

A hardware/software system for electrophysiology "supersessions" in marmosets

Jens-Oliver Muthmann^{1*}, Aaron J. Levi¹, Hayden C. Carney¹, Alexander C. Huk^{1*}

*For correspondence:

huk@utexas.edu (ACH);
ollimuh@utexas.edu (JOM)

¹The University of Texas at Austin

Abstract We introduce a straightforward, robust method for recording and analyzing spiking activity over timeframes longer than a single session, with primary application to the marmoset (*Callithrix jacchus*). Although in theory the marmoset's smooth brain allows for broad deployment of powerful tools in primate cortex, in practice marmosets do not typically engage in long experimental sessions akin to rhesus monkeys. This potentially limits their value for detailed, quantitative neurophysiological study. Here we describe chronically-implanted arrays with a 3D arrangement of electrodes yielding stable single and multi-unit responses, and an analytic method for creating "supersessions" combining that array data across multiple experiments. We could match units across different recording sessions over several weeks, demonstrating the feasibility of pooling data over sessions. This could be a key tool for extending the viability of marmosets for dissecting neural computations in primate cortex.

Introduction

The marmoset has drawn attention as a complementary nonhuman primate model system for visual neuroscience. While the dominant primate model system in neuroscience, the rhesus monkey (*Macaca mulatta*), has the advantage of (relatively) rich cognitive abilities, a large body and robust physiology, and an aggressive work ethic, their large and convoluted (gyrified) brains currently limit the number of techniques that can be applied for measurements of neural activity. Thus, despite their excellent trainability for complex tasks and willingness to engage in lengthy experimental sessions, the scale and variety of neurophysiological questions that can be addressed have been somewhat limited by practical constraints. Recently, the common marmoset (*Callithrix jacchus*) has emerged as a complementary primate model system because of their smooth (lissencephalic) cortex, opening up a much larger number of cortical areas to the use of large-scale chronically implanted electrode arrays (in addition to other techniques). However, a major current concern for adopting the awake behaving marmoset for detailed quantitative studies is their tendency to perform far fewer trials per session compared to macaques. Such a behavioral limitation would result in correspondingly smaller amounts of neural data (and hence, statistical power) per experiment, undercutting the other advantages of the species, and likely limiting their applicability as a powerful neurophysiological complement to the sorts of quantitative neuroscience work done in macaques.

To redress this fundamental potential limitation, we have developed a straightforward, user-friendly tool for recording from large-scale arrays in marmosets while surmounting the relatively short behavioral sessions performed by this smaller (and more delicate) species. First, we report successful long-term electrophysiological recordings using a new type of multi-electrode array for which primate use has not yet been reported in publication to our knowledge, but which is com-

42 commercially available. These “3D” arrays are available with customizable electrode spacing not just
43 across a 2D grid, but also along the depth of individual shanks. The arrays yielded good quality
44 single-unit (SUA) and multi-unit (MUA) activity, as demonstrated in two different marmoset cortical
45 areas (area MT, and the posterior parietal cortex, PPC). Second, we introduce a transparent
46 means for identifying activity recorded on these arrays, not just within individual sessions, but —
47 importantly — *across* sessions. This integration of hardware and software solutions allowed for
48 data from the same unit to be combined over multiple behavioral sessions, into what we termed
49 “supersessions.” This brings the statistical power of awake-behaving marmoset neurophysiology
50 closer to that of macaques on a per-unit basis, while still allowing for larger scale recordings and/or
51 powerful complementary tools, such as patch-clamp and optogenetics, that are more challenging
52 to perform in macaques.

53 Here, we describe both the physiological and computational components of this tool and dem-
54 onstrate its potential usefulness for transcending the behavioral limitations of marmosets into the
55 realm of detailed, quantitative assessments of neural activity at large scales. Furthermore, the
56 tool we introduce here is intentionally straightforward, meaning it can be readily implemented by
57 others, as well as extended when ongoing updates to hardware and software emerge. We conclude
58 by describing current limitations and how updates to this tool could further improve it.

59 To provide a bit more detail before delving into the results, we found that implanting commer-
60 cially-available 3D “N-form arrays” (ModularBionics, Berkeley, CA, USA) resulted in high quality,
61 stable unit activity in marmosets. In our hands and experiences, this reflected a significant step
62 forward in neural recording success, as two prior attempts using more common types of 2D planar
63 arrays (Utah, Black rock systems) yielded lower-quality outcomes (one successful insertion without
64 detectable spikes and one with spiking activity for about three months after implantation). Al-
65 though our goal was simply to record neural activity and not to mechanistically understand why a
66 particular array style works better or worse, our hypothesis is that there is a reduced initial damage
67 due to the lower number of shanks of the N-form array, allowing to avoid vasculature and permit-
68 ting a slow insertion style. In contrast to single shanks and arrays with a single row of shanks, we
69 believe that long-term stability is improved by a better fixation of the brain tissue, reducing chronic
70 respiratory micromotion (*Prodanov and Delbeke, 2016*), while eventually compromising a smaller
71 brain volume for blood circulation than the larger 2D planar arrays.

72 Given the success of the neural recording hardware in yielding qualitatively impressive neural
73 activity over long time periods, we asked whether such recordings would yield a broad sample
74 of neurons that change from experiment to experiment or if they would yield longer recordings
75 of the same neurons. In the first case, we could ask how neural responses generalize across the
76 population, but would overestimate generalization if we recorded from the a substantial subset
77 of neurons from day to day, but did not recognize that in our analyses. In the second case, we
78 could obtain longer recordings for individual units and hence a higher statistical power. We thus
79 designed a method to systematically compare and match (distributions of) spike waveforms across
80 sessions. Our method identifies units from individual sessions independently, and then integrates
81 spike clusters from new recordings into known, existing ones identified in prior sessions. Analyses
82 of units can therefore be performed over multiple experimental sessions.

83 In order to achieve a representation of spike shapes that was robust to potentially varying noise
84 levels and/or forms across experimental sessions, we extracted simple properties of spike shapes
85 in a narrow window around their peak. This was achieved by matching a family of predefined
86 templates on a GPU to yield a parametric representation of local excursions in the raw voltage
87 traces, which included conventional unit spiking activity, spike events from weaker or more distant
88 neural sources, and noise. Unit isolation was performed as a multivariate classification problem,
89 similar to conventional approaches (*Pachitariu et al., 2016; Rossant et al., 2016; Chung et al., 2017;*
90 *Hilgen et al., 2017; Jun et al., 2017a; Lee et al., 2017; Chaure et al., 2018; Diggelmann et al., 2018;*
91 *Yger et al., 2018*). In our method, we did not threshold spikes during a detection step, but clustered
92 shapes of local minima in the voltage traces. The resulting clusters were then matched across

93 recording sessions. Although we are not deeply attached to this particular spike sorting approach,
94 we provide it as a robust, intuitive starting point, which we validated against a more sophisticated
95 and complex spike-sorting package. Its simplicity also allows for online views of sorting results
96 during experiments, which could be useful for experimental decisions even if more sophisticated
97 sorting routines are employed post hoc.

98 Finally, in addition to laying out the hardware and software that allows for supersession-style
99 electrophysiology in marmosets with chronic recording arrays, we also provide starting-point quan-
100 tifications of the performance of this system. These metrics confirm the applicability of this sys-
101 tem to many conventional neurophysiological experiments given the performance level that arises
102 from the current arrays and implantation style, as well as the spike sorting algorithm. However,
103 the greater value of these metrics is in future use, as they will allow for comparisons of relative per-
104 formance (in matters such as falsely-matched units across sessions) as array technology changes,
105 as surgical procedures are refined, and as different spike sorting algorithms are applied.

106 Taken together, this work puts forth a synthesis of commercially-available hardware and intu-
107 itive software that allows experimenters to overcome one of the major limitations of the marmoset
108 as a model species by introducing the concept of supersessions. More generally, this framework
109 may support better integration of work done in marmosets and macaques, allowing these two
110 awake-behaving primate preparations to have greater scientific overlap and thus to more solidly
111 allow for their relative strengths and weaknesses to be considered.

112 Results

113 Neural activity apparent for more than 9 months on chronically-implanted 3D ar- 114 rays

115 We recorded single and multi-unit (hereafter, "unit") activity in the brains of 2 marmosets, one with
116 a 3D N-form array in and around the middle temporal area (MT), the other with an identical array
117 placed in posterior parietal cortex (PPC). For both arrays (Figure 1 A, B, respectively), we were able
118 to record spiking activity starting a week after insertion. Activity lasted for a duration of at least
119 9 months, as depicted in Figure 1 (top rows). Figure 1 (second rows) show, in comparison, the
120 relatively short durations of individual recording sessions (approximately a half hour to an hour).
121 These durations likely reflect a lower bound on how long marmosets will work, as they were largely
122 determined by the animal's preponent motivation to engage in various visual tasks with no fluid
123 or food restriction.

124 Signal amplitudes (Figure 1, third rows) were fairly constant over long periods of time, per-
125 haps with the first two weeks after implantation yielding smaller signals before stabilizing (i.e., first
126 few recording sessions, visible at the very left of the plots). A gradual decline in signal amplitude
127 was further apparent after about 7 months for marmoset J. Detected events (see Methods) had a
128 wide amplitude range of relatively sparse (0.1 – 10 Hz) events, indicative of spiking activity (Figure
129 1, bottom rows). Taken together, these descriptions of the behavior of the animals and the signals
130 from the electrode arrays lay the groundwork for attempting to stitch together data from multiple,
131 subsequent recording sessions. The next critical step would be identifying unit activity that could
132 conservatively be identified across such sessions.

133 Spike clusters overlap in consecutive sessions

134 Our goal was to identify spikes from the same units across recording sessions. This required mea-
135 sures that would be robust to noise, in the sense that spikes from other neurons would not perturb
136 or distort characterization and identification of a given unit. To that aim, we focused our analysis
137 on a very short temporal window, including only the depolarization phase of a spike, represented
138 by a local minimum in the raw voltage traces.

139 For each local minimum (i.e., putative spike) in the raw voltage trace, we determined: (a) ampli-
140 tude, measured as the dot product with a template (of unit power), expressed in standard devia-

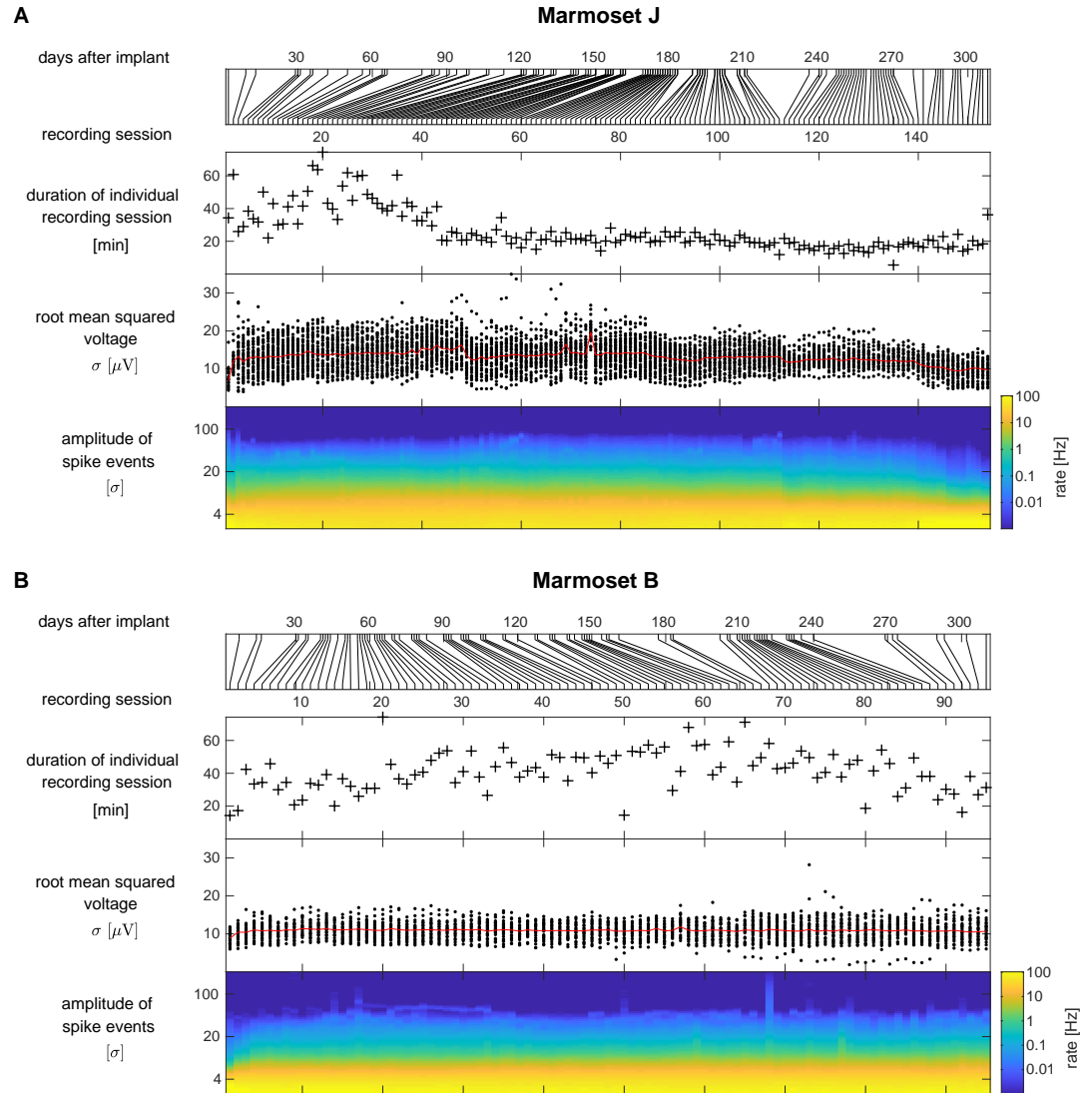


Figure 1. Long-term stability of arrays. **(A)** marmoset J. Top panel: Illustration when individual recording sessions were performed. For clarity, the plots below and in subsequent Figures reflect individual recording sessions rather than time. Second row: Durations of electrophysiological recordings in individual sessions. Third row: Root-mean-squared voltage fluctuations of the common averaged, 300 Hz high-pass filtered data (scatter plots for active electrodes, average shown in red). Bottom row: Amplitude histograms of detected events, averaged across electrodes. **(B)** Same statistics for marmoset B.

Figure 1-source data 1. Source data to generate this Figure

141 tions (σ), as calculated on the high-pass filtered voltage traces; (b) width, measured as the full width
142 at half minimum; and (c) symmetry, measured as the ratio of its falling and rising phase durations
143 (i.e., a 1 : 2 ratio means that recovering back to baseline took twice as long as reaching the voltage
144 minimum).

145 These parameterized shape characterizations of the units were put into 3D-histograms (mar-
146 ginals shown in Figure 2 A) for each recording session, and clustered using a watershed algorithm
147 (see Methods for details). This procedure yielded shape clusters (cyan markers in Figure 2 A) for
148 every session in a common coordinate system to allow for cross-session comparisons of spike
149 shapes. Shape clusters between consecutive sessions often looked very similar, and so we further
150 tested whether they likely reflected spikes from the same or from different units.

151 Specifically, if the brain tissue was held in place by the 16 electrode shanks of the array such
152 that relative movements between the electrodes and the sampled neurons rarely happened, we
153 would always record from the same neurons and see identical spike shapes. Otherwise, if there
154 were substantial shifts in relative position between brain and electrodes, both amplitude and spike
155 shape would shift with movement, and we would be unable to track units across a large number
156 of sessions.

157 We were indeed able to systematically match units across recordings. This was done quantita-
158 tively, using the Jensen-Shannon divergence as a distance measure in the histogram shape space
159 (allowing for small amplitude shifts under a penalty). Figure 2 B shows an example of tracking the 3
160 units observed on February 1 across multiple sessions. Cluster 1 provides an example of a clearly
161 isolated unit with very large spikes with distinctive features, which lasted for about 5 weeks. For
162 this cluster, averaged spike shapes were very similar across recording sessions, with smaller am-
163 plitudes for the initial and final recordings (Figure 2 C, cluster 1). Cluster 2 represents a cluster
164 with more modest amplitude spikes and relatively common spike shapes, resulting in somewhat
165 more variable sorting performance. While being reasonably well-isolated from January 29 to Febru-
166 ary 1, it is contaminated to a variable degree with spikes from different units in other sessions
167 and couldn't be separated from another cluster in two intermediate recording sessions. Cluster 3
168 had low spike amplitudes, but would be considered a decent multi-unit cluster from January 29 to
169 February 1. For the other sessions, there is a small local maximum in the shape histograms, but
170 the cluster would be considerably contaminated with unclassified, smaller amplitude spikes. Given
171 that larger amplitude clusters slowly (and independently) drift over time, we can assume that the
172 same happens to units in this cluster, making it difficult to obtain exact matches across recordings.
173 But, the relatively moderate firing rate of the cluster would suggest that few units with defined
174 shapes were involved, distinguishing it from unclassified spikes.

