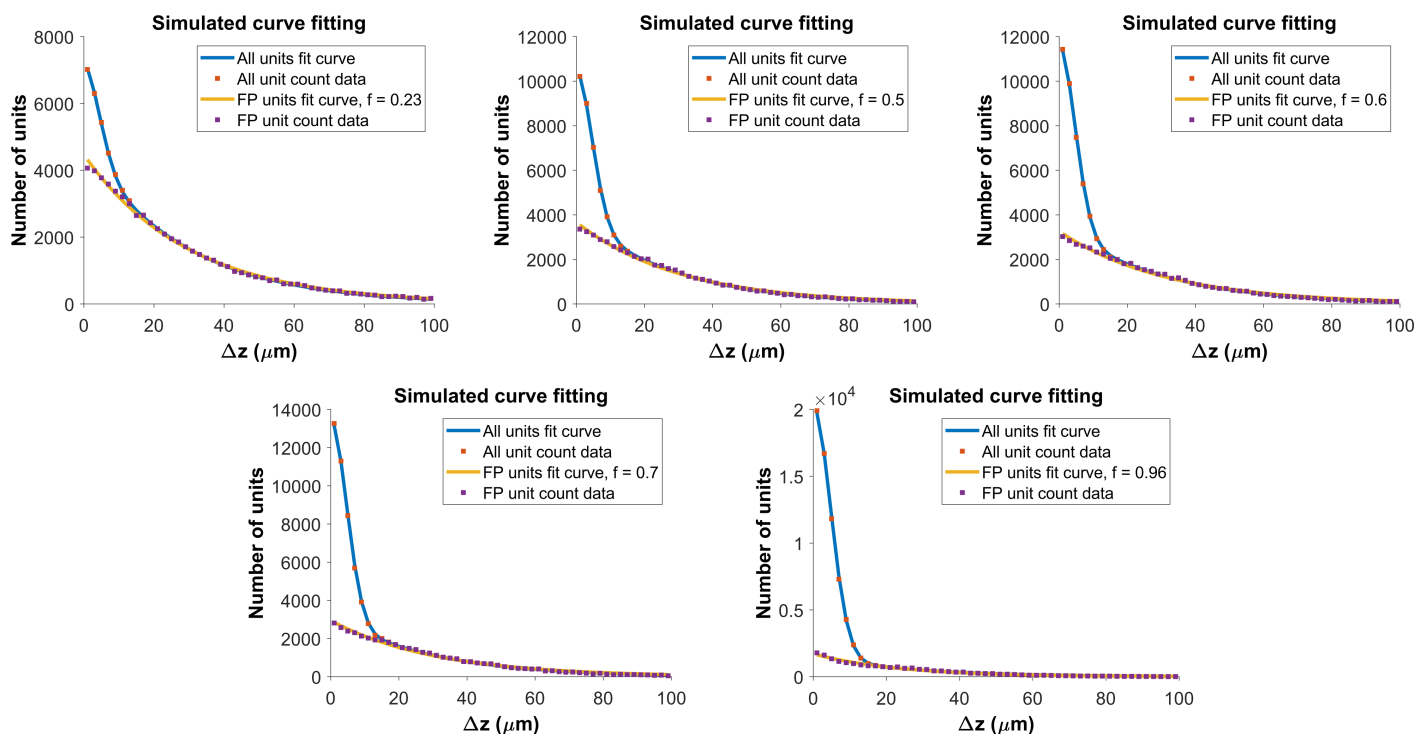


Fig. S8: Fits of z-distance distributions from the Monte Carlo simulations. The five panels correspond to: $f = 0.23, 0.5, 0.6, 0.7$ and 0.96 .

486 To fit experimental data, we first fit the z-distance distribution of the reference units to obtain
 487 the width σ of the folded Gaussian in the first term of *Equation 3*. With σ fixed, we then fit the
 488 z-distance distribution of all KSGood units to *Equation 3* to obtain the width of the exponential and
 489 f . Then we can estimate the false positive rate by integrating $P(\Delta | H)$ and $P(\Delta | \sim H)$ up to the
 490 z-distance threshold. The fraction of false positives as a function of z-distance threshold is shown
 491 in *Figure 4a*, in the bottom panel.

492 Finally, to test model fitting using no information from the reference units, we fit the same z-
 493 distance data allowing the width of the folded Gaussian to vary. *Figure S9*. Panels a and b show
 494 the distribution on the same dataset fit with and without fixing the folded Gaussian distribution
 495 width. The resulting false positive rate from the no-reference fit at threshold $z = 10\mu\text{m}$ is larger
 496 than that from the fit using reference data, so the procedure gives a conservative estimate of
 497 the accuracy.

498 Panel c of *Figure S9* shows the model fit to data from an unrelated dataset acquired from mouse
 499 prefrontal cortex using a Neuropixels 1.0 probe.³⁵ The similar shape of the distribution and a 29%
 500 false positive rate suggest that this method can be generalized.

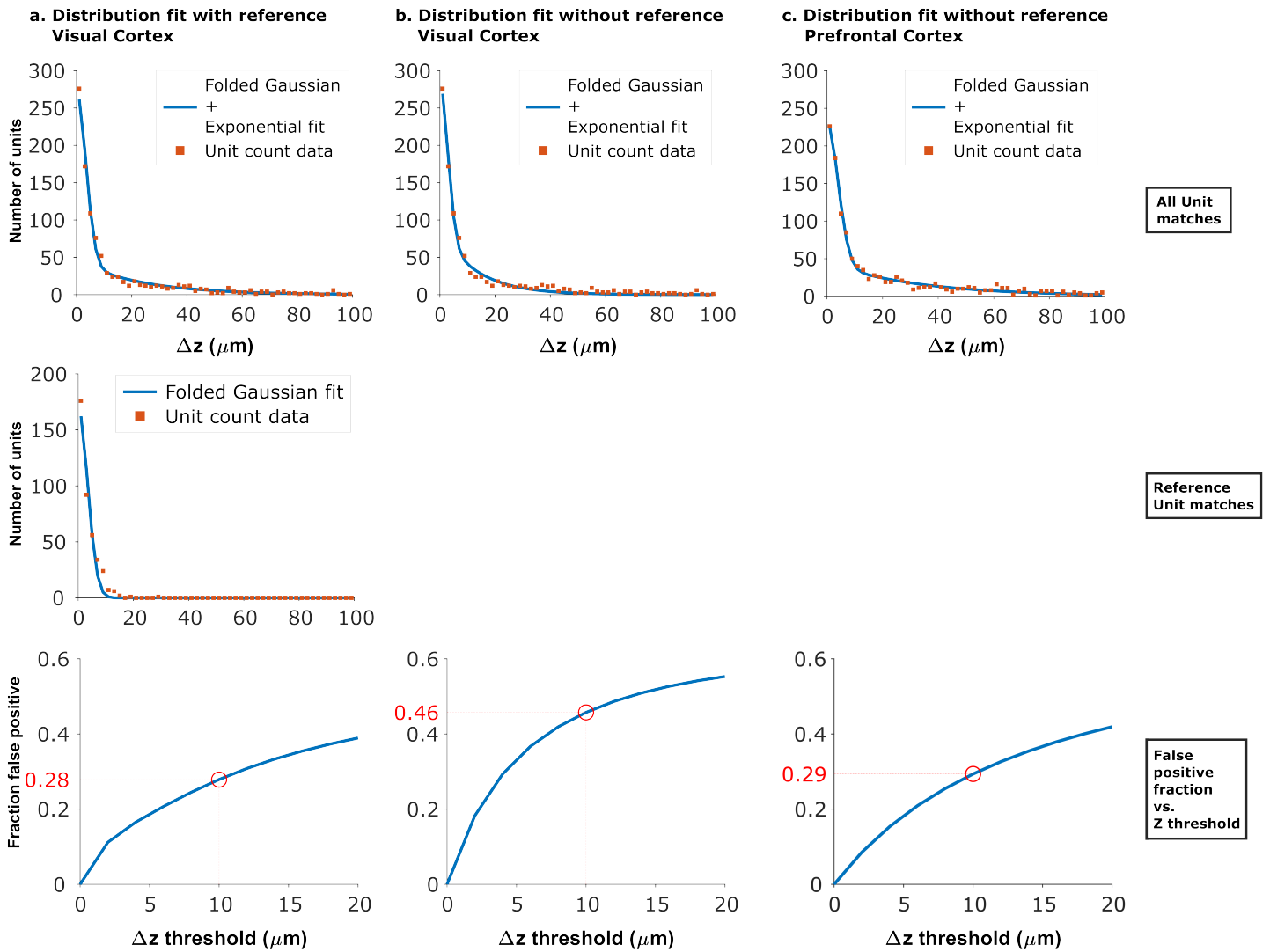


Fig. S9: z-distance distribution fit comparison: a. Distribution fit with 3 parameters, where the z-distribution for true hits is estimated from the reference units. The same as figure 4a. b. Distribution fit with 4 parameters, using no reference information. c. Distribution fit of a dataset in prefrontal cortex using Neuropixels 1.0, using no reference information.³⁵

501 **8.5 Recovery rate vs. time between recordings**

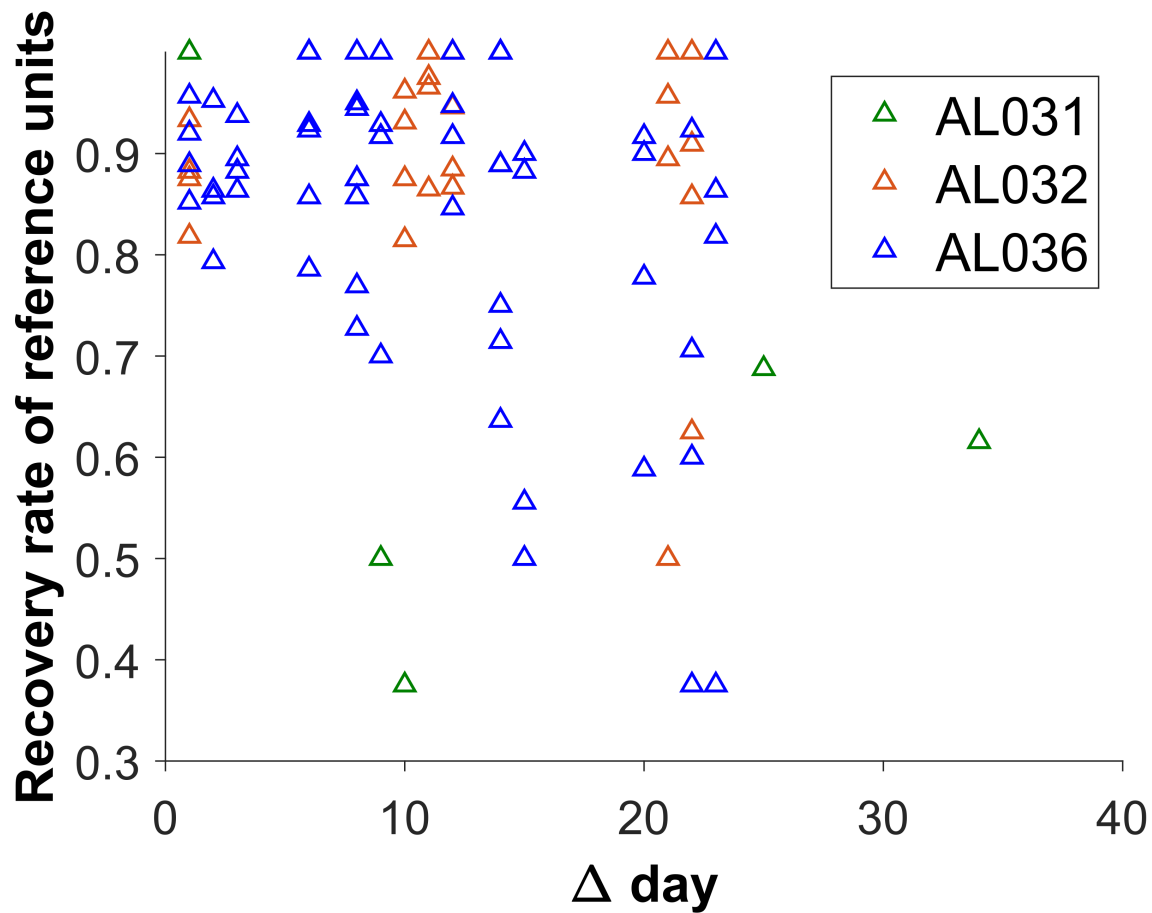


Fig. S10: The reference unit recovery rate for recordings spanning durations. Each triangle represents the matching results of two datasets. Animal AL031 has 6 sets of matched units, with one outlier removed. Animal AL032 has 24 sets of matched units. Animal AL036 has 60 sets of matching. The recovery rate is lower for longer durations.