175 These three example clusters from a brief phase of recording demonstrate both the successes
176 and the challenges of this approach, leaving the real work to be quantifying the overall perfor-
177 mance and aligning particular scientific questions with corresponding tradeoffs between unit iso-
178 lation, data per unit, and number of total units. For example, for the assessment of basic physio-
179 logical mapping and tuning in cortical areas with known columnar architecture, a mixture of singe
180 units and tuned multi-units is often scientifically acceptable, and this approach could provide a
181 wide array of such units, which is important for thorough functional assays. At the other extreme,
182 questions regarding interneuronal correlations can require confidently isolated single units; this
183 approach would provide a smaller number of units, but a large amount of data per unit (as ac-
184 quired across sessions), which could provide critical statistical power for these sorts of detailed
185 questions.

186 In conclusion, our main result is that matching simple shape statistics of spike waveforms across
187 several recording sessions using N-form arrays in marmosets is feasible, and for some units this
188 consecutive recording is possible over notably long periods of time (> 1 month). This grants us
189 the capacity to combine data from multiple experimental days, which we deem "supersessions".
190 Having demonstrated feasibility, we now turn to the issues of validating and quantifying the per-
191 formance of this system.

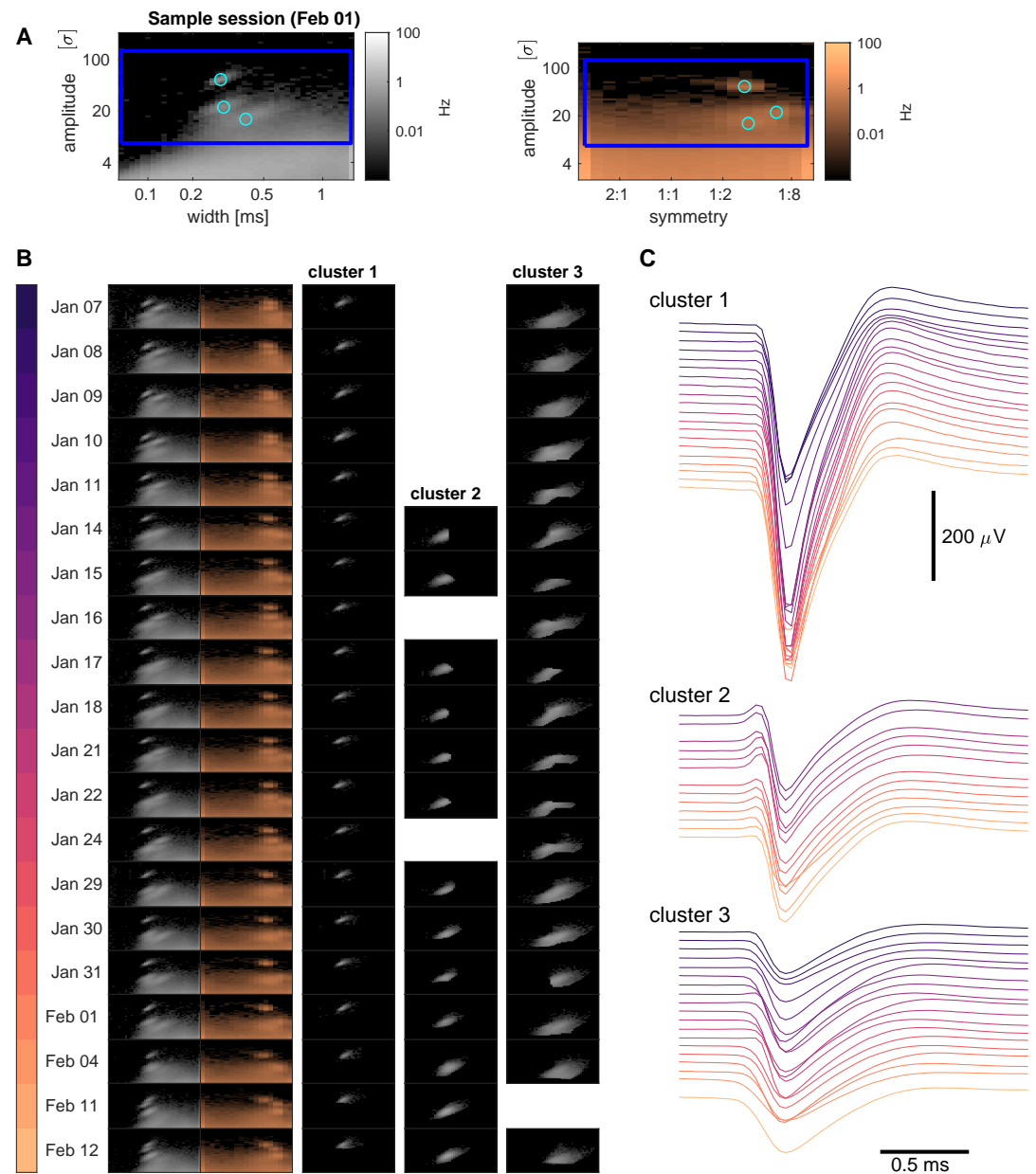


Figure 2. Example of merging clusters across sessions. **(A)** Histograms for amplitudes and widths (left panel) or symmetries (right panel) of detected events on February 1. Regions outlined in blue are shown for a range of dates in (B), using the same color code and axes. Cyan circles mark the three clusters detected in this session. **(B)** Left: marginal histograms of local maxima for 20 consecutive recording sessions, labeled with dates. Right: temporal matches of the 3 clusters found on February 1. **(C)** Waterfall plots of average spike shapes, for dates as color-coded in (B). Data from marmoset B.

Figure 2-source data 1. Source data to generate this Figure

192 **Tuning properties on individual electrodes are stable across sessions**

193 We further confirmed the stability of the measured "supersession" neuronal activity by evaluating
194 the cross-session consistency of physiological tuning properties. This evaluation was done for the
195 MT array implanted in marmoset J, where we were able to confirm that several sites on the array
196 showed directionally-tuned activity in response to moving dots in the left visual field (as expected
197 when recording from area MT in the right hemisphere).

198 The MT electrodes recorded strongly tuned multi-unit activity, so we focused on MUA super-
199 sessions for this analysis. We again used our parameterized representation of spike shapes to
200 determine a region of interest (Figure 3 A, E, outlined in black) in spike shape space with strong
201 directional tuning across recording sessions (Figure 3 A, E). This was feasible because tuning on a
202 given electrode was consistent across a wide range of spike shapes (Figure 3 B, F). For the two MUA
203 sites shown as examples, the direction tuning curves measured were stable over almost 3 weeks.
204 This stability of physiological properties, built on top of the stability of spike shapes themselves,
205 further strengthens the case for the validity and viability of supersessions.

206 We therefore created supersessions across these sessions that exhibited stable tuning and
207 spike shapes, which allowed us to combine larger amounts of data for a single analysis. As an ex-
208 ample here, we show that supersessions allow us to resolve the detailed time course of responses
209 to individual motion directions at a high temporal resolution (Figure 3 C, G). Note that transient
210 aspects of the motion-driven response were very short and consisted of only a few spikes per trial,
211 such that averages across many trials were beneficial. To illustrate this effect, we show the same
212 analysis for responses obtained in a single session (Figure 3 I-K). Averaging over the temporal re-
213 sponses, we then obtained tuning curves for individual sessions (Figure 3 D, H, L).

214 In this example, tuning was stable for considerably longer than one week. This demonstrates
215 not only that shape clusters with high amplitudes were stable across sessions, but also that func-
216 tional properties of low-amplitude activity were conserved across many sessions. Furthermore,
217 being able to combine 10 or more sessions provides an order-of-magnitude increase in trial count
218 that, even assuming some degree of lower-quality unit isolation, should counterweight the rela-
219 tively short individual behavioral sessions. We delve into this issue in more depth at the end of the
220 results sections.

221 **Most units in a given recording were observed for several sessions**

222 Having established stability of both spike waveforms and physiological tuning, we now turn to
223 report a more comprehensive statistical description of recording stability and our ability to distin-
224 guish spike shape clusters (i.e., to isolate one unit from another). A summary of all tracked units
225 across recording sessions is shown in Figure 4. Spike clusters were regions in 3D-shape-histograms,
226 consisting of a set of voxels, which could be divided into boundary voxels (adjacent to a voxel out-
227 side the cluster) and center voxels. If the average spike count in boundary voxels was less than 3/4
228 of the average density in center voxels, clusters were considered as "better-isolated" and shown
229 in darker colors in Figure 4.

230 We further distinguished clusters that lasted for shorter numbers of sessions (<5, orange) and
231 longer numbers of sessions (blue, ≥ 5), as many of the short-lived units had low amplitudes and
232 were less reliably detected.

233 We found that a large proportion of units in a given recording survived for multiple recording
234 sessions (histograms in Figure 4, blue vs. orange), especially when they were considered as better-
235 isolated (Figure 4, darker colors).

236 A more detailed visualization of the survival of individual units is shown in the upper half of
237 both panels in Figure 4. This plot can resolve whether the appearance or disappearance of units
238 between two sessions happened locally (i.e., affecting only some individual units), or globally (i.e.,
239 affecting most, if not all, units across the array). To further see whether the temporal separation
240 (i.e., number of days) between consecutive sessions was a major factor for the loss (/turnover) of
241 units, we visualized the relation between the number of long lasting units lost and the temporal

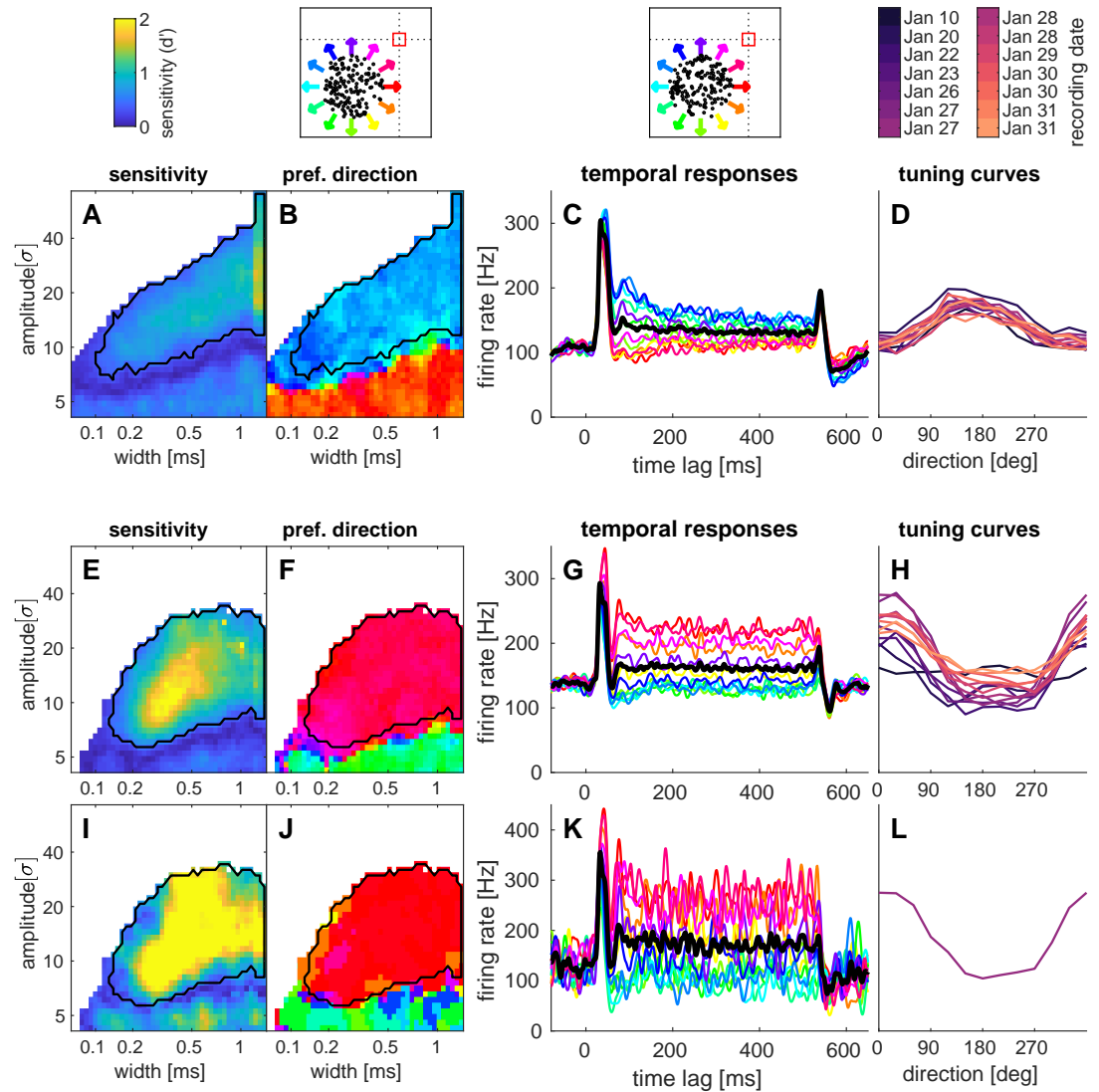


Figure 3. Examples of direction tuning on two electrodes. Top: Legends and stimuli for the examples below. Moving dots were presented at (-15,-15) degrees from the fixation point (red square). (A) Sensitivity indices and (B) maximum response directions as a function of spike shapes. (across sessions, corrected for a cross session baseline effect). The region outlined in black was used for further analysis. (C) Temporal firing rate responses, averaged across sessions and shown for individual tuning directions (colored lines, black line: avg. response, 4041 trials). (D) Tuning curves obtained for individual recording sessions (labeled above, some dates had a morning and afternoon session). (E - H) Same analysis for a second example electrode. (I - L) Tuning observed in a single session (January 27 afternoon session, 254 trials). Recordings in area MT (marmoset J).

Figure 3-source data 1. Source data to generate this Figure

242 separation between the two sessions when the loss occurred (Figure 4, insets). Although larger
243 temporal separations tended to correlate with a higher turnover of units, substantial unit turnover
244 could also occur even with very short temporal separations between sessions.

245 This analysis also highlights a difference between the two animals: while there are several dis-
246 tinct time points of high turnover in marmoset J (Figure 4 A, dotted lines mark disappearances of
247 more than 16 long-term units between consecutive sessions, likely indicative of discrete changes
248 in electrode array position), no such events could be identified in marmoset B (Figure 4 B, dotted
249 lines mark disappearances of the maximum of 5 long-term units, likely indicative of only smaller
250 and/or more gradual changes in array position within the brain). Although we are not sure why
251 the array stability was different in the two animals, this does show that: (a) our analysis scheme
252 is capable of revealing changes and differences in stability; and (b) regardless of whether an array
253 was stable over longer or short terms with or without distinct temporal changes, it is possible to
254 follow units across supersessions in both regimes.

255 We further quantified how often the algorithm would incorrectly classify two units as being the
256 same, by attempting to merge clusters found on different channels. While such chance matches
257 (Figure 4 – Figure supplements 1 and 2) were unable to explain the number and longevity of units
258 we observed, they did vary considerably across clusters, as some spike shapes were more likely to
259 be found in the data.

260 Alternatively to asking how well units matched across sessions, we could ask how much long-
261 term units varied over time. Specifically, we were interested in the variability (or coefficient of
262 variation) of properties which were rather neuron and less network specific. Spike shapes or spike
263 amplitudes (Figure 4 – Figure supplements 3 A and 4 A) were used in the process of merging units
264 across sessions and variability would therefore be biased to lower values. Spiking statistics was not
265 used in this process, and we estimated firing rates (Figure 4 – Figure supplements 3 B and 4 B), as
266 they would not be drastically influenced by experimental conditions. As independent measures, we
267 examined spiking statistics at a fast timescale, arguing that intrinsic neuronal dynamics would be
268 more relevant for the dynamics of bursting behavior than the local network activity. We estimated
269 the maximum instantaneous spike rate in a 50 ms temporal window after a spike, relative to the
270 firing rate of a unit (referred to as 'burstiness', (Figure 4 – Figure supplements 3 C and 4 C), and the
271 time to reach 75% of this rate, which we refer to as 'relative refractory period' (Figure 4 – Figure
272 supplements 3 D and 4 D).

273 All these measures are expected to fluctuate (due to different behavioral conditions, different
274 levels of recording noise, homeostatic changes in neuronal properties and stochastic errors in the
275 estimates), but would on average be even more different between different neurons. We therefore
276 quantified how much of the variability of these four measures was found across sessions in the
277 same unit, as fraction of the variability across sessions and units (Figure 4 – Figure supplements 3 E
278 and 4 E). While we have no ground truth data for how much variability to expect, we report these
279 numbers here and note that further studies would be required with better constrained marmoset
280 behavior or at least longer recordings in individual sessions, especially for interval statistics at a
281 fast temporal scale. We note that in all cases, most of the variance observed across the population
282 was explained by unit identity.