502 8.6 Reference unit count and the EMD cost matrix

503 In animal AL036, there is a large decrease in the number of reference units after the second dataset, likely due to a large physical shift of the probe relative to the tissue. It is important to be able to
 504 detect such discontinuities to eliminate datasets from consideration. We find that the discontinuity
 505 can be detected in the EMD mean cost, location mean cost and waveform mean cost. The pairwise
 506 values for the costs are shown in *Figure S11*.

507
 508 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
 509 cost, reject H_0 , $p = 6 \times 10^{-7}$; location mean cost, reject H_0 , $p = 6 \times 10^{-5}$; waveform mean cost, reject
 510 H_0 , $p = 5 \times 10^{-7}$). To show that days 3-9 come from the same distribution, we compare odd and
 511 and even rows using the same test. All three cost values show no significant difference between
 512 odd and even days (accept H_0 , $p = 0.92$).

513
 514 Because days 1-2 are significantly different from 3-9, we eliminated them from our analysis.

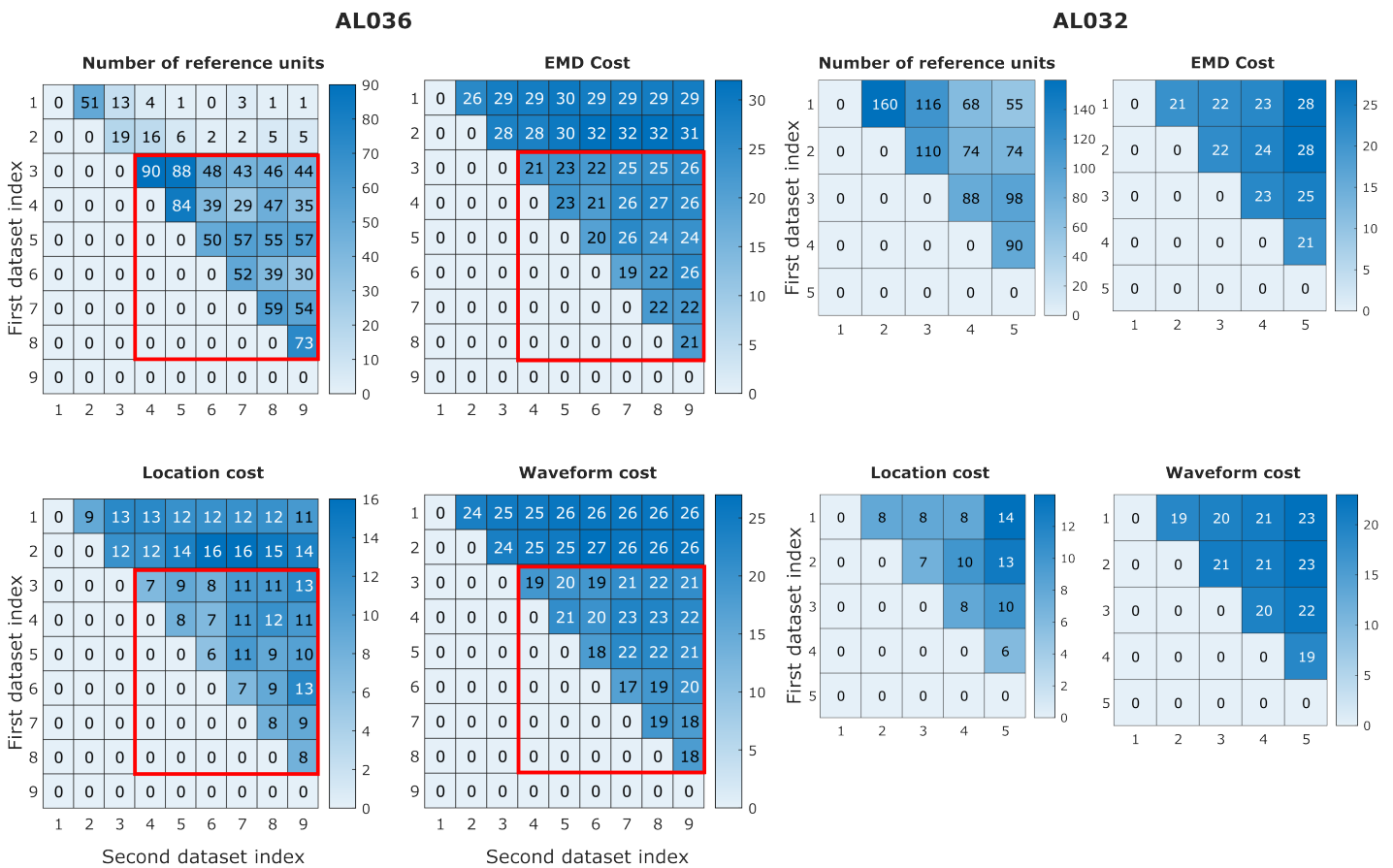


Fig. S11: Reference unit counts and normalized EMD cost for each pair of datasets recorded by the same shank. For animal AL036 (left), we excluded the first two datasets and all of their matching results (first two rows of each matrix on the left) based on the low reference unit counts. Following analysis on their matching EMD cost, location-only cost and waveform-only cost suggest a significant difference compared to the following days (datasets in the red rectangles). We infer that the first two datasets were recorded from a different population than later days. The other matrices show similar information for animal AL032 for reference. To show the relative magnitude of EMD cost in related datasets versus unrelated datasets, we calculated the cost between unrelated datasets with similar unit count (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 those between related datasets (between 20-30).

515 **8.7 Recovery rate vs. the EMD cost**

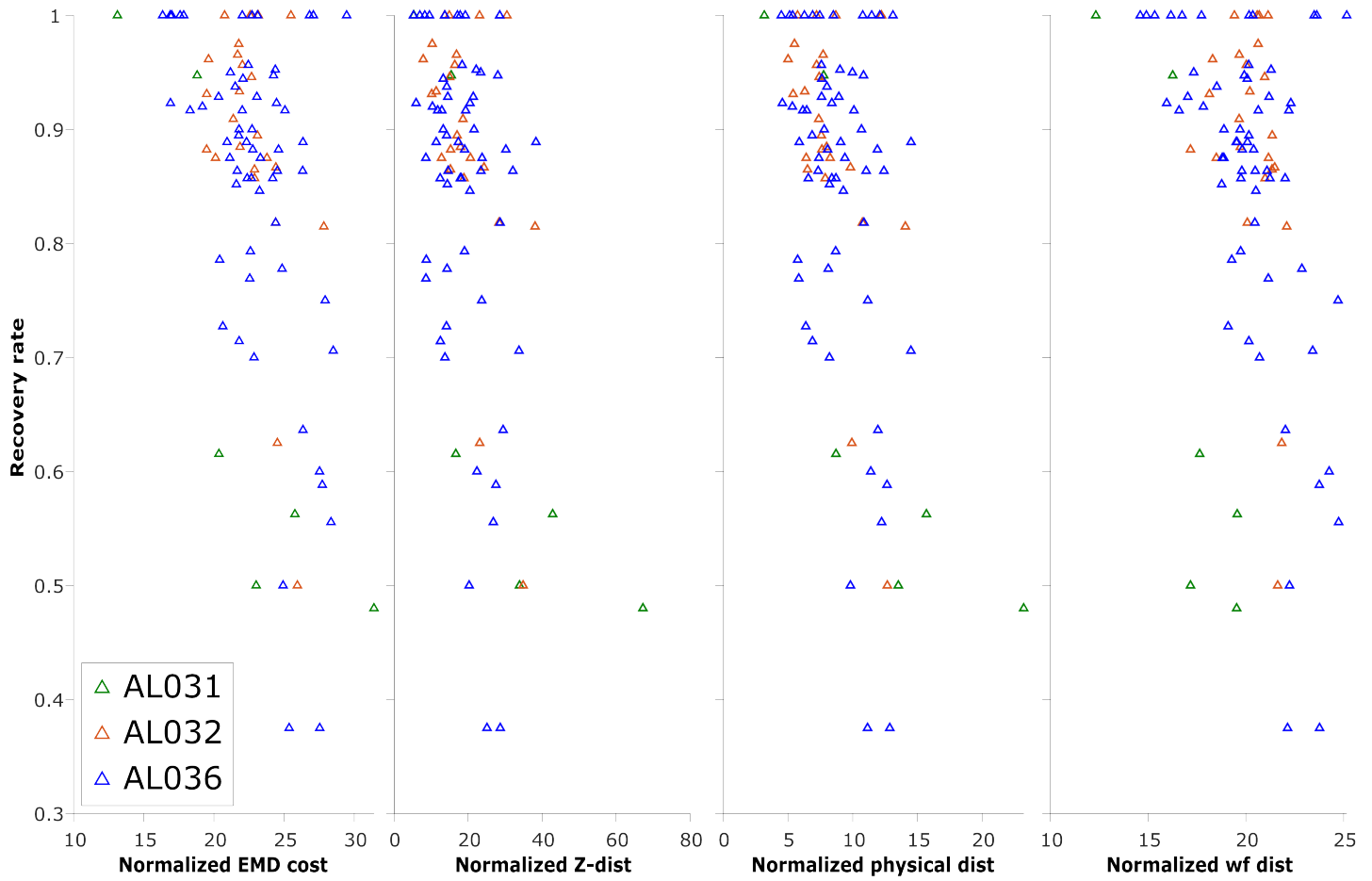


Fig. S12: The normalized EMD cost (unitless), z distance (μm), physical distance (μm), and waveform distance (unitless) and the corresponding recovery rate in pairwise matches of all to all pairs of recordings, on each shank. Each triangle represents the recovery rate in a pair of datasets. Animal AL031 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 cost in the range 20-30. Note that the EMD cost is not predictive of recovery rate.