283 Figure 5 shows descriptive histograms of the basic properties of all detected shape clusters
284 (grayscale background). We distinguished clusters that survived short-term (upper row) and long-
285 term (lower row). Several basic relations become apparent from visual inspection. First, the spread
286 (avg. diameter) and firing rates of clusters tended to be larger for smaller amplitude waveforms,
287 likely reflecting the effects of merging overlapping shapes from multiple units. Second, large am-
288 plitude waveforms were generally more skewed than those with low amplitudes, likely reflecting
289 our descriptive approach's ability to identify the basic shape of individual unit waveforms. Third,
290 waveforms from the array in MT tended to be narrower than those from the PPC array (two sided
291 Wilcoxon rank sum test, short-term units: $p=2e-20$, median widths 0.28 ms vs. 0.40 ms and long-
292 term units: $p=4e-19$ median widths 0.24 ms vs. 0.32 ms), perhaps revealing a biophysical difference

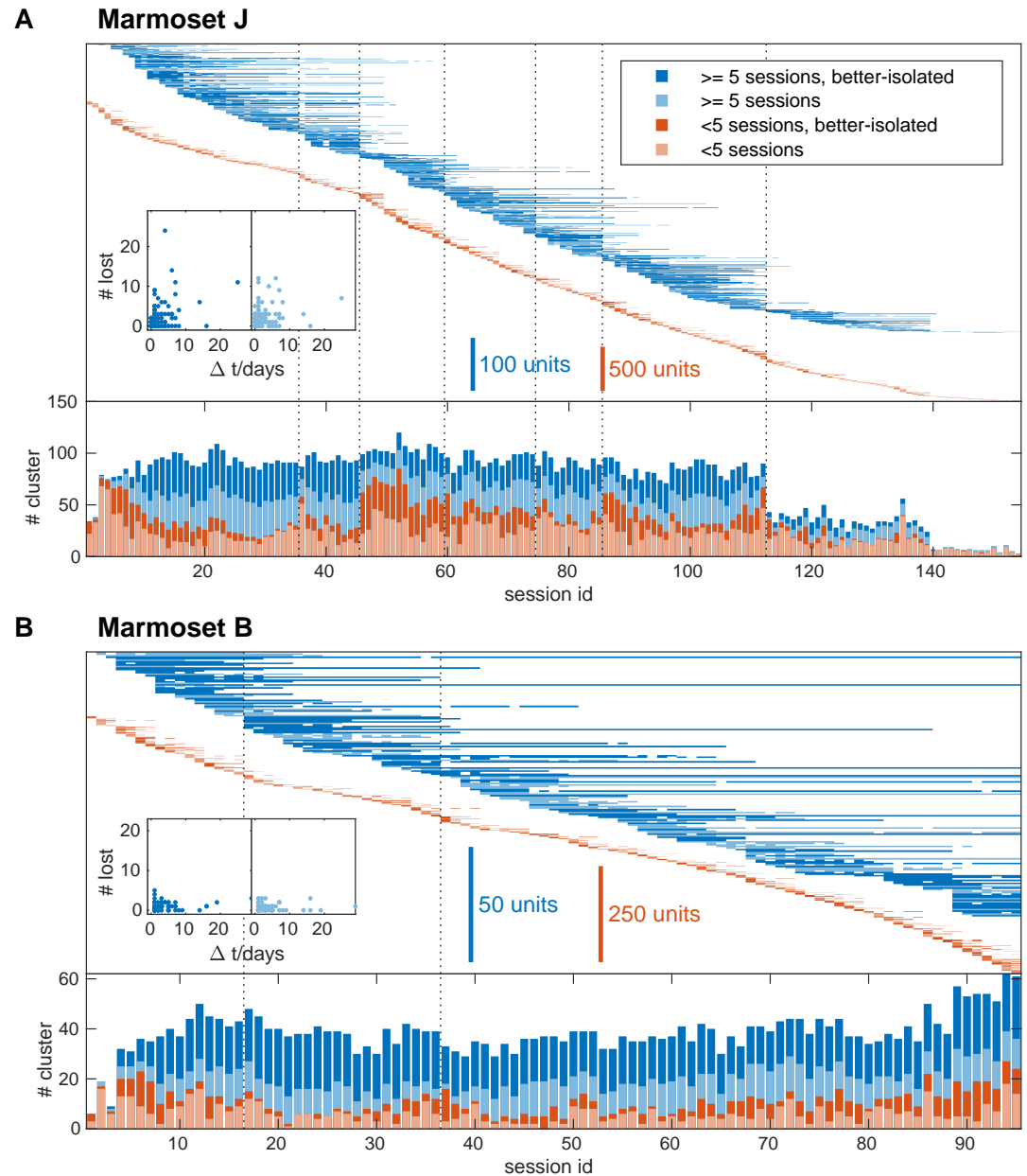


Figure 4. The majority of clusters survives for multiple sessions. **(A)** Clusters detected in recordings of area MT (marmoset J). Top: temporal pattern of long-term (at least 5 sessions, blue) and short lived (< 5 sessions, orange) clusters. Better-isolated clusters are shown in darker shades. Dotted lines mark times when more than 16 long-term units were lost. Inset: Number of disappearing units as a function of the temporal gap between two recording sessions. Bottom: Number of clusters in each session. **(B)** Same plots for recordings in PPC (marmoset B), except that dotted lines mark times when the highest observed number (five) of long-term units were lost.

Figure 4-Figure supplement 1. False discovery rate estimates for marmoset J.

Figure 4-Figure supplement 2. False discovery rate estimates for marmoset B.

Figure 4-Figure supplement 3. Long-term statistics for marmoset J.

Figure 4-Figure supplement 4. Long-term statistics for marmoset B.

Figure 4-source data 1. Source data to generate this Figure and the associated Figure supplements

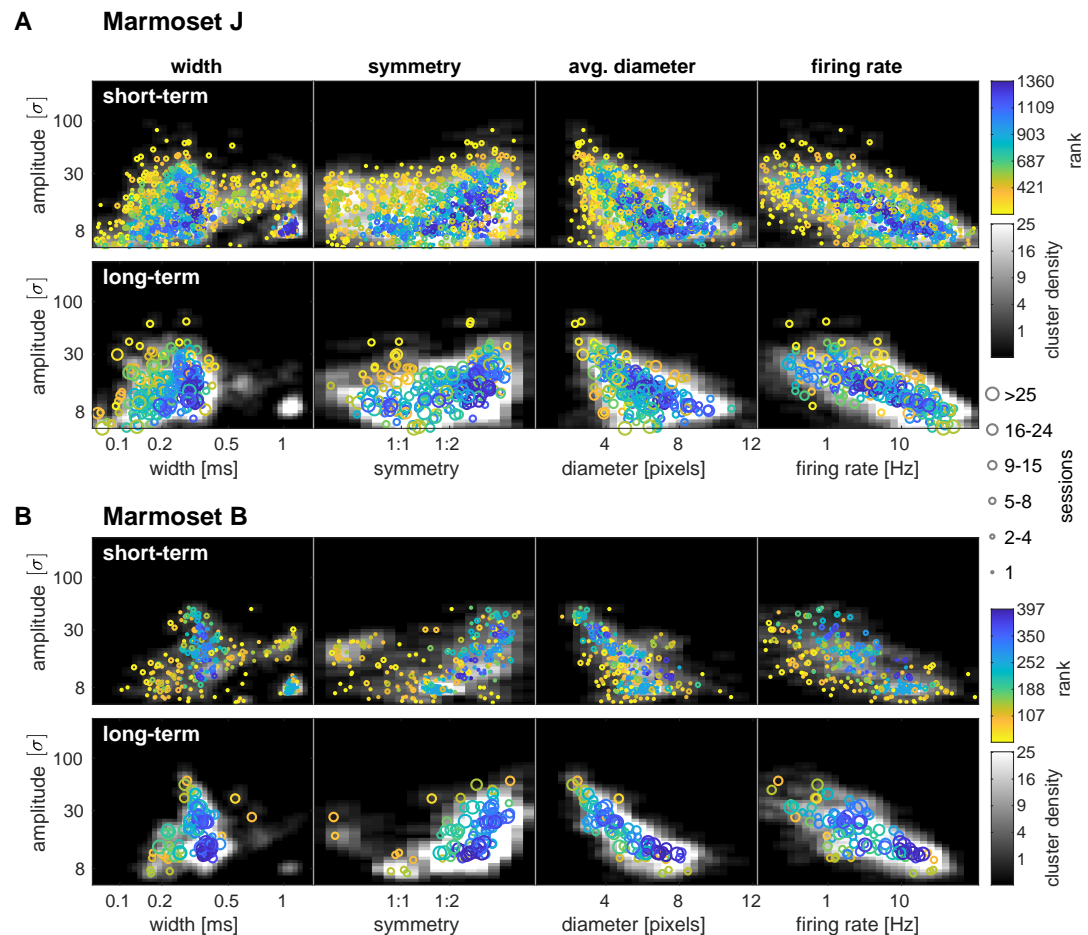


Figure 5. Detected shape clusters are similar (at a population level) when observed for multiple sessions. **(A)** Clusters detected in all recordings and electrodes of area MT (marmoset J). Grayscale represents the density of all detected clusters without merging them across sessions. Colored circles represent individual, better-isolated clusters, merged across sessions. These were ranked according to the corresponding overall density of clusters (i.e. grayscale background) and this ranking is shown in color. Specifically, properties of clusters depicted in yellow were rarely observed and those in blue were commonly found in the data. Clusters surviving less than (top row) and at least (bottom row) 5 sessions are plotted separately for clarity. **(B)** Same analysis for recordings in PPC (marmoset B).

Figure 5-source data 1. Source data to generate this Figure

293 that our approach is capable of picking up.

294 Viewing these basic descriptive plots, we also wondered whether long term matches of spike
 295 clusters might be a result of detecting different units that just happen to produce similar shapes.
 296 To test this, we estimated how likely a given cluster might be mistaken for a different cluster by
 297 counting the clusters with similar spike shapes from all recording sessions. We then ranked better-
 298 isolated clusters according to the number of similar shaped clusters. The resulting rank a cluster
 299 had in the sorted array is depicted in color in Figure 5. A low rank corresponds to isolated units and
 300 a low likelihood to detect the same cluster by chance (Figure 5, yellow/green circles), and a high
 301 rank means that the corresponding spike shapes were frequently observed (Figure 5, blue circles).

302 Sorting clusters in this way allows us to investigate whether clusters with commonly observed
 303 spike shapes would show a bias in long-term survival. We observed that many clusters with unique
 304 shapes survived less than 5 sessions (Figure 5, yellow circles). However, we also noticed that many
 305 of these clusters had uncommonly wide or narrow spike widths or very low firing rates. We there-

306 fore performed a second ranking, which only included units with an average width between 0.1 –
307 0.5 ms and an average firing rate above 0.5 Hz and assigned the excluded units the ranks of the next
308 lowest ranked included unit. This was not done to exclude units from our analysis of the relation
309 between spike waveform uniqueness and lifetime, but to group them more evenly.

310 In order to assess whether clusters with more or less common waveform shapes might show a
311 difference in their lifespans, we analyzed cluster survival, excluding different amounts of the most
312 common cluster shapes. Due to the limited amount of data, we visualized the expected additional
313 lifetime at a given age, assuming a constant probability to lose a cluster in each session. Figure 6
314 shows that this assumption is reasonable, as the expected lifetime does not change dramatically
315 after 5 sessions. Importantly, except for clusters with the 10% most uncommon shapes, the rate at
316 which spike clusters were lost over time did not depend on how common the spike shapes of that
317 cluster were. This is good news, as it does not appear that the longevity of units over sessions is
318 strongly confounded by the appearance and disappearance of units which happen to have similar
319 spike shapes.

320 This analysis also revealed an interesting difference between the two animals: For the array in
321 PPC, cluster survival was about twice as long as for the array in area MT. Although there were more
322 clusters observed for the MT array, we also observed greater variations in signal amplitude and we
323 gradually lost signal in the later recordings of that array (Figure 1 A). We therefore infer that the
324 observed effect could have been due to a higher degree of general instability of the MT array over
325 time.

326 **Supersessions provide the power to estimate spatial and temporal aspects of re-** 327 **sponses across sessions**

328 Finally, we tested whether clearly isolated units could be matched across multiple sessions to as-
329 sess their spatial and temporal properties. We therefore performed generic receptive field map-
330 ping assays at regular intervals over multiple experimental sessions. As proof of concept, here, we
331 describe an example in which both spatial receptive fields and temporal dynamics of responses
332 were estimated using supersession data.

333 Figure 7 shows two example units. The first unit had well isolated, high amplitude spike shapes
334 (Figure 7 C,E) and a pronounced refractory period (Figure 7 F) for at least 6 recording sessions (fir-
335 ing rate (1.7 ± 0.2) Hz; avg. spike count per trial (400 ms) 0.7 ± 0.4 overall and 1.5 ± 0.5 for stimuli in
336 the receptive field). It consistently responded transiently to stimuli in the left visual field, 50-80 ms
337 after stimulus onset. The second example (Figure 7 G-L) shows a unit with an amplitude gradu-
338 ally increasing and decreasing across sessions. Corresponding to an increase in SNR and lower
339 contamination by false detections averaged spike shapes became sharper for sessions with large
340 spikes (Figure 7 K). This unit had a much faster response around 40 ms, consisting of about 1 spike
341 per trial (and eventually a slightly elevated sustained activity during stimulus presentation). In both
342 of these cases, the response properties of the unit would have been difficult to determine using
343 only a single session's worth of data, due to the low absolute number of spikes recorded. For ex-
344 ample, the total number of spikes recorded in the first 400 ms in the receptive field of the unit in
345 a single session was just 20-80 spikes, the total number of spikes across all trials about twice that
346 amount. But by evaluating data across sessions, the supersession data shows that these units had
347 clearly-localized receptive fields.

348 We further investigated how these examples would generalize to a larger population of units
349 with substantial inhomogeneity in both receptive fields and signal-to noise ratio. For this analysis,
350 we found 172 units that were recorded across at least 4 sessions in which we mapped receptive
351 fields. In order to see whether there was consistency in responses across sessions, we estimated
352 receptive field locations for individual sessions and calculated a 'sensitivity index' to quantify the
353 strength of the spatial tuning.

354 Units that were spatially selective generally had receptive fields that were clustered in a small
355 region of the lower left visual field (Figure 7 – Figure supplement 1 A,B). Importantly, we found

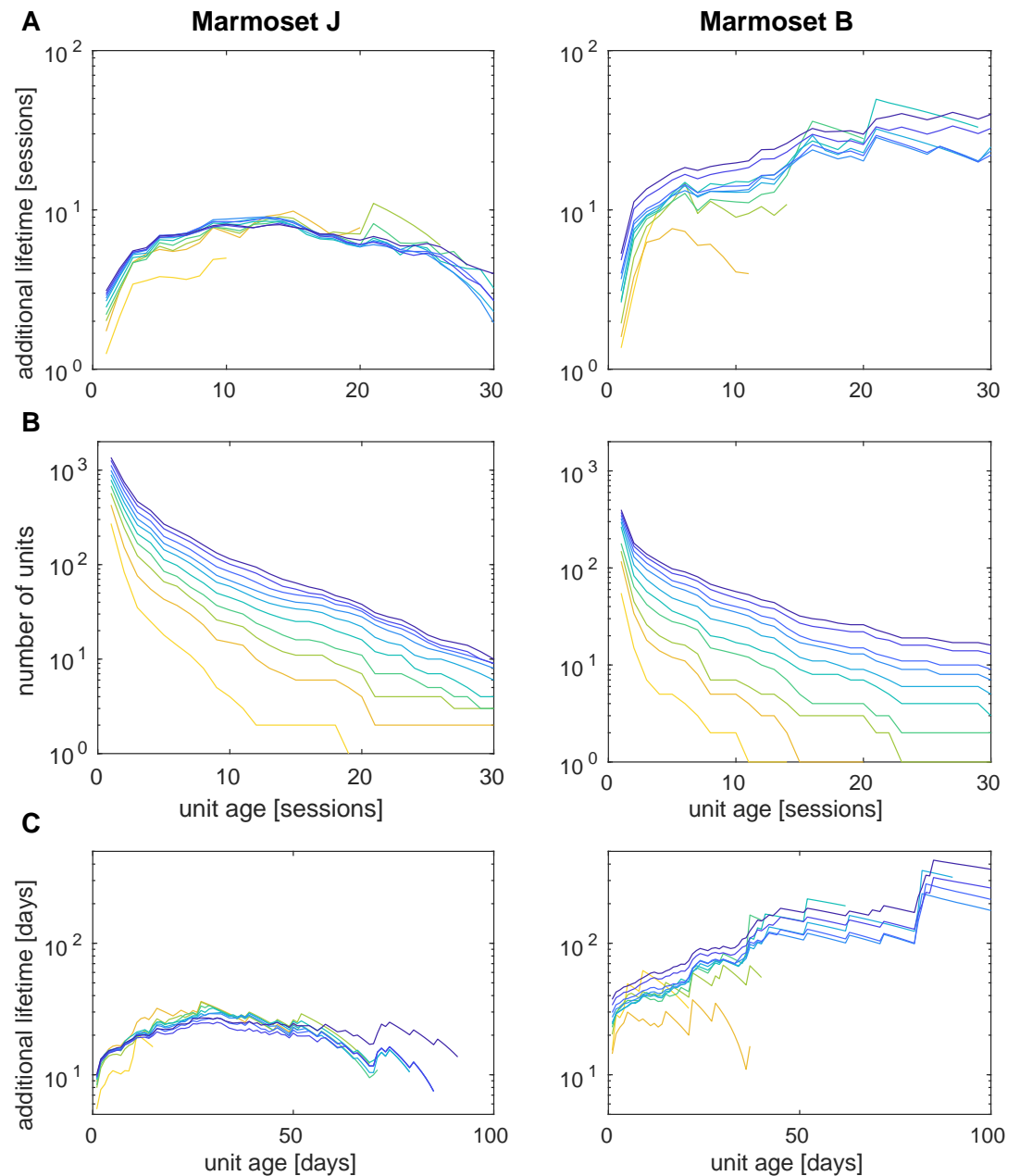


Figure 6. Cluster survival is not an effect of common spike shapes. **(A)** Estimated additional lifetime of clusters after surviving the number of sessions indicated on the x-axis. Coloured lines correspond to the fraction of clusters included in the analysis (steps of 10%, as in Figure 5), where the most yellow curve corresponds to only including the 10% most uncommon shapes. **(B)** Number of units observed for a minimum lifetime. **(C)** Same as in (A) when measured in days rather than sessions. Recordings in area MT (marmoset J, left column) and PPC (marmoset B, right column).