516 8.8 Reference unit ratio

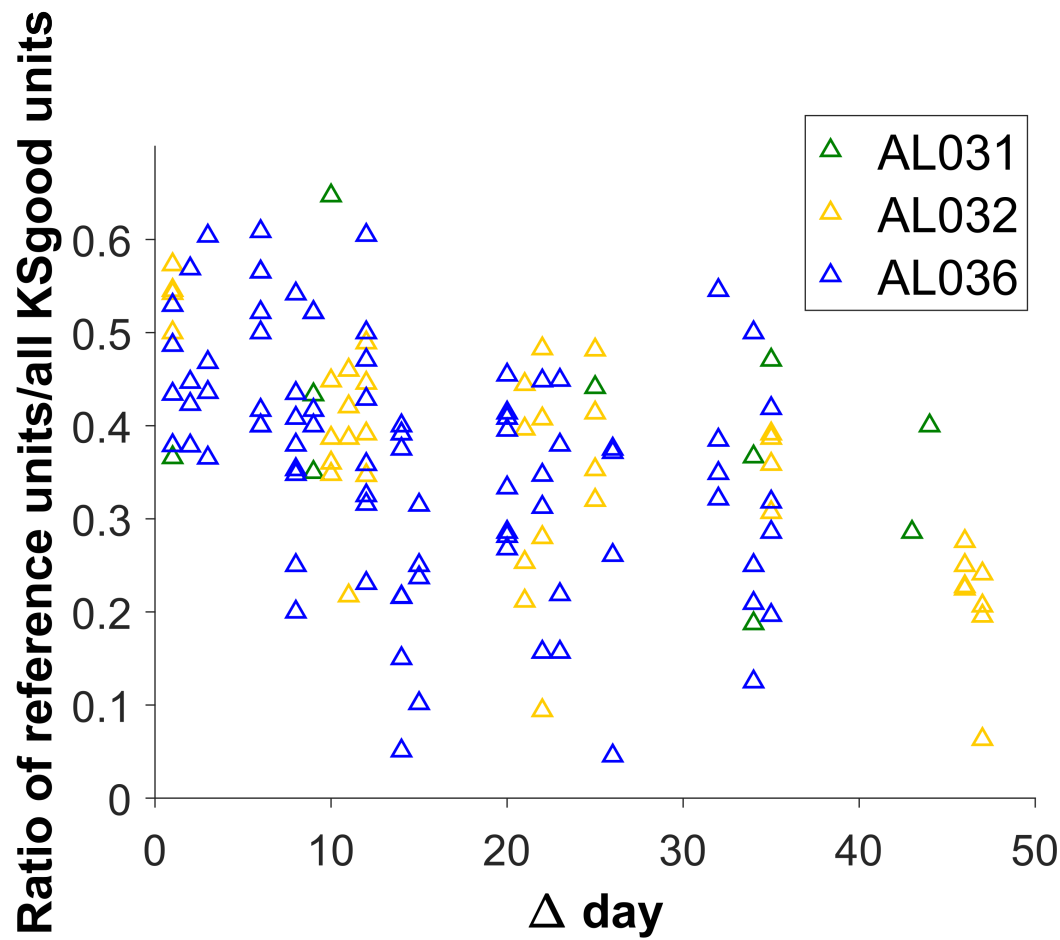


Fig. S13: The ratio of number of reference units to number of KSgood units decreases for pairs of datasets with larger time intervals. However, the variability of the number of reference units is generally large for all time intervals.

517 **8.9 Parameter tuning: L2-weight vs. Recovery rate**

Recovery rate across subjects v.s. waveform metrics weight

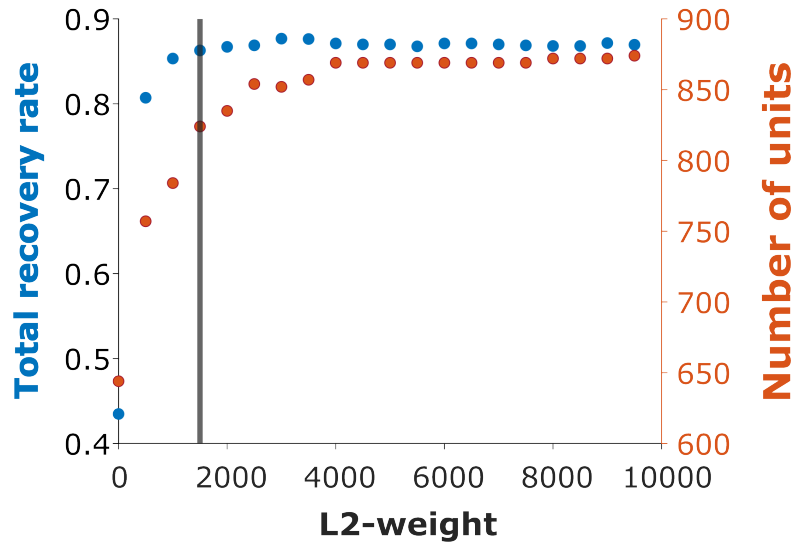


Fig. S14: We varied the weight ω in *Equation 4* used to combine the physical and waveform distances in increments of 500. The vertical line indicates weight = 1500, where the overall recovery rate = 86.29%. The maximum recovery rate = 87.68% occurs at weight = 3000. We chose weight = 1500 for all subsequent analysis.

518 **8.10 Reference unit counts**

519 The number of KSgood units in each dataset and number of reference units between a later
520 dataset and the first dataset in animals AL031 and AL032 are shown here.

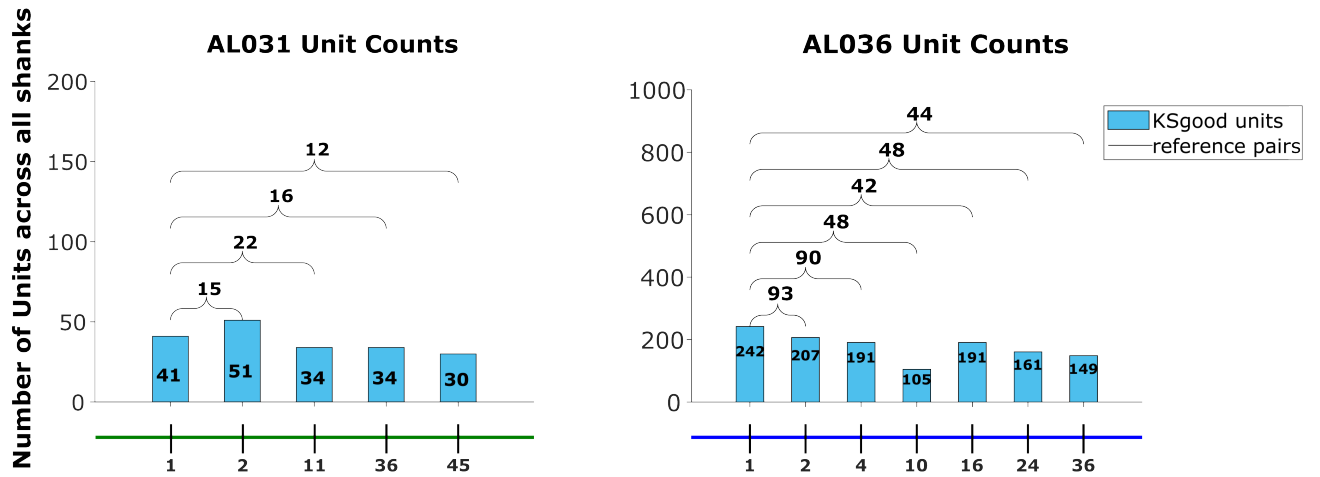


Fig. S15: The Kilosort-good and reference unit counts for the animals AL031 and AL036, as shown for animal AL032 in Figure 5.