Figure 6-source data 1. Source data to generate this Figure

356 that receptive fields were even better localized across sessions in individual units than across the
357 population of equally well or better tuned units (Figure 7 – Figure supplement 1 C). In addition, we
358 saw that the strength of tuning, (quantified as ‘sensitivity index’, see Methods) generally matched
359 between sessions (Figure 7 – Figure supplement 1 D).

360 While this final analysis outlines a strategy to perform analyses on multi-session and multi-unit
361 data and quantifies consistencies in receptive fields across sessions, we don’t have an obvious ref-
362 erence or gold standard that these numbers could be compared to. These results rather demon-
363 strate what is currently possible, with available data. We do believe that this approach will only
364 improve quantitatively, as array technology continues to improve and yield higher-quality data.

365 Discussion

366 Modern neurophysiological studies in primates require increasingly large amounts of data, either
367 because the parameter space of relevant stimuli or behaviors grows richer (and hence, data are
368 distributed across a larger number of conditions), or because the goal of the experiment itself
369 is to measure more detailed aspects of population activity (and hence, more data are required
370 to estimate higher order statistics). Here, we established the potential of chronically-implanted
371 3D electrode arrays, coupled with a simple unit identification scheme, to allow for the creation of
372 supersession datasets that transcend the standard limitations of marmoset behavior within indi-
373 vidual experimental sessions. We found that high quality activity was evident on this type of array
374 for many months, that a mixture of stable SUA and MUA data could be collected spanning multi-
375 ple individual sessions, and that these supersessions yielded stable physiological characterizations
376 that were more detailed than those from single sessions.

377 Recording performance

378 With the goal of making the marmoset more strongly viable for detailed quantitative studies, we
379 aimed to develop an analysis pipeline that would be robust to different levels of recording quality,
380 measuring single-unit activity where possible, but at the same time considering multi-unit activity.
381 When applying this analysis to data recorded from implanted electrode arrays over the course of
382 more than 9 months and averaging across all recording sessions, we obtained 28 better-isolated
383 units/array/session. For individual arrays, these averages were 32 and 23 for marmoset J and B,
384 respectively, 20 and 18 of which would be seen across a span of five or more sessions. In addi-
385 tion, we found another 40 and 16 multi-unit clusters per array per session for marmosets J and B,
386 respectively; 18 and 9.5 sessions of these multi-unit clusters lasting for five sessions or more).

387 In comparison, previous reports of recording stability using planar (2D) ‘Utah’ arrays in ma-
388 caques (*Dickey et al., 2009; Vaidya et al., 2014; Fraser and Schwartz, 2011*) focused on single unit ac-
389 tivity, which strengthened their claims to be able to track individual units, but at the cost of discard-
390 ing multi-unit activity. Values reported in those prior studies were at most 137 units/array/session,
391 but with large variations across arrays and with decreasing number over time, the average val-
392 ues were closer to 30 units/array/session. In addition, most recordings were done in the first two
393 months after implantation, possibly implying a quicker falloff in signal quality than we encountered
394 with different arrays, and making the comparison to our unit identification and quality less direct.

395 Although a complete comparison between these types of array is beyond the scope of this
396 proof-of-concept tool introduction, we believe it is likely that the variations in performance ob-
397 served with ‘Utah’ arrays in macaques were larger than for the 3D arrays we used. In fact, in mar-
398 mosets, arrays with similar sizes as the ones used in this study (but with fewer electrode contacts)
399 have been reliably implanted and often measured spiking activity for months (*Debnath et al., 2018*).

400 We conclude this comparison by noting that we recorded from a similar number of units as
401 reported for the larger 96 channel ‘Utah’ arrays (*Dickey et al., 2009; Vaidya et al., 2014; Fraser and*
402 *Schwartz, 2011*), but from a smaller region of the brain, largely thanks to the denser 3D geometry
403 of the arrays. This is another advantage on the hardware side of this tool, as it allows for larger

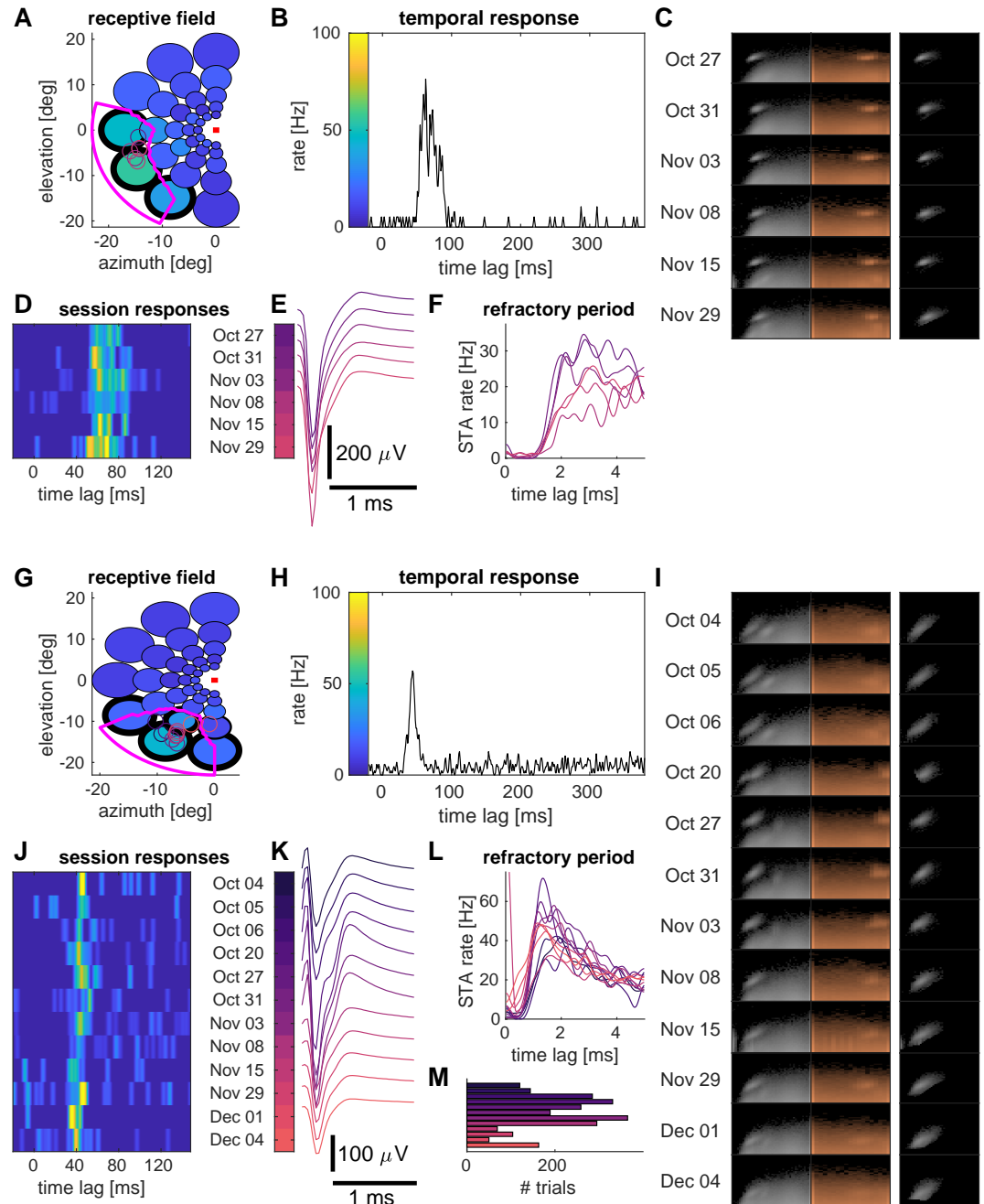


Figure 7. Examples of receptive fields of two units near area MT. **(A)** Maximum firing rates in response to presentation of a disk of moving dots (diameter scaled by 1/2 for clarity; colors indicates firing rate) at a given location in the visual field (fixation spot indicated by a red square). The receptive field (region where the interpolated firing rate exceeded a threshold; see Methods) is outlined in magenta. Colored circles represent estimates of receptive field locations for individual recording sessions. **(B)** Average firing rate for the three conditions (around the RF) outlined in black in (A). **(C)** Marginal shape histograms (as in Figure 2). **(D)** Close-up for firing rates shown in (B) for each recording session. **(E)** Averaged spike shapes. **(F)** Spike triggered averaged firing rates show a refractory period after spikes. **(G-L)** Same as (A-F) for a different unit. **(M)** Total number of trials per session. Colors indicate recording dates (sessions) and firing rates, respectively, and are matched across panels. Recordings near area MT (marmoset J).

Figure 7-Figure supplement 1. Statistics for aggregate data.

Figure 7-source data 1. Source data for this Figure and the associated Figure supplement

404 scale recordings within small brain areas in the marmoset— arrays built for larger primate brains
405 will often sparsely sample within a single area, spanning their footprint over many adjacent areas.

406 **Long-term stability of units**

407 The 3D array recordings had excellent long-term stability, which is a novel and important result for
408 studies using marmosets. The feasibility of long term recordings is itself not totally unprecedented,
409 as there are multiple approaches that align with our observations in a number of species. Here we
410 review some examples, not just to bolster the case that long term stable recordings can be made
411 in a number of species, but to point to the broader potential adoption of the supersession analysis
412 approach we have introduced.

413 For example, *Jackson and Fetz (2007)* used microwires and studied stability of single units in
414 continuous recordings using a window discriminator, and found single units surviving for up to
415 17 days in a one year experiment, where microwires were moved periodically to different neu-
416 rons to improve signal quality. More systematic experiments addressing long-term stability of in-
417 dividual units were done with 'Utah' arrays by matching spike waveforms and inter-spike interval
418 histograms across recording sessions (*Dickey et al., 2009; Vaidya et al., 2014*), eventually in combi-
419 nation with correlations and firing rates (*Fraser and Schwartz, 2011*) to increase statistical power.
420 While comprising relatively small numbers of units and recording sessions, these studies demon-
421 strated a few single units being recorded for months, suggesting that there was likely no relative
422 movement between the electrodes and the neural tissue. *Linderman et al. (2006)* used continu-
423 ous recordings to study short-term changes of spike amplitudes and reported moderate amplitude
424 fluctuations in two example units.

425 The N-form arrays we used had the same spacing between shanks as the 'Utah' type of array
426 — albeit with a higher density of recording sites along a shank, and far fewer total shanks. Even
427 though the N-form arrays comprised only 16 shanks, we found a similar long-term stability for
428 well-isolated single units, suggesting that this number of shanks is sufficient to mitigate substantial
429 array drift. The smaller "bed of nails" also permits a slow insertion method, which we hypothesize
430 is important for avoiding damage associated with ballistic insertion methods, especially important
431 in the smaller and more delicate marmoset brain.

432 In assessing the usefulness of supersession unit data, we used relatively relaxed criteria for unit
433 selection. Given this liberal approach, we did not focus on comparing session-scale average spike
434 waveforms (as these are sensitive to varying amounts of other-spike contamination and noise), but
435 rather distributions of a parametric representation of spikes, where contamination could be con-
436 sidered as a mostly flat, additive component. Likewise, we dropped the comparison of inter-spike
437 interval histograms, firing rates and correlations. While these can provide useful information about
438 unit identity, they rely on a high SNR and good isolation of units in every single session and might
439 even depend on the animal's engagement in experiments. To avoid discarding large amounts of
440 good data without further inspection, we argue that these measures might best be used for post-
441 hoc tests. Spike shapes themselves proved to be reasonably informative about cluster identity, and
442 for short experimental sessions and low firing rates, multiple sessions may be required to obtain
443 useful second order estimates.

444 Recent studies in rodents have been very successful in long-term tracking of neuronal activ-
445 ity. However, this performance was in large part made possible by increasing the density of elec-
446 trode contacts, and therefore the number of observables available for spike sorting. Specifically,
447 *Okun et al. (2016)* successfully sorted concatenated data for a small number of sessions and im-
448 mobile NeuroNexus silicon probes with 4-8 tetrodes (slow insertion). Tetrode recordings in mouse
449 (*Dhawale et al., 2017*) have been used for continuous tracking over weeks. Continuous tracking
450 seems required here due to larger fluctuations in electrical coupling of neurons to electrodes. Re-
451 cent work with high density arrays (*Chung et al., 2019*) in rats showed smaller fluctuations and
452 allowed sorting segments of data and linking these together. Other recent high-density record-
453 ing techniques using ultraflexible mesh electronics (*Fu et al., 2016, 2017*) and silicon high-density

454 arrays (*Jun et al., 2017b*) have not yet been systematically studied for unit longevity. In primates,
455 heptodes have been used in acute recordings, in marmoset cerebellum (*Sedaghat-Nejad et al.,*
456 *2019*) and in macaques *Kaneko et al. (2007)*, and single unit tracking was done in the latter case.

457 In terms of stability of units, the following general picture emerges: wires and tetrodes drift
458 within days, but stability is better when they are left in place without an attached micromanipulator
459 *Okun et al. (2016)* or when they are continuously tracked (*Dhawale et al., 2017*), approaches which
460 can yield stability for days to weeks. Multiple shanks likely reduce electrode drift and units can be
461 tracked for weeks to months ('Utah' arrays potentially for months if no degrading signal quality,
462 *Vaidya et al. (2014); Fraser and Schwartz (2011)*), while ultraflexible, polymer based electrodes
463 might remain stable even longer. Our results fit well into this picture.

464 **Implications for experimental planning and spike sorting methods**

465 Long-term stability offers the potential to generate detailed characterizations of neuronal behav-
466 ior, but it also requires more careful experimental planning. In the two sections below, we high-
467 light conceptual differences for experimental planning and spike sorting compared to the classical
468 single-session approach.

469 **Experimental Planning**

470 While the general long-term stability and the observation of single- and multi-unit activity did sup-
471 port more data-rich analyses than would have been possible from a single session, the fashion in
472 which units ended up being sampled across recordings crucially affects the planning of possible ex-
473 periments. If, at one extreme, we had recorded from a different set of neurons in every recording
474 session, we would have ended up with a large sample of recorded neurons, but not more data per
475 unit. Such a scenario would allow us to estimate distributions of neuronal behavior in a given area.
476 At the other extreme, if we were to always record from the same set of neurons, we would end up
477 with a small sample, but would be able to measure their responses in many different conditions
478 and further quantify the higher-order statistical interactions between them.

479 In reality, we found ourselves in a fruitful middle regime: Units were recorded for variable du-
480 rations, in which a small fraction of units both appeared and was lost between recording sessions.
481 This process was not entirely random, as we saw that most units disappeared during the initial ses-
482 sions after their appearance. This means that the chance for a unit to survive for another session
483 increased with the number of sessions that this neuron had already been observed. Hence, if we
484 were to ask which of the units we would most likely observe in a future session, the best bet would
485 be those units that were already observed for the most sessions in the past.

486 The variable lifetimes of units also provide an additional tool for raising the standard for isola-
487 tion. Restricting an analysis to only long-lasting units would likely reduce the chance of including
488 less clearly isolated units. Such units may not be found in some of the recordings due to variations
489 in signal amplitude.

490 The exact timescales at which units were lost between sessions varied slightly across our two
491 test arrays/animals. However, there may be two different mechanisms involved: while we found a
492 relatively low, constant turnover of units on both arrays, in marmoset J we additionally saw a few
493 events where a large fraction of units was lost between subsequent recordings (Figure 4). These
494 events could not be explained by a long temporal gap between the recordings, suggesting a rela-
495 tively fast mechanism for that, with a timescale of hours to days (as opposed to weeks and months).

496 We believe that these findings can impact the planning of experiments using chronic arrays.
497 In the classical single session approach, experimenters devote part of the experimental time for
498 general characterization of receptive fields and tuning of neurons, in order to target a neuron and
499 adapt the stimulus properties to efficiently sample responses, avoiding stimuli without an expected
500 effect on the neuron's firing behavior. In the case of chronic array recordings, we record from
501 many neurons with potentially different receptive fields and tuning properties, suggesting the use
502 of more general stimuli, e.g. sampling a larger visual area and different tuning directions. Especially

503 when studying interactions between a small number of units, one should keep in mind that some of
504 these units may disappear during the course of an experiment and it would be advisable to start
505 with a larger group of candidate units. In this regard, chronic arrays would be ideally suited for
506 continuous tasks and naturalistic stimuli (e.g. *Huk et al. (2018)*; *Knöll et al. (2018)*), which efficiently
507 sample a large parameter space, allowing for simultaneous characterization of units with different
508 tuning properties.

509 If, however, an experimental design requires finding persistent units in order to adapt focused
510 studies to suit their tuning, we recommend choosing units that have already been observed for
511 at least 3 sessions, as these units have a high chance to survive the next sessions. In our experi-
512 ments, such units had a conditional (additional) lifespan of 6 and 14 sessions (for marmoset J and
513 B, respectively, cf. Figure 6 A). Likewise, studies of changes in firing behaviour of single units across
514 sessions (e.g. while an animal is learning a task, or after drug treatment) are in principle feasible.
515 However, such experiments can usually not be repeated in the same animal, and few units will be
516 clearly isolatable, resulting in a rather inefficient use of the acquired data. In this case, the sug-
517 gested approach is to perform several consecutive studies on an animal, which is possible given
518 the longevity of the arrays used here.

519 Importantly, we have shown that it is feasible to combine data across multiple sessions to infer
520 tuning properties of neurons from multiple sessions. When looking at a population of recorded
521 units, we would encounter a relatively high variability in both signal-to-noise ratio and physiological
522 properties across the population. Such variations would generally result in different requirements
523 on the amount of data needed for statistical tests (e.g. a weak tuning requires more data to deter-
524 mine a receptive field). It was therefore useful to sort units according to their tuning strength, and
525 to perform a relatively focused analysis to specifically detect changes in receptive field locations
526 with high statistical power, using data from single sessions. This strategy would then allow to ask
527 the more detailed questions for data pooled across sessions in a second step.

528 The same type of analysis should be possible for inter-neuronal correlations. Our results also
529 highlight that, in many cases, it would be incorrect to assume that units with similar spike shapes
530 recorded on the same electrode in subsequent sessions would correspond to different neurons.

531 We conclude that chronically implanted electrode arrays allow for both sampling of a large set
532 of neurons and detailed analysis of a few long-term units, but different timescales need to be con-
533 sidered when planning experiments. If the objective is to sample the population of neurons across
534 a brain area, experimental sessions could be separated by a month to take advantage of appear-
535 ance and disappearance of neurons on the array. If instead the objective is a detailed analysis of
536 a smaller set of neurons and their interactions, daily recordings for 2-4 weeks are ideal.

537 Features of the spike sorting method

538 We adopted a modular strategy for spike sorting, where individual sessions were processed inde-
539 pendently and could be iteratively merged to form 'supersessions'. In this way, experimenters can
540 perform preanalyses as data are generated and determine receptive fields and tuning properties
541 of neurons to guide stimulus selection as well as monitor recording quality. This modular approach
542 further facilitates excluding particularly noisy segments in individual sessions, which might impair
543 or bias the clustering algorithm.

544 The primary reason for eschewing existing spike sorting methods was a general concern about
545 robustness when stationarity assumptions were not met across recording sessions. This is a known
546 challenge to even cutting-edge algorithms (*Jun et al., 2017a*). We instead chose a simple paramet-
547 ric representation that was designed to be robust to noise and artifacts, which can differ from
548 session to session. Our focus was on characterizing the peak of the depolarization phase using
549 unimodal templates where the SNR would be highest. While spike shapes can be strongly bimodal,
550 depending on the relative position of the electrode and neuron, the shapes for spikes with highest
551 amplitudes near the soma have been shown to be largely unimodal in theoretical studies (*Lindén*
552 *et al., 2011*; *Quiñan Quiroga, 2009*; *Camuñas-Mesa and Quiroga, 2013*). As we recorded spikes on

553 single electrodes and could expect a large number of neurons in the vicinity of an electrode (*Pe-*
554 *dreira et al., 2012*), high amplitude spikes would be easiest to separate from other units. This
555 situation would certainly be different for high-density probes. The process of estimating parame-
556 ters of the spike shapes was essentially an optimization. We would shift a template temporally at
557 sub-sampling resolution and change its width and symmetry to best match a local minimum in the
558 raw voltage traces. In practice, this step was implemented by running the raw data through a large
559 filter bank on a GPU.

560 Our spike sorting approach did not solve the problem of overlapping spikes. However, it greatly
561 reduced the problem as the time interval needed for detection was reduced to the width of the
562 spike and thus, due to zero padding, much smaller than the the width of the templates in the fil-
563 ter bank. In addition, for cases where overlapping spikes exist, we should see them in the shape
564 histograms as somewhat isolated shapes that are a bit wider and of higher amplitude than an ad-
565 jacent cluster. In our data, we did not find evidence for significant numbers of overlapping spikes
566 near isolated clusters. Overlapping spikes would generally lead to wider and larger observed spike
567 shapes, and such shapes would be reflected as asymmetries in the histograms, where larger and
568 wider than average spikes would be found with a low probability. We didn't observe such asymme-
569 tries, so we can conclude that overlapping spikes were small enough that they wouldn't affect the
570 observed spike shapes to a greater extent than noise. This situation was different for low amplitude
571 events which could not be separated into distinct clusters, but clearly showed stimulus dependent
572 modulations (as in Figure 3 C, G). These events would necessarily overlap in many cases, as their
573 baseline rate was in the order of 100 Hz and peak rates in single trials therefore likely an order of
574 magnitude higher. Hence, firing rate estimates for low amplitude spikes should be read as a lower
575 bound, providing useful (slightly distorted) information about tuning in sustained responses, while
576 truncating transient responses.

577 In this work, we used the parametric representation of local minima as a spike sorting method.
578 But we could certainly perform spike sorting with an existing method and obtain these parametric
579 representations for spikes in order to subsequently match spike clusters across recording sessions.
580 Likewise, as current sorting techniques are validated with respect to stability over long time frames,
581 it would be straightforward to replace our sorting approach. However, our sorting approach could
582 still be used for fast, online assessments of recording quality, neuronal yield and tuning properties
583 as it does not require manual curation.

584 **Application to data**

585 In many cases, we observed that shape clusters appeared and disappeared gradually over time,
586 such that the observed spike amplitudes were highest around the middle of their lifetime. We
587 could thus have a situation where some shape clusters of a given unit were clearly isolated single
588 unit activity, and others were contaminated (e.g. Figure 7 I). Although this effect means that some
589 of the unit data from 'supersessions' is less well-isolated than conventional single-session data, the
590 framework can also be used to estimate the impact of contamination for a given analysis, and
591 hence to determine in a principled manner how high an isolation standard is required.

592 To give an example how such analysis could look, assume that we have a number of sessions (W)
593 where a unit was well-isolated, and some sessions (C), where the same unit was contaminated with
594 low amplitude spikes from other neurons and some of its spikes were lost due to low amplitudes.
595 We would then pool data from each group (W and C) of sessions to obtain a larger sample size and
596 estimate firing rates and interspike interval histograms.

597 Assuming that low amplitude spikes from other neurons are uncorrelated (alternatively, the
598 interspike interval distribution of low amplitude spikes could be estimated with sufficient data)
599 and uniformly distributed, we would fit the ISI histograms of group C as a linear combination of
600 the ISI histogram of group W and a uniform distribution. The component explained by the uniform
601 distribution could then be translated into an estimate of the spike count for the low amplitude
602 spikes from other neurons (i.e., dividing the rate of the uniform component by spike count of

603 group C and multiply with the total recording duration of group C). To obtain an estimate of the
604 number of spikes missed in group C due to low spike amplitudes, one can multiply the difference in
605 firing rates between group W and C with the total recording duration of group C and add the spike
606 count for the low amplitude spikes determined above. After doing a given analysis separately for
607 groups W and C, one could then compare the results and see how they are affected for a known
608 contamination and signal loss.

609 Furthermore, if one looked into the datasets of group W, one would likely find spikes that are
610 statistically similar to the contaminating spikes in group C, simply by identifying identically shaped
611 spikes at much lower amplitudes. Therefore, it is possible to create surrogate datasets with known
612 contamination (and, by removing spikes, signal loss) and treat them as a model to predict effects
613 on a given analysis. The above analysis would then provide independent data to test this model.

614 Apart from spike clusters, our sorting approach also gives access to low amplitude spikes that
615 do show tuned responses to visual stimulation, but likely arise from a multitude of units with a con-
616 tinuum of corresponding spike shapes (e.g. Figure 3). For the purpose of decoding neural activity,
617 such low amplitude spikes can be of great value. In fact, results from other groups indicate that
618 lowering the detection threshold increased the performance of a decoder despite losing informa-
619 tion about the neuronal identity (*Trautmann et al., 2019; Kloosterman et al., 2013; Todorova et al.,*
620 *2014*). Our work suggests that we can define a detection threshold (or region of interest) post-hoc,
621 based on responsiveness to stimuli known to drive neural activity. We refer to this activity as multi-
622 unit hash (MUH), creating a third category alongside with MUA, which should form clusters that are
623 separable from MUA, and SUA which would additionally show a clear refractory period. We need
624 to stress here that MUH is still distinct from the 'unsorted spikes' often left behind by most sorting
625 algorithms.

626 In summary, we were able to create 'supersessions' for individual units on a timescale of sev-
627 eral days to a few weeks. This allows for more statistical power than a single session's worth of
628 data can provide, and hence could put the awake marmoset preparation more on par with that of
629 macaques. This is important because the marmoset is also a "pivot species" to richer and more
630 powerful techniques that are more difficult to apply to the macaque. Such supersessions do re-
631 quire reconsidering the design of experiments to handle the comings-and-goings of identified units.
632 Such experiments will likely have a long term structure where basic characterization of neural re-
633 sponse properties is performed approximately once a week, with the remainder of experimental
634 data collection being dedicated to more sophisticated experiments.

635 **Methods and Materials**

636 **Electrophysiology preparation**

637 Two marmosets were implanted with N-Form arrays (Modular Bionics, Berkeley, CA, USA) in area
638 MT (marmoset J) or PPC (marmoset B). Prior to placing the chronically implanted array, we drilled a
639 grid of 9 burr-holes over and surrounding the desired brain area based on stereotaxic coordinates
640 from *Paxinos et al. (2012)*. We performed extracellular recordings using single tungsten electrodes
641 in each burr-hole to fine tune the placement of the array based on the physiological response.
642 The MT array was placed based on high response to direction of motion, while the LIP array was
643 placed based on high eye-movement related activity. A small craniotomy and duratomy were made
644 surrounding the desired area for array placement.

645 The N-form array was mounted on a stereotax arm and manually lowered till tips of the shanks
646 had entered the brain. The brain dimpled slightly, then the tissue relaxed around the implant.
647 The array was then slowly lowered until the baseplate was just above the brain's surface. The
648 array was stabilized and sealed with KwikCast before being closed entirely with dental cement and
649 acrylic. The array connectors were enclosed in a custom 3D-printed box embedded in the acrylic
650 implant.

651 Animal procedures described in this study were approved by the UT Austin Institutional Care

652 and Use Committee (IACUC, Protocol AUP-2017-00170). All of the animals were handled in strict
653 accordance with this protocol.

654 The N-form arrays (Modular Bionics, Berkeley, CA, USA) consisted of a 4x4 grid of electrode
655 shanks, spaced by 400 μm . Each shank was 1.5 mm long and had 4 electrode contacts, one at its
656 tip, and three more at 250 μm , 375 μm and 500 μm distance from the tip. Extracellular signals were
657 recorded at all 64 electrode contacts with sampling rate of 30 kHz, using the OpenEphys recording
658 system (*Siegle et al., 2017*). For marmoset J, seven of the electrode contacts were found damaged
659 after the surgery and ignored for further analyses.

660 **Visual tasks and stimuli**

661 All stimuli were presented using custom MATLAB (Mathworks) code with the Psychophysics Tool-
662 box (*Brainard, 1997*) and a Datapixx I/O box (Vpixx) for precise temporal registration of stimulus,
663 behavioral, and electrophysiological events (*Eastman and Huk, 2012*).

664 Marmosets were trained to fixate a central dot in the presence of peripheral visual stimuli. The
665 animals fixated the dot within a window of 1.5 degree radius for the whole trial to obtain liquid
666 reward in the form of marshmallow juice. If the marmoset broke fixation, the trial was aborted.
667 Fixation was acquired and held for 200 ms before a stimulus appeared.

668 To measure MT receptive fields, we presented a circular cloud of randomly moving dots for
669 350 ms at one of 35 different screen locations during controlled fixation. The diameter of the stim-
670 ulus aperture scaled with the eccentricity of its center.

671 To measure direction tuning, we presented coherent motion in 12 possible directions at a fixed
672 location based on previously measured receptive fields. Each trial contained motion in one direc-
673 tion for a duration of 500 ms.

674 For PPC recordings, marmosets were trained to perform a memory guided saccade task. The
675 animals fixated the central dot while a target dot was briefly flashed at a random location in the pe-
676 riphery. After a delay of 400-1000 ms, the central dot was extinguished and the marmosets received
677 liquid reward for saccades to the remembered location of the target. Memory guided saccades are
678 well known to generate PPC activity in primates (*Andersen et al., 1990*). The task itself was not part
679 of the investigations in this work. We outline it here as context for the behavioral engagement of
680 the animal in the experiments and to emphasize its potential to drive neuronal activity in PPC.

681 On average, recording durations of individual sessions were (26 ± 13) min for marmoset J and
682 (41 ± 12) min for marmoset B.

683 **Pre-processing**

684 We filtered a 60 Hz component out of the raw data for each electrode using a custom made al-
685 gorithm. We also performed common average referencing by subtracting (projections onto) the
686 median of high-pass filtered signals over all electrodes from each channel. We further up-sampled
687 data to 60 kHz before feeding into Kilosort (*Pachitariu et al., 2016*). For this, values between sam-
688 ples were obtained by linear interpolation and values at samples were smoothed with a [1/6 2/3
689 1/6] smoothing kernel to obtain a uniform variance across data points for the case of Gaussian
690 white noise.

691 **Spike sorting**

692 Code for the spike sorting pipeline is available at <https://github.com/HukLab/SuperSessioning> and
693 will further be made available within the SpikeInterface project (<https://github.com/SpikeInterface>,
694 *Buccino et al. (2020)*).

695 We aimed at jointly sorting spike data from tens of recording sessions (marmoset J: N=154,
696 marmoset B: N=95) under the following constraints:

- 697 1. Marmosets were head-fixed, but able to move their bodies within the chair, creating tempo-
698 rally variable amounts of noise in the data.

- 699 2. Electrodes were separated by at least $\geq 125 \mu\text{m}$ and spikes were not generally expected to be
700 seen on multiple electrodes.
- 701 3. We observed only few separable units (0-3) per electrode.
- 702 4. There was no apparent electrode drift within recording sessions.
- 703 5. Spike clusters needed to be matched across recordings.

704 If spike shapes are known, then template matching would be the best way to detect spikes. How-
705 ever, if spikes are to be sorted, information in the raw data needs to be used to separate spike
706 clusters, and especially to separate them from fluctuations in the background noise level and low-
707 amplitude events of neuronal origin. A good sorting algorithm therefore needs to make estimates
708 that are maximally invariant when subjected to noise. Potential issues are:

- 709 1. Baseline estimate: errors could change the match of bimodal templates. This may especially
710 become a problem when the noise level is temporally varied.
- 711 2. Sampling frequency and temporal resolution for peak detection: Misaligned spikes differ in
712 shape. This can be resolved by upsampling the data, but results in longer templates.
- 713 3. Temporally overlapping spikes: Need to be detected and fitted.

714 To address these three issues, we generated a bank of unimodal templates (essentially triangles
715 with a tip rounded off by a cosine function) which varied in phase (to effectively yield 180 kHz sam-
716 pling frequency), width and symmetry (see examples in Figure 8 B), covering a wide range of pos-
717 sible shapes. Each template was normalized to have an energy (sum of squared entries) of one.
718 Using this bank of templates in a template matching strategy reduces baseline errors, temporal
719 misalignment and the chance of fitting overlapping spikes, but does sacrifice some detection power
720 (when compared to using templates generated from the data, about 10% of the signal power).