521 **8.11 Example reference and putative chains**

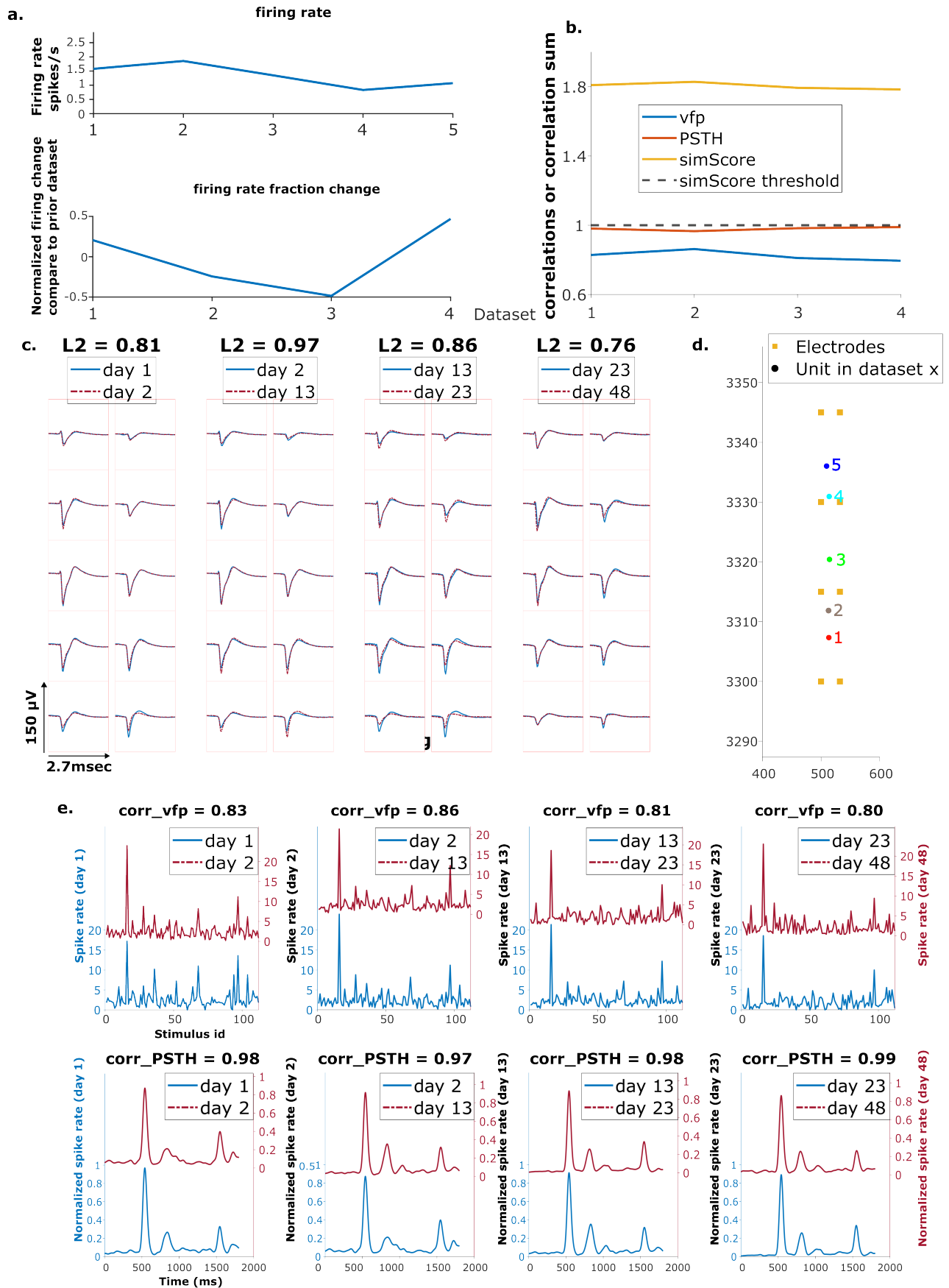
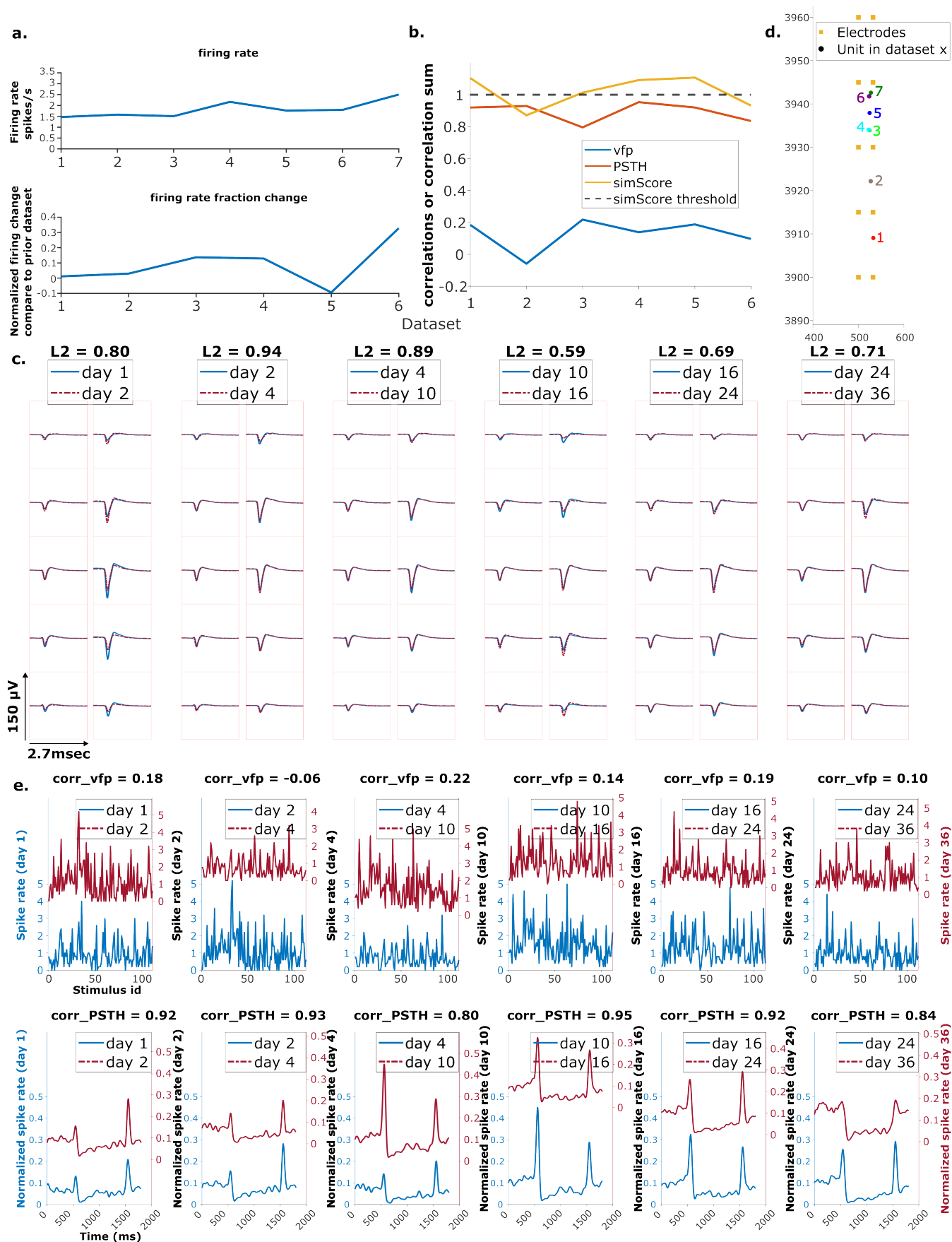