721 We determined local maxima (in time and width, but global in symmetry to avoid double de-
722 tectations) for the match (dot product) between our templates and the preprocessed voltage traces.
723 In this setting, we were fitting the peak of the depolarization phase of a spike. While an error in
724 the baseline estimate would have an effect on the detected spike power, it would have little effect
725 on both the estimated spike width and symmetry. Temporally overlapping spikes were less likely
726 as the temporal interval for detection was restricted the duration of the depolarization phase (i.e.
727 0.5 ms or less) and a linear combination fitting was not necessary in our recordings. Note that we did
728 not capture the repolarization phase of a spike at all, however, we argue that due to smoothness
729 constraints, the shape of the repolarization phase covaried with its symmetry, and its duration was
730 hard to estimate due to potential drifts in baseline. Matching a large set of potential templates was
731 computationally expensive, but also well suited to run on a GPU. Our implementation ran about
732 twice as long as recording the data for 64 electrodes sampled at 30 kHz. Marginal histograms of
733 shapes obtained for an example recording are shown in Figure 8 C.

734 Clusters of spike shapes were then determined with a density based approach, using the wa-
735 tershed algorithm, which required some amount of smoothing and a step to reduce global density
736 gradients. In more detail, we aimed constraining the number of spikes to average, rather than
737 setting a fixed kernel size for smoothing. For a given number of spikes, we could then estimate
738 the radius required to find that number of spikes, and the watershed algorithm would yield clus-
739 ters. This approach tends to fail when there are global gradients in spike density and further use
740 extremely small volumes for high spike densities. Therefore, we instead determined the area of a
741 number of spikes that scaled sublinearly with a local firing rate baseline R . This baseline was es-
742 timated by smoothing with a trivariate Hanning kernel (width 13 bins, truncated first and last bin,
743 and sheared a bit, by 0.5 bins in amplitude per bin in width, such that larger spike widths would be
744 combined with lower amplitudes, to reduce a potential bias due to spike clusters, which were often
745 tilted in the opposite direction). We applied a sublinear scaling and added a small offset to that
746 baseline to determine a firing rate (and therefore the number of spikes), given by $0.015 \text{ Hz} + 0.7R^{0.9}$
747 for which we determined the required radius. We excluded areas from the analysis for which that

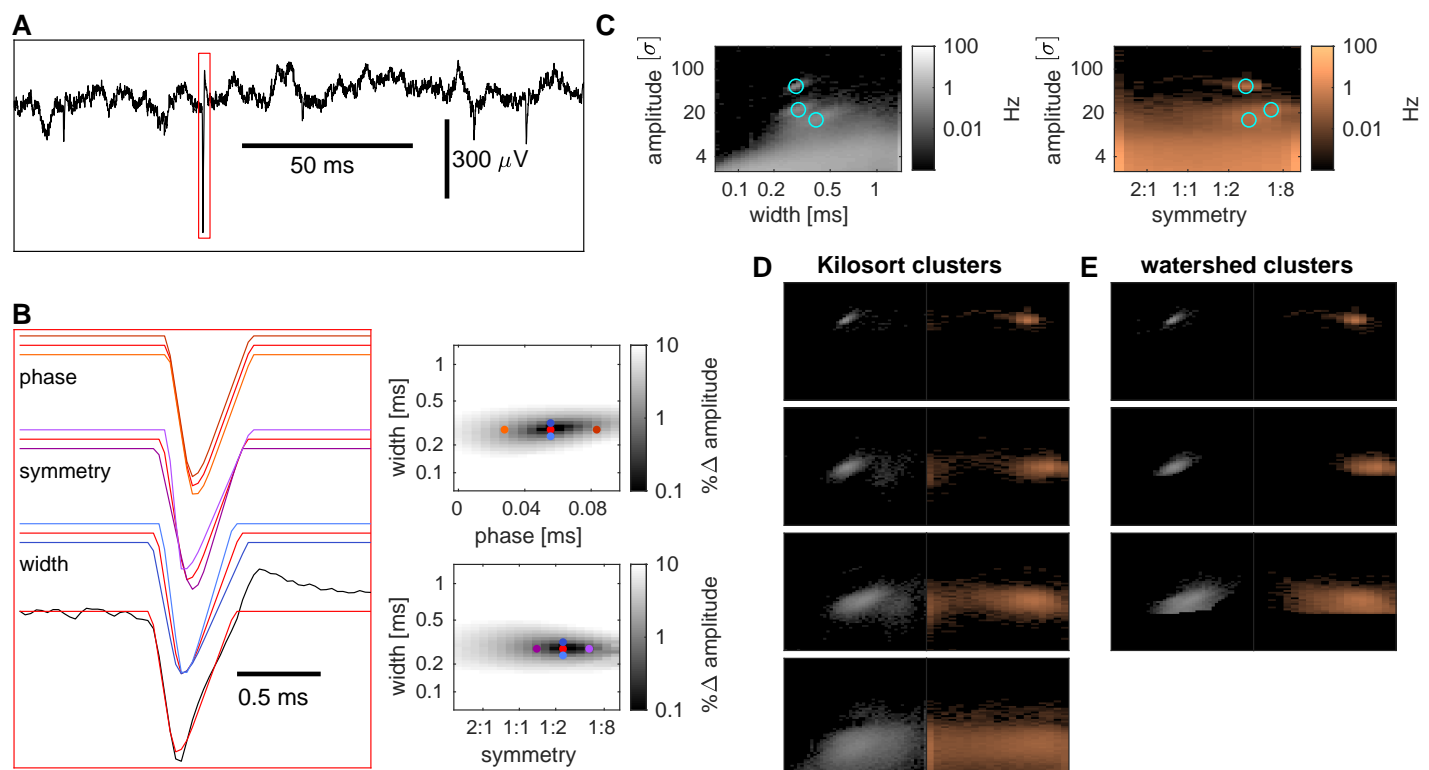


Figure 8. Spike detection and sorting. Raw voltage traces from single electrodes (**A**) are matched in a sliding window to a set of triangular, unimodal templates (examples in **B**, upper left) differing in width, symmetry and phase offset. Local maxima of template - raw trace matches in this parameter space (right plots, dots colored as in left panel) are then detected as putative spikes with a shape characterized by the corresponding width, symmetry and signal power (dot product of template and raw trace). (**C**) Histograms of shapes for an example electrode and recording (marginal distributions). Locations of clusters determined by a watershed algorithm and a recording are marked with cyan circles. (**D**) Shapes of events detected by Kilosort on the same electrode, grouped into clusters by an automated procedure. (**E**) Clusters determined by the watershed algorithm (corresponding to the cyan circles in (**C**)).

Figure 8-source data 1. Source data to generate this Figure

748 radius was larger than 5 bins. To avoid instances where the watershed algorithm would turn indi-
 749 vidual voxels into clusters, we determined a sliding median across 3x3x3 voxels. We further note
 750 that there is a dependency between the recording duration and the resolution of this method (i.e.
 751 higher resolution for longer recordings).

752 For clusters obtained from the watershed algorithm (using a three-dimensional 18-connected
 753 neighborhood), we excluded clusters that systematically had extreme values for spike width or
 754 symmetry or very low amplitudes. Specifically, we ensured that clusters had their center at least
 755 half a standard deviation above the lowest or below the highest bin. As there were many events
 756 with wide shapes, we lowered the exclusion threshold for wide spikes to half a standard deviation
 757 below the second highest bin. For amplitudes we included clusters with an amplitude of at least
 758 half a standard deviation above 2.7σ for lowest spike widths and 5.9σ for highest widths (linear
 759 cutoff in the histograms).

760 To show that these spike clusters indeed corresponded to units found in a conventional spike
 761 sorting approach, we sorted spikes with a widely used spike sorting algorithm (Kilosort, *Pachitariu*
 762 *et al.* (2016)). For that, we used a low threshold for splitting clusters in the Kilosort algorithm
 763 and extracted the shapes of the corresponding spikes from our template matching strategy. This
 764 allowed us to perform the manual step of merging clusters in an automated procedure, using the
 765 Jensen-Shannon divergence between shape histograms as a distance metric.

766 We obtained three dimensional histograms of shape parameters for spikes from each Kilosort

767 cluster (Figure 8 D). We compared Kilosort clusters to clusters obtained by running the watershed
768 algorithm on shape histograms and found a good match for high amplitude clusters (Figure 8 E).
769 The latter clusters were (by construction) better localized in our histograms and we decided to use
770 them instead of Kilosort clusters in the following analyses.

771 Possible extensions

772 We implemented the spike sorting for the case of single, isolated electrodes. An extension to dense
773 arrays is beyond the scope of this article, but we will briefly discuss potential implementation issues
774 here.

- 775 1. Linear arrays/stereotrodes: can be treated as another dimension, like the phase. This just
776 requires one to set a spatial extent of spikes, creating spatially shifted templates. With this
777 method, one could determine maxima at each time frame for each spatial shift, and do a
778 recursive maximization in a second step to obtain spatially isolated maxima.
- 779 2. Spatial grids: memory constraints on the GPU will currently require chunking the array into
780 rows of electrodes.

781 Our current implementation does not include a template generation and matching step, poten-
782 tially resulting in suboptimal detection performance. A potential improvement, while still avoiding
783 the baseline issue, could be to generate templates, smooth them with a kernel and generate tem-
784 plate versions with different widths and phases by interpolation. We would need to normalize the
785 templates to unit power and reduce positive (repolarization) parts of the templates (e.g. divide
786 by 2), to reduce a potential baseline effect. Then we would replace the predefined templates of
787 a given cluster (obtained from the watershed algorithm) with these templates, while keeping the
788 other predefined templates as alternative options (for events that do not match a particular tem-
789 plate). Next, we could rerun the detection with the modified set of templates, considering events
790 which are best matching the inserted templates as spikes.

791 Cross-session merges

792 We computed pairwise Jensen-Shannon divergences between existing clusters from the previous
793 2 sessions and clusters from the current session allowing for small shifts in amplitude, width and
794 symmetry for a penalty. Specifically, we did multiply the Jensen-Shannon divergence with the in-
795 verse of Hanning kernels with a half-width of 7 (for amplitude) and 3 (width and symmetry) bins.
796 Each cluster from the current session was then merged with the existing cluster with the smallest
797 Jensen-Shannon divergence if this was below a threshold of $0.3 \ln(2)$, otherwise it was labeled as a
798 new cluster. To allow for slow temporal drifts, the merged cluster was then assigned a shape den-
799 sity equal to the average of the previous and current density (resulting in effective down-weighting
800 of earlier densities).

801 Motion direction tuning

802 Tuning of spiking activity to the motion direction of a visual stimulus was examined as a function of
803 the width and amplitude of spike shapes, rather than for well isolated clusters, to systematically in-
804 vestigate how much of the low amplitude events was affected by visual stimulation. To this aim, we
805 marginalized over the symmetry parameter of spike shapes, and used a sliding window of 5x5 pix-
806 els for amplitudes and widths, to obtain samples of spikes around each spike width and amplitude.
807 For a temporal window from 20-470 ms from stimulus onset, we computed mean and standard de-
808 viation of spike counts for trials from each stimulus condition, excluding the 3 highest and lowest
809 spike counts from the analysis for robustness of the estimate. The difference between opposing
810 motion directions in the stimulus was then divided by the root mean squared standard deviations
811 to obtain a sensitivity index for each direction. We maximized the sensitivity index across motion
812 directions and, for a sample session, visualized the argument of the maximum as tuned direction
813 in Figure 3J and the maximum value as sensitivity index in Figure 3I. To average these sensitivity in-
814 dices and directions across sessions, we treated the tuning in each session as vectors in the tuned

815 direction with a length equal to the sensitivity index, and averaged them, to obtain an interpolated
816 tuning direction and averaged sensitivity index, shown in Figures 3 B, F and A, E, respectively. To
817 obtain a region of interest for analysis of all stimulus dependent events found on a given electrode,
818 we thresholded the averaged sensitivity indices at 0.3 and determined connected regions exceed-
819 ing this threshold. The largest connected region was then used as a region of interest (outlined in
820 Figures 3 A,B,E,F,I,J) for the cross session analysis performed in Figure 3 C, D, G, H, as well as the
821 single session spike time histograms and tuning curves in Figure 3 K, L. All spikes within that region
822 of interest were used to compute spike time histograms with a bin width of 1 ms and temporally
823 smoothed with an 20 ms wide Hanning kernel (Figure 3 C,G,K).

824 To see how tuning responses at a given electrode site change across sessions, we determined
825 tuning curves for each session (Figure 3 D,H,L). Theoretically, a drift in firing rate or sensitivity could
826 signal a change in coupling between neurons and the electrode, eventually caused by z-drift. Like-
827 wise, due to the spatial organization of area MT, a change in phase could reflect a lateral movement
828 of the electrode.

829 **Cluster survival**

830 Spike shapes were very similar for a large fraction of clusters. It could be that clusters only ap-
831 peared to last across sessions, but in fact represented multiple different clusters that just hap-
832 pened to have matching shapes. Therefore we wanted to test for a bias in longevity for units with
833 common spike shapes. We computed histograms of amplitudes, widths, symmetry and volume of
834 shape clusters, and the average of these quantities for each better-isolated unit across sessions.
835 We then ranked units according to the local density of shape clusters. A lot of short-lived units had
836 uncommonly wide or narrow spike widths or very low firing rates. We therefore performed a rank-
837 ing, which only included units with an average width between 0.1 – 0.5 ms and an average firing
838 rate above 0.5 Hz and assigned the excluded units the ranks of the next lowest ranked included
839 unit. This was not done to exclude units from our analysis of the relation between spike waveform
840 uniqueness and lifetime, but to group them more evenly. For all units with ranks smaller than a
841 given percentile, we then estimated the conditional probability that a unit was lost in the subse-
842 quent session after having survived at least until that session (N). With l_i denoting the measured
843 lifetimes of units, and Θ the Heaviside step function, that probability estimate was

$$844 \hat{p}_N = 1 - \frac{\sum_i (l_i - N - 1) \Theta(l_i - N - 1)}{\sum_i (l_i - N) \Theta(l_i - N)}. \quad (1)$$

845 It assumes that after the N-th session, unit losses are described by a Poisson process with a fixed
846 rate. The estimated additional lifetime (in sessions) τ_N was then given by

$$847 \hat{\tau}_N = -\frac{1}{\ln(\hat{p}_N)} \quad (2)$$

848 and shown in Figure 6 A. The same analysis (replacing ‘sessions’ by ‘days’) was performed to assess
849 temporal lifetimes.

850 **Receptive fields**

851 Firing rate responses were averaged across sessions and smoothed using a 41 ms Hanning kernel.
852 Maximum responses were obtained for each stimulus condition and visualized. The receptive field
853 was then determined as the region where the spatially interpolated response exceeded a threshold
854 of twice the interquartile range above the median across conditions. Data were insufficient for
855 estimating the size of the receptive field for individual sessions. To visualize the cross-session
856 variation of receptive field locations, we assumed periodic boundary conditions and calculated the
857 circular mean eccentricity and direction (colored circles in Figure 7 A, G). Temporal firing responses
858 of individual sessions (Figure 7 D, J) were smoothed using an 18 ms Hanning kernel.

857 **Figure supplements**

858 **Figure 4 — False discovery rate estimates**

859 Spike shapes from different neurons can be similar, or even indistinguishable. To estimate how
860 often we would falsely match a cluster from different units, we tried to match each cluster with
861 clusters found on different channels within 3 sessions before and after its detection. The fraction
862 of chance matches obtained from pairwise comparisons was then scaled by the number of clusters
863 found on the same electrode to obtain an expected number of chance matches. This estimate
864 assumes that the cluster would in fact be absent in the subsequent recording session (and would be
865 lower otherwise). We further determined the dissimilarity threshold at which each pair of clusters
866 would be matched to obtain a threshold dependence of the (pairwise) fraction of chance matches.

867 **Figure 4 — Statistics for long-term units**

868 We determined variations in spike amplitude, rate and inter-spike intervals for long-term units
869 by estimating relative standard deviations. For inter-spike intervals, specifically, we focused on
870 short intervals, as these would more likely reflect intrinsic dynamics of a single neuron, rather
871 than overall network behavior or stimulus dependent responses. In addition, these would also be
872 more robust to potential contamination with noise.

873 We computed spike triggered spike count histograms in an interval from 0.2 - 50 ms after a
874 spike. The first 0.2 ms were ignored as it would merely reflect noise in a few particularly noisy
875 sessions, which were not the subject of this analysis. The histograms were converted into firing
876 rates, smoothed using a 2 ms Hanning window, and normalized by the estimated firing rate of a
877 given session, yielding an instantaneous, relative firing rate. Bursts of spikes would be reflected by
878 an increased instantaneous firing rate shortly after a spike. For quantification, we measured the
879 maximum of the instantaneous, relative firing rate, which was referred to as 'burstiness' in Figure
880 4 – Figure supplements 3 C and 4 C. As an estimate for a relative refractory period, we computed
881 the temporal lag after a spike required to reach 3/4 of this maximum instantaneous firing rate.