Fig. S16: An example of reference chain. a. Above: Firing rates of this neuron on each day. Below: Firing rate fractional change compared to the previous day. b. Visual response similarity (yellow line), PSTH correlation (orange line), and visual fingerprint correlation (blue line). The similarity score 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.

522



523 References

- 524 [1] Carmena JM, Lebedev MA, Henriquez CS, Nicolelis MAL. Stable Ensemble Performance
525 with Single-Neuron Variability during Reaching Movements in Primates. *J Neurosci*.
526 2005;25:10712–10716. <https://doi.org/10.1523/JNEUROSCI.2772-05.2005>.
- 527 [2] Huber D, Gutnisky DA, Peron S, O'Connor DH, Wiegert JS, Tian L, et al. Multiple dynamic rep-
528 resentations in the motor cortex during sensorimotor learning. *Nature*. 2012;484:473–478.
529 <https://doi.org/10.1038/nature11039>.
- 530 [3] Liberti WA, Markowitz JE, Perkins LN, Liberti DC, Leman DP, Guitchoyants G, et al. Unstable
531 neurons underlie a stable learned behavior. *Nat Neurosci*. 2016;19:1665–1671. <https://doi.org/10.1038/nn.4405>.
- 532
- 533 [4] Clopath C, Bonhoeffer T, Hübener M, , Rose T. Variance and invariance of neuronal long-term
534 representations. *Phil Trans R Soc*. 2017;372. <https://doi.org/10.1098/rstb.2016.0161>.
- 535 [5] Dhawale AK, Poddar R, Wolff SB, Normand VA, Kopelowitz E, Ölveczky BP. Automated long-
536 term recording and analysis of neural activity in behaving animals. *eLife*. 2017;6:e27702. <https://doi.org/10.7554/eLife.27702>.
- 537
- 538 [6] Jensen KT, Harpaz NK, Dhawale AK, Wolff SBE, , Ölveczky BP. Long-term stability of single
539 neuron activity in the motor system. *Nat Neurosci*. 2022;25:1664–1674. <https://doi.org/10.1038/s41593-022-01194-3>.
- 540
- 541 [7] Steinmetz NA, Aydin C, Lebedeva A, Okun M, Pachitariu M, Bauza M, et al. Neuropixels
542 2.0: A miniaturized high-density probe for stable, long-term brain recordings. *Science*.
543 2021;372:eabf4588. <https://doi.org/10.1126/science.abf4588>.
- 544 [8] Luo TZ, Bondy AG, Gupta D, Elliott VA, Kopec CD, Brody CD. An approach for long-term, multi-
545 probe Neuropixels recordings in unrestrained rats. *eLife*. 2020;9. <https://doi.org/10.7554/eLife.59716>.
- 546
- 547 [9] Harris KD, Quiroga RQ, Freeman J, Smith SL. Improving data quality in neuronal population
548 recordings. *Nature Neuroscience*. 2016;19:1165–1174. <https://doi.org/10.1038/nn.4365>.
- 549 [10] Buzsáki G. Large-scale recording of neuronal ensembles. *Nature Neuroscience*. 2004;7:446–
550 451. <https://doi.org/10.1038/nn1233>.
- 551 [11] Brown EN, Kass RE, Mitra PP. Multiple neural spike train data analysis: state-of-the-art and
552 future challenges. *Nature Neuroscience*. 2004;7:456–461. <https://doi.org/10.1038/nn1228>.
- 553 [12] Quiroga RQ, Panzeri S. Extracting information from neuronal populations: information theory
554 and decoding approaches. *Nature Reviews Neuroscience*. 2009;10:173–185. <https://doi.org/10.1038/nrn2578>.
- 555
- 556 [13] Harris KD. Neural signatures of cell assembly organization. *Nature Reviews Neuroscience*.
557 2005;6:399–407. <https://doi.org/10.1038/nrn1669>.
- 558 [14] Quiroga RQ, Nadasdy Z, Ben-Shaul Y. Unsupervised Spike Detection and Sorting with Wavelets
559 and Superparamagnetic Clustering. *Neural Computation*. 2004;16:1661–1687. <https://doi.org/10.1162/089976604774201631>.
- 560
- 561 [15] Chah E, Hok V, Della-Chiesa A, Miller JH, O'Mara SM, , et al. Automated spike sorting algo-
562 rithmbased on Laplacian eigenmaps and k -means clustering. *J Neural Eng*. 2011;8:016006.
563 <https://doi.org/10.1088/1741-2560/8/1/016006>.