882 As a summary statistic, we computed the fraction of the total variance across all clusters (from
883 either group of long-term units), that the variation within units (and across sessions) could explain.
884 This analysis was performed with logarithmized values in order to more equally weight clusters
885 with lower averages.

886 **Figure 7 — Statistics for aggregate data**

887 This analysis aimed at testing whether receptive field locations of identified units were consistent
888 over time. Due to the retinotopic organization of area MT and the small size of the array, we
889 expected similar receptive field locations across the array. Importantly, our sampling of space was
890 relatively sparse and not perfectly homogeneous (few (i.e. 0-10) trials per condition). Additionally,
891 there were few spikes per trial, as we analysed spiking in a short temporal window from 20 to
892 120 ms after stimulus onset.

893 To obtain a robust estimate of RF location with a high spatial resolution, we converted the sam-
894 pled eccentricity and direction to unit vectors on a circle, to perform circular statistics (compute a
895 resultant vector and compare to a uniform Poisson noise model). This approach may distort actual
896 RF locations, but in the same manner for every dataset, and can therefore be used for comparing
897 responses across sessions at a higher resolution. Specifically, we estimated receptive field loca-
898 tions by mapping the 5x7 grid of stimulus eccentricities and directions to circular variables equally
899 spaced on unit circles. Summing up response vectors for different stimuli allowed forming a re-
900 sultant vector with approximate multivariate Gaussian distribution for uniform responses (as null
901 hypothesis), with a variance given by half the number of spikes in each of the 4 dimensions.

902 To account for different trial numbers for different conditions, we smoothed responses and
903 trial numbers across directions and eccentricities using a [0.25 0.5 0.25] kernel (to ensure that
904 there were no conditions without trials). We normalized each condition to reflect an average, per
905 trial spike count and computed its variance under the assumption of probabilistic firing. Variances

906 were then summed across conditions and divided by 2 (2 dimensions) to obtain an approximation
907 of the variance of (each dimension of) the resultant vector under the null hypothesis.

908 Comparing the resultant vector with the null hypothesis yields two numbers:

909 (1) a sensitivity index, specific for a given receptive field location and independent of the number
910 of trials. When treating the null hypothesis as a noise model and the resultant vector as the signal;
911 both would have a variance of half the number of spikes, and hence the sensitivity index would be
912 the length of the resultant vector divided by the square root of half the number of spikes. To obtain
913 a sensitivity index independent of the number of trials, spike counts and resultant vectors were
914 averaged across trials, allowing to compare individual sessions with the cross-session average.

915 (2) a p-value for accepting the null hypothesis of no spatial modulation. The half squared length
916 of the resultant vector, divided by the total number of spikes is Chi-squared distributed with 4
917 degrees of freedom under the null hypothesis. Computing percentiles yielded p-values for each
918 session.

919 It is a curiosity that units with larger sensitivity indices (Figure 7 A,B, red) tended to have re-
920 ceptive fields closer to the center of the region of detected receptive fields from the population
921 than units with lower sensitivity indices (Figure 7 A,B, blue). We do not have an explanation for this
922 observation, and neither did we have the statistical power to examine it in more detail.

923 Acknowledgments

924 This work was supported by the US BRAIN Initiative (U01 NS094330) to ACH, the University of Texas
925 at Austin (College of Natural Sciences Catalyst Award) to ACH, the National Institute on Drug Abuse
926 (T32 DA018926) to AJL and HCC and the National Eye Institute (T32 EY021462) to AJL. We thank John
927 P. Liska for comments on the manuscript.

928 Competing interests

929 The authors declare that no competing interests exist.

930 References

- 931 **Andersen RA**, Bracewell RM, Barash S, Gnadt JW, Fogassi L. Eye position effects on visual, memory, and saccade-
932 related activity in areas LIP and 7a of macaque. *Journal of Neuroscience*. 1990 Apr; 10(4):1176–1196. <https://www.jneurosci.org/content/10/4/1176>, doi: 10.1523/JNEUROSCI.10-04-01176.1990, publisher: Society for
933 Neuroscience Section: Articles.
- 934
- 935 **Brainard DH**. The Psychophysics Toolbox. *Spatial Vision*. 1997 Jan; 10(4):433–436. [https://brill.com/view/](https://brill.com/view/journals/sv/10/4/article-p433_15.xml)
936 [journals/sv/10/4/article-p433_15.xml](https://brill.com/view/journals/sv/10/4/article-p433_15.xml), doi: 10.1163/156856897X00357, publisher: Brill Section: Spatial Vision.
- 937 **Buccino AP**, Hurwitz CL, Garcia S, Magland J, Siegle JH, Hurwitz R, Hennig MH. Spikelnterface, a uni-
938 fied framework for spike sorting. *eLife*. 2020 Nov; 9:e61834. <https://doi.org/10.7554/eLife.61834>, doi:
939 [10.7554/eLife.61834](https://doi.org/10.7554/eLife.61834), publisher: eLife Sciences Publications, Ltd.
- 940 **Camuñas-Mesa LA**, Quiroga RQ. A Detailed and Fast Model of Extracellular Recordings. *Neural Computation*.
941 2013 Mar; 25(5):1191–1212. https://doi.org/10.1162/NECO_a_00433, doi: 10.1162/NECO_a_00433.
- 942 **Chaure FJ**, Rey HG, Quiroga R. A novel and fully automatic spike-sorting implementation with variable
943 number of features. *Journal of Neurophysiology*. 2018 Jul; 120(4):1859–1871. [https://journals.physiology.org/](https://journals.physiology.org/doi/full/10.1152/jn.00339.2018)
944 [doi/full/10.1152/jn.00339.2018](https://journals.physiology.org/doi/full/10.1152/jn.00339.2018), doi: 10.1152/jn.00339.2018, publisher: American Physiological Society.
- 945 **Chung JE**, Joo HR, Fan JL, Liu DF, Barnett AH, Chen S, Geaghan-Breiner C, Karlsson MP, Karlsson M, Lee
946 KY, Liang H, Magland JF, Pebbles JA, Tooker AC, Greengard LF, Tolosa VM, Frank LM. High-Density,
947 Long-Lasting, and Multi-region Electrophysiological Recordings Using Polymer Electrode Arrays. *Neu-*
948 *ron*. 2019 Jan; 101(1):21–31.e5. <http://www.sciencedirect.com/science/article/pii/S0896627318309930>, doi:
949 [10.1016/j.neuron.2018.11.002](http://www.sciencedirect.com/science/article/pii/S0896627318309930).
- 950 **Chung JE**, Magland JF, Barnett AH, Tolosa VM, Tooker AC, Lee KY, Shah KG, Felix SH, Frank LM, Greengard LF. A
951 Fully Automated Approach to Spike Sorting. *Neuron*. 2017 Sep; 95(6):1381–1394.e6. [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0896627317307456)
952 [com/science/article/pii/S0896627317307456](http://www.sciencedirect.com/science/article/pii/S0896627317307456), doi: 10.1016/j.neuron.2017.08.030.

- 953 **Debnath S**, Prins NW, Pohlmeier E, Mylavarapu R, Geng S, Sanchez JC, Prasad A. Long-term stability of neural
954 signals from microwire arrays implanted in common marmoset motor cortex and striatum. *Biomedical*
955 *Physics & Engineering Express*. 2018 Aug; 4(5):055025. <https://doi.org/10.1088/2057-1976/aada67>, doi:
956 10.1088/2057-1976/aada67.
- 957 **Dhawale AK**, Poddar R, Wolff SB, Normand VA, Kopelowitz E, Ölveczky BP. Automated long-term recording and
958 analysis of neural activity in behaving animals. *eLife*. 2017 Sep; 6:e27702. <https://doi.org/10.7554/eLife.27702>,
959 doi: 10.7554/eLife.27702.
- 960 **Dickey AS**, Suminski A, Amit Y, Hatsopoulos NG. Single-Unit Stability Using Chronically Implanted Multielec-
961 trode Arrays. *Journal of Neurophysiology*. 2009 Aug; 102(2):1331–1339. [https://journals.physiology.org/doi/](https://journals.physiology.org/doi/full/10.1152/jn.90920.2008)
962 [full/10.1152/jn.90920.2008](https://journals.physiology.org/doi/full/10.1152/jn.90920.2008), doi: 10.1152/jn.90920.2008.
- 963 **Diggelmann R**, Fiscella M, Hierlemann A, Franke F. Automatic spike sorting for high-density microelectrode
964 arrays. *Journal of Neurophysiology*. 2018 Sep; 120(6):3155–3171. [https://journals.physiology.org/doi/full/10.](https://journals.physiology.org/doi/full/10.1152/jn.00803.2017)
965 [1152/jn.00803.2017](https://journals.physiology.org/doi/full/10.1152/jn.00803.2017), doi: 10.1152/jn.00803.2017, publisher: American Physiological Society.
- 966 **Eastman KM**, Huk AC. PLDAPS: A Hardware Architecture and Software Toolbox for Neurophysiology Requiring
967 Complex Visual Stimuli and Online Behavioral Control. *Frontiers in Neuroinformatics*. 2012; 6. [https://www.](https://www.frontiersin.org/articles/10.3389/fninf.2012.00001/full)
968 [frontiersin.org/articles/10.3389/fninf.2012.00001/full](https://www.frontiersin.org/articles/10.3389/fninf.2012.00001/full), doi: 10.3389/fninf.2012.00001.
- 969 **Fraser GW**, Schwartz AB. Recording from the same neurons chronically in motor cortex. *Journal of Neuro-*
970 *physiology*. 2011 Dec; 107(7):1970–1978. <https://journals.physiology.org/doi/full/10.1152/jn.01012.2010>, doi:
971 10.1152/jn.01012.2010.
- 972 **Fu TM**, Hong G, Viveros RD, Zhou T, Lieber CM. Highly scalable multichannel mesh electronics for stable chronic
973 brain electrophysiology. *Proceedings of the National Academy of Sciences*. 2017 Nov; 114(47):E10046–
974 E10055. <https://www.pnas.org/content/114/47/E10046>, doi: 10.1073/pnas.1717695114.
- 975 **Fu TM**, Hong G, Zhou T, Schuhmann TG, Viveros RD, Lieber CM. Stable long-term chronic brain mapping at
976 the single-neuron level. *Nature Methods*. 2016 Oct; 13(10):875–882. [https://www.nature.com/articles/nmeth.](https://www.nature.com/articles/nmeth.3969)
977 [3969/](https://www.nature.com/articles/nmeth.3969), doi: 10.1038/nmeth.3969.
- 978 **Hilgen G**, Sorbaro M, Pirmoradian S, Muthmann JO, Kepiro IE, Ullo S, Ramirez CJ, Encinas AP, Maccione A,
979 Berdondini L, Murino V, Sona D, Zanacchi FC, Sernagor E, Hennig MH. Unsupervised Spike Sorting for Large-
980 Scale, High-Density Multielectrode Arrays. *Cell Reports*. 2017 Mar; 18(10):2521–2532. [https://www.cell.com/](https://www.cell.com/cell-reports/abstract/S2211-1247(17)30236-X)
981 [cell-reports/abstract/S2211-1247\(17\)30236-X](https://www.cell.com/cell-reports/abstract/S2211-1247(17)30236-X), doi: 10.1016/j.celrep.2017.02.038, publisher: Elsevier.
- 982 **Huk A**, Bonnen K, He BJ. Beyond Trial-Based Paradigms: Continuous Behavior, Ongoing Neural Activity, and
983 Natural Stimuli. *Journal of Neuroscience*. 2018 Aug; 38(35):7551–7558. [https://www.jneurosci.org/content/38/](https://www.jneurosci.org/content/38/35/7551)
984 [35/7551](https://www.jneurosci.org/content/38/35/7551), doi: 10.1523/JNEUROSCI.1920-17.2018, publisher: Society for Neuroscience Section: TechSights.
- 985 **Jackson A**, Fetz EE. Compact Movable Microwire Array for Long-Term Chronic Unit Recording in Cerebral Cortex
986 of Primates. *Journal of Neurophysiology*. 2007 Nov; 98(5):3109–3118. [https://journals.physiology.org/doi/full/](https://journals.physiology.org/doi/full/10.1152/jn.00569.2007)
987 [10.1152/jn.00569.2007](https://journals.physiology.org/doi/full/10.1152/jn.00569.2007), doi: 10.1152/jn.00569.2007.
- 988 **Jun JJ**, Mitelut C, Lai C, Gratiy SL, Anastassiou CA, Harris TD. Real-time spike sorting platform for high-density
989 extracellular probes with ground-truth validation and drift correction. *bioRxiv*. 2017 Jan; p. 101030. <https://www.biorxiv.org/content/10.1101/101030v2>, doi: 10.1101/101030, publisher: Cold Spring Harbor Laboratory
990 [Section: New Results](https://www.biorxiv.org/content/10.1101/101030v2).
991
- 992 **Jun JJ**, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA, Andrei A, Aydın
993 C, Barbic M, Blanche TJ, Bonin V, Couto J, Dutta B, Gratiy SL, Gutnisky DA, Häusser M, Karsh B, Ledochow-
994 itsch P, et al. Fully integrated silicon probes for high-density recording of neural activity. *Nature*. 2017 Nov;
995 551(7679):232–236. <https://www.nature.com/articles/nature24636>, doi: 10.1038/nature24636.
- 996 **Kaneko H**, Tamura H, Suzuki SS. Tracking Spike-Amplitude Changes to Improve the Quality of Multi-
997 neuronal Data Analysis. *IEEE Transactions on Biomedical Engineering*. 2007 Feb; 54(2):262–272. doi:
998 10.1109/TBME.2006.886934.
- 999 **Kloosterman F**, Layton SP, Chen Z, Wilson MA. Bayesian decoding using unsorted spikes in the rat hippocam-
1000 pus. *Journal of Neurophysiology*. 2013 Oct; 111(1):217–227. [https://journals.physiology.org/doi/full/10.1152/](https://journals.physiology.org/doi/full/10.1152/jn.01046.2012)
1001 [jn.01046.2012](https://journals.physiology.org/doi/full/10.1152/jn.01046.2012), doi: 10.1152/jn.01046.2012.
- 1002 **Knöll J**, Pillow JW, Huk AC. Lawful tracking of visual motion in humans, macaques, and marmosets in a natu-
1003 ralistic, continuous, and untrained behavioral context. *Proceedings of the National Academy of Sciences of*
1004 *the United States of America*. 2018; 115(44):E10486–E10494. doi: 10.1073/pnas.1807192115.