- 564 [16] Pachitariu M, Steinmetz N, Kadir S, Carandini M, Harris KD.: Kilosort: realtime spike-sorting
565 for extracellular electrophysiology with hundreds of channels. Preprint at <https://www.biorxiv.org/content/10.1101/061481v1>.
566
- 567 [17] Carlson D, Carin L. Continuing progress of spike sorting in the era of big data. *Current Opinion*
568 *in Neurobiology*. 2019;55:90–96. <https://doi.org/10.1016/j.conb.2019.02.007>.
- 569 [18] Jun JJ, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, et al. Fully integrated silicon
570 probes for high-density recording of neural activity. *Nature*. 2017;551:232–236. <https://doi.org/10.1038/nature24636>.
571
- 572 [19] Hall NJ, Herzfeld DJ, Lisberger SG. Evaluation and resolution of many challenges of neural
573 spike sorting: a new sorter. *Journal of Neurophysiology*. 2021;126:2065–2090. <https://doi.org/10.1152/jn.00047.2021>.
574
- 575 [20] Tolias AS, Ecker AS, Siapas AG, Hoenselaar A, Keliris GA, Logothetis NK. Recording Chron-
576 ically From the Same Neurons in Awake, Behaving Primates. *Journal of Neurophysiology*.
577 2007;98:3780–3790. <https://doi.org/10.1152/jn.00260.2007>.
- 578 [21] Swindale NV, Spacek MA. Spike sorting for polytrodes: a divide and conquer approach. *Fron-*
579 *tiers in Systems Neuroscience*. 2014;8. <https://doi.org/10.3389/fnsys.2014.00006>.
- 580 [22] Bar-Hillel A, Spiro A, Stark E. Spike sorting: Bayesian clustering of non-stationary data. *Journal*
581 *of Neuroscience Methods*. 2006;157:303–316. <https://doi.org/10.1016/j.jneumeth.2006.04.023>.
- 582 [23] Lee J, Mitelut C, Shokri H, Kinsella I, Dethe N, Wu S, et al. YASS: Yet Another Spike Sorter ap-
583 plied to large-scale multi-electrode array recordings in primate retina. 2020;p. 10712–10716.
584 Preprint at <https://www.biorxiv.org/content/10.1101/2020.03.18.997924v1>. <https://doi.org/10.1101/2020.03.18.997924>.
585
- 586 [24] Chung JE, Magland JF, Barnett AH, Tolosa VM, Tooker AC, Lee KY, et al. A Fully Automated
587 Approach to Spike Sorting. *Neuron*. 2017;95:1381–1394.e6. <https://doi.org/10.1016/j.neuron.2017.08.030>.
588
- 589 [25] Chung JE, Joo HR, Fan JL, Liu DF, Barnett AH, Chen S, et al. High-Density, Long-Lasting, and Multi-
590 region Electrophysiological Recordings Using Polymer Electrode Arrays. *Neuron*. 2019;101:21–
591 31.e5. <https://doi.org/10.1016/j.neuron.2018.11.002>.
- 592 [26] Vasil'eva LN, Badakva AM, Miller NV, Zbova LN, Roshchin VY, Bondar IV. Long-Term Record-
593 ing of Single Neurons and Criteria for Assessment. *Neuroscience and Behavioral Physiology*.
594 2016;46:264–269. <https://doi.org/10.1007/s11055-016-0227-8>.
- 595 [27] Rokni U, Richardson AG, Bizzi E, Seung HS. Motor Learning with Unstable Neural Representa-
596 tions. *Neuron*. 2007;54:653–666. <https://doi.org/10.1016/j.neuron.2007.04.030>.
- 597 [28] Lewicki MS. A review of methods for spike sorting: the detection and classification of neural ac-
598 tion potentials Michael S Lewicki. *Network*. 1998;9:R53–78. <https://doi.org/10.1088/0954-898X/9/4/001>.
599
- 600 [29] Colonell J.: ecephys spike sorting. GitHub. [https://github.com/jenniferColonell/ecephys_spike_](https://github.com/jenniferColonell/ecephys_spike_sorting)
601 [sorting](https://github.com/jenniferColonell/ecephys_spike_sorting).
- 602 [30] Cohen S. FINDING COLOR AND SHAPE PATTERNS IN IMAGES (Stanford University, Palo Alto,
603 1999). 1999;.
- 604 [31] Bertrand NP, Charles AS, Lee J, Dunn PB, , Rozell CJ. Efficient Tracking of Sparse Signals via
605 an Earth Mover's Distance Dynamics Regularizer. *IEEE*. 2020;27:1120–1124. <https://doi.org/10.1109/LSP.2020.3001760>.
606

- 607 [32] Boussard J, Varol E, Lee HD, Dethé N, Paninski L. Three-dimensional spike localization and
608 improved motion correction for Neuropixels recordings. *NeurIPS Proceedings*. 2021; <https://doi.org/10.1101/2021.11.05.467503>.
609
- 610 [33] Sauerbrei BA, Guo JZ, Cohen JD, Mischiati M, Guo W, Kabra M, et al. Cortical pattern generation
611 during dexterous movement is input-driven. *Nature*. 2020;577:386–391. <https://doi.org/10.1038/s41586-019-1869-9>.
612
- 613 [34] Stringer C, Pachitariu M, Steinmetz N, Carandini M, Harris KD. High-dimensional geometry
614 of population responses in visual cortex. *Nature*. 2019;571:361–365. <https://doi.org/10.1038/s41586-019-1346-5>.
615
- 616 [35] Böhm C, Lee AK.: Functional specialization and structured representations for space and time
617 in prefrontal cortex. Preprint at <https://www.biorxiv.org/content/10.1101/2023.01.16.524214v1>.