- 1005 **Lee J**, Carlson D, Shokri H, Yao W, Goetz G, Hagen E, Batty E, Chichilnisky EJ, Einevoll G, Paninski L. YASS: Yet
1006 Another Spike Sorter. *bioRxiv*. 2017 Jun; p. 151928. <https://www.biorxiv.org/content/10.1101/151928v1>, doi:
1007 10.1101/151928, publisher: Cold Spring Harbor Laboratory Section: New Results.
- 1008 **Linderman MD**, Gilja V, Santhanam G, Afshar A, Ryu S, Meng TH, Shenoy KV. Neural Recording Stability of
1009 Chronic Electrode Arrays in Freely Behaving Primates. In: *2006 International Conference of the IEEE Engineering
1010 in Medicine and Biology Society*; 2006. p. 4387–4391. doi: 10.1109/IEMBS.2006.260814, ISSN: 1557-170X.
- 1011 **Lindén H**, Tetzlaff T, Potjans TC, Pettersen KH, Grün S, Diesmann M, Einevoll GT. Modeling the spatial reach of
1012 the LFP. *Neuron*. 2011; 72(5):859–872. doi: 10.1016/j.neuron.2011.11.006.
- 1013 **Okun M**, Lak A, Carandini M, Harris KD. Long Term Recordings with Immobile Silicon Probes in the Mouse Cor-
1014 tex. *PLoS ONE*. 2016 Mar; 11(3). <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4784879/>, doi: 10.1371/jour-
1015 nal.pone.0151180.
- 1016 **Pachitariu M**, Steinmetz N, Kadir S, Carandini M, D HK. Kilosort: realtime spike-sorting for extracellular elec-
1017 trophysiology with hundreds of channels. *bioRxiv*. 2016 Jun; p. 061481. [https://www.biorxiv.org/content/10.
1018 1101/061481v1](https://www.biorxiv.org/content/10.1101/061481v1), doi: 10.1101/061481.
- 1019 **Paxinos G**, Watson C, Petrides M, Rosa M, Tokuno H. *The Marmoset Brain in Stereotaxic Coordinates*. El-
1020 sevier Academic Press; 2012. <https://espace.curtin.edu.au/handle/20.500.11937/40725>, accepted: 2017-01-
1021 30T14:45:05Z.
- 1022 **Pedreira C**, Martinez J, Ison MJ, Quian Quiroga R. How many neurons can we see with current spike sorting
1023 algorithms? *Journal of Neuroscience Methods*. 2012 Oct; 211(1):58–65. [http://www.sciencedirect.com/science/
1024 article/pii/S0165027012002749](http://www.sciencedirect.com/science/article/pii/S0165027012002749), doi: 10.1016/j.jneumeth.2012.07.010.
- 1025 **Prodanov D**, Delbeke J. Mechanical and Biological Interactions of Implants with the Brain and Their Impact on
1026 Implant Design. *Frontiers in Neuroscience*. 2016; 10. [https://www.frontiersin.org/articles/10.3389/fnins.2016.
1027 00011/full](https://www.frontiersin.org/articles/10.3389/fnins.2016.00011/full), doi: 10.3389/fnins.2016.00011, publisher: Frontiers.
- 1028 **Quian Quiroga R**. What is the real shape of extracellular spikes? *Journal of Neuroscience Meth-
1029 ods*. 2009 Feb; 177(1):194–198. <http://www.sciencedirect.com/science/article/pii/S0165027008005797>, doi:
1030 10.1016/j.jneumeth.2008.09.033.
- 1031 **Rossant C**, Kadir SN, Goodman DFM, Schulman J, Hunter MLD, Saleem AB, Grosmark A, Belluscio M, Denfield
1032 GH, Ecker AS, Tolias AS, Solomon S, Buzsáki G, Carandini M, Harris KD. Spike sorting for large, dense elec-
1033 trode arrays. *Nature Neuroscience*. 2016 Apr; 19(4):634–641. <https://www.nature.com/articles/nn.4268>, doi:
1034 10.1038/nn.4268, number: 4 Publisher: Nature Publishing Group.
- 1035 **Sedaghat-Nejad E**, Herzfeld DJ, Hage P, Karbasi K, Palin T, Wang X, Shadmehr R. Behavioral training of mar-
1036 mosets and electrophysiological recording from the cerebellum. *Journal of Neurophysiology*. 2019 Aug;
1037 122(4):1502–1517. <https://journals.physiology.org/doi/full/10.1152/jn.00389.2019>, doi: 10.1152/jn.00389.2019.
- 1038 **Siegle JH**, López AC, Patel YA, Abramov K, Ohayon S, Voigts J. Open Ephys: an open-source, plugin-based
1039 platform for multichannel electrophysiology. *Journal of Neural Engineering*. 2017 Jun; 14(4):045003. [https:
1040 //doi.org/10.1088/1741-2552/14/4/045003](https://doi.org/10.1088/1741-2552/14/4/045003), doi: 10.1088/1741-2552/14/4/045003.
- 1041 **Todorova S**, Sadtler P, Batista A, Chase S, Ventura V. To sort or not to sort: the impact of spike-sorting on
1042 neural decoding performance. *Journal of Neural Engineering*. 2014 Aug; 11(5):056005. [https://doi.org/10.
1043 1088/1741-2552/11/5/056005](https://doi.org/10.1088/1741-2552/11/5/056005), doi: 10.1088/1741-2552/11/5/056005.
- 1044 **Trautmann EM**, Stavisky SD, Lahiri S, Ames KC, Kaufman MT, O’Shea DJ, Vyas S, Sun X, Ryu SI, Gan-
1045 guli S, Shenoy KV. Accurate Estimation of Neural Population Dynamics without Spike Sorting. *Neu-
1046 ron*. 2019 Jul; 103(2):292–308.e4. <http://www.sciencedirect.com/science/article/pii/S0896627319304283>, doi:
1047 10.1016/j.neuron.2019.05.003.
- 1048 **Vaidya M**, Dickey A, Best MD, Coles J, Balasubramanian K, Suminski AJ, Hatsopoulos NG. Ultra-long term
1049 stability of single units using chronically implanted multielectrode arrays. In: *2014 36th Annual Inter-
1050 national Conference of the IEEE Engineering in Medicine and Biology Society*; 2014. p. 4872–4875. doi:
1051 10.1109/EMBC.2014.6944715, ISSN: 1558-4615.
- 1052 **Yger P**, Spampinato GL, Esposito E, Lefebvre B, Deny S, Gardella C, Stimberg M, Jetter F, Zeck G, Picaud S, Duebel
1053 J, Marre O. A spike sorting toolbox for up to thousands of electrodes validated with ground truth recordings
1054 in vitro and in vivo. *eLife*. 2018 Mar; 7:e34518. <https://doi.org/10.7554/eLife.34518>, doi: 10.7554/eLife.34518,
1055 publisher: eLife Sciences Publications, Ltd.

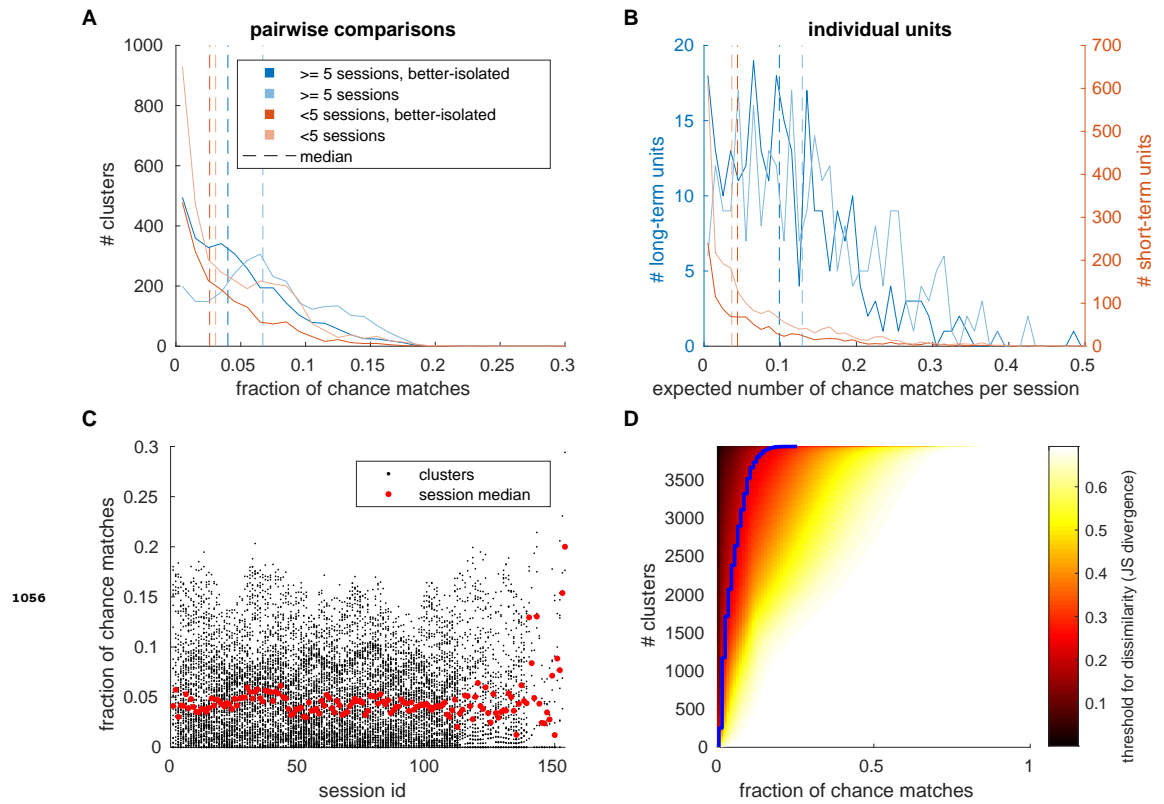


Figure 4-Figure supplement 1. False discovery rate estimates for marmoset J. Spike shapes from different neurons can be similar, or even indistinguishable. To estimate how often we would falsely match a cluster from different units, we tried to match each cluster with clusters found on different channels within 3 sessions before and after its detection. (A) Histograms of the fraction chance matches in pairwise comparisons. Units were classified as in Figure 4 and the corresponding histograms were colored accordingly. Dashed lines mark median values. (B) Histograms of the average expected number of chance matches per session, when accounting for the number of detected clusters on the same electrode. (C) Pairwise false discovery rates across recording sessions. Red dots depict median values for each session. (D) False discovery rates in dependence of the dissimilarity threshold (blue line depicts threshold used in this work). Clusters were sorted according to the fraction of chance matches when using a fixed threshold.

1056

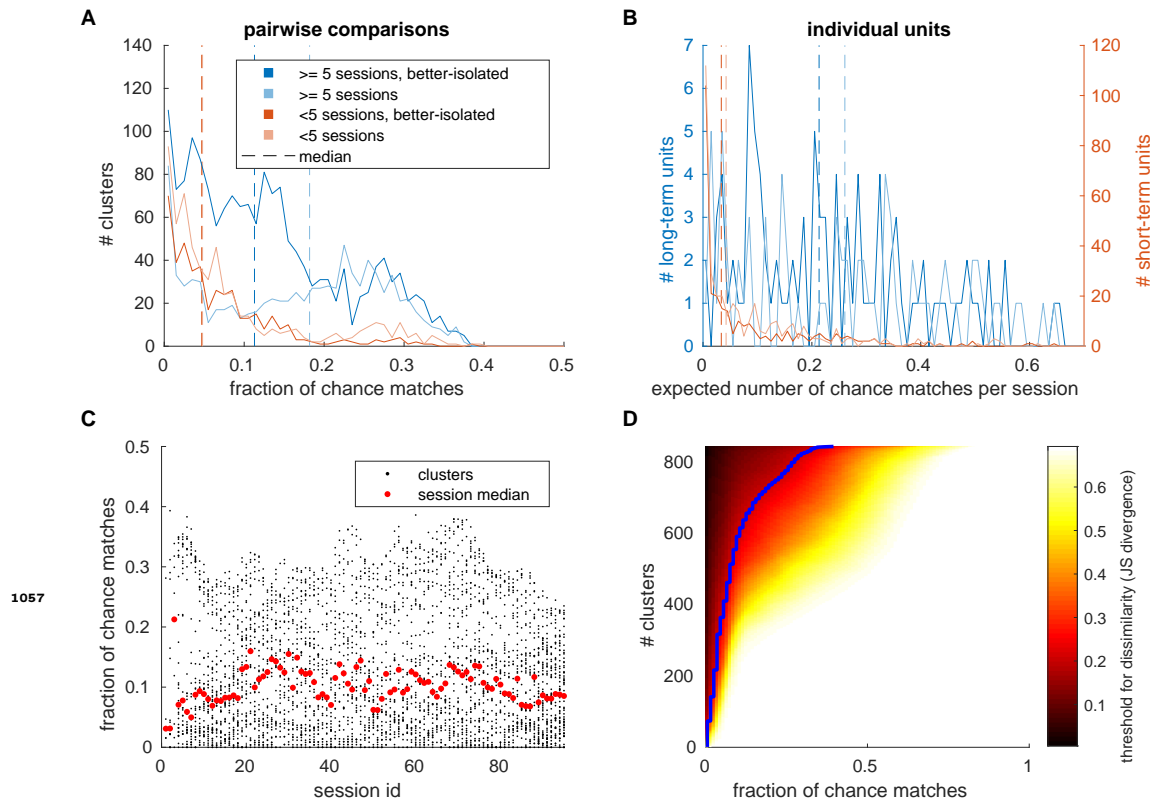


Figure 4-Figure supplement 2. False discovery rate estimates for marmoset B. Spike shapes from different neurons can be similar, or even indistinguishable. To estimate how often we would falsely match a cluster from different units, we tried to match each cluster with clusters found on different channels within 3 sessions before and after its detection. (A) Histograms of the fraction chance matches in pairwise comparisons. Units were classified as in Figure 4 and the corresponding histograms were colored accordingly. Dashed lines mark median values. (B) Histograms of the average expected number of chance matches per session, when accounting for the number of detected clusters on the same electrode. (C) Pairwise false discovery rates across recording sessions. Red dots depict median values for each session. (D) False discovery rates in dependence of the dissimilarity threshold (blue line depicts threshold used in this work). Clusters were sorted according to the fraction of chance matches when using a fixed threshold.

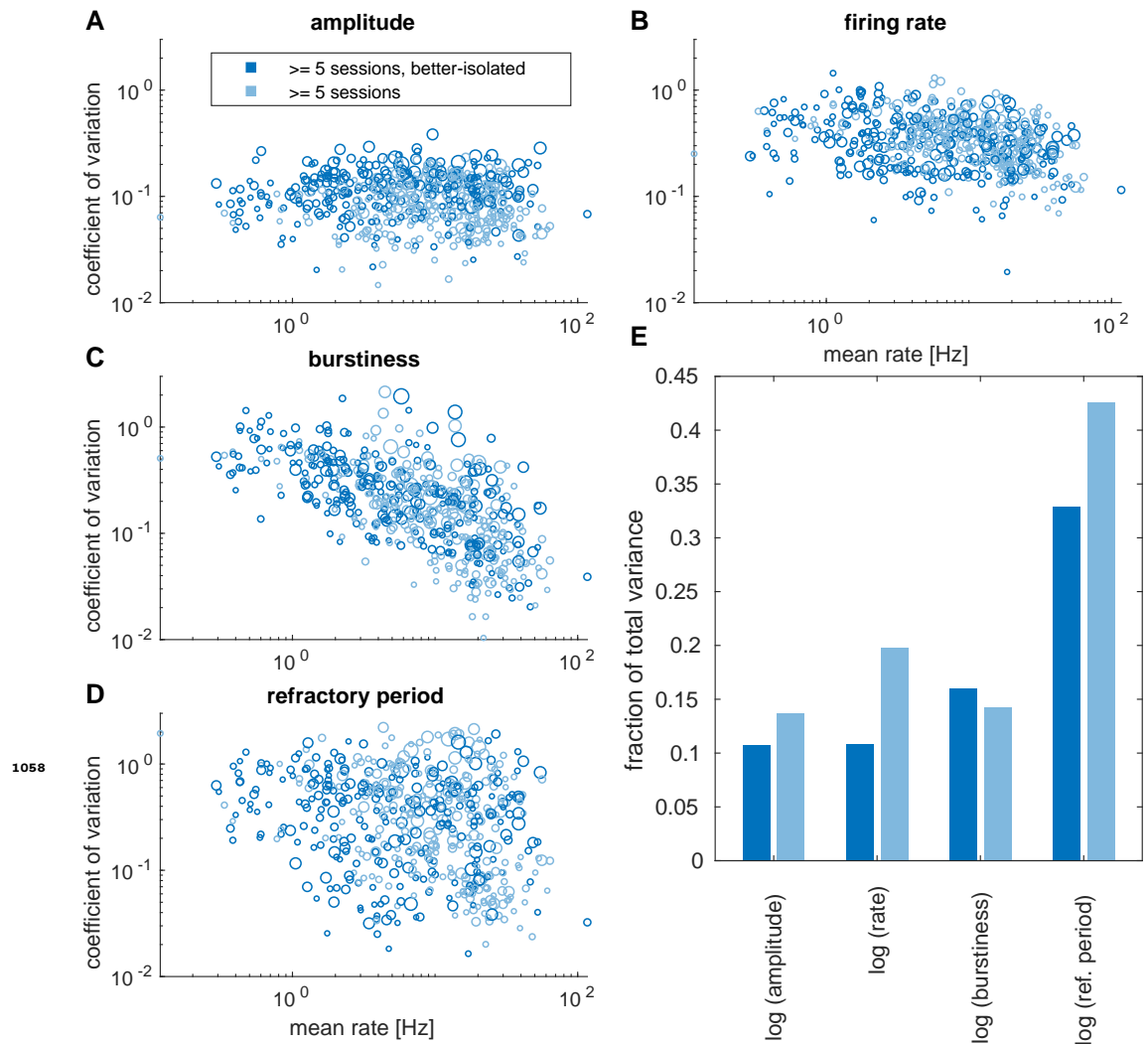


Figure 4-Figure supplement 3. Long-term statistics for marmoset J. (A) Relative amplitude variations of long-term clusters. Larger symbols represent clusters observed in more experimental sessions, darker shades correspond to better-isolated units (as in Figure 4). (B) Relative firing rate variations. (C-D) Averages and variability of relative spike triggered averaged firing rates. To quantify the propensity of spiking in a short window after a spike, we computed spike triggered spike count histograms in an interval from 0.2 - 50 ms after a spike. These were converted into firing rates, smoothed using a 2 ms Hanning window, and normalized by the estimated firing rate of a given session. The maximum relative spike triggered firing rate was termed 'burstiness', and its variability for individual units is shown in (C). A high value would correspond to an increased chance of firing shortly after a spike, and a value around one would reflect no burst firing. As an estimate for a relative refractory period (variability shown in (D)), we computed the temporal lag after a spike required to reach 3/4 of this maximum firing rate. (E) Fraction of the total variance explained by within unit and across session variability. In order to more equally weight clusters with lower averages, this analysis was performed on a logarithmic scale.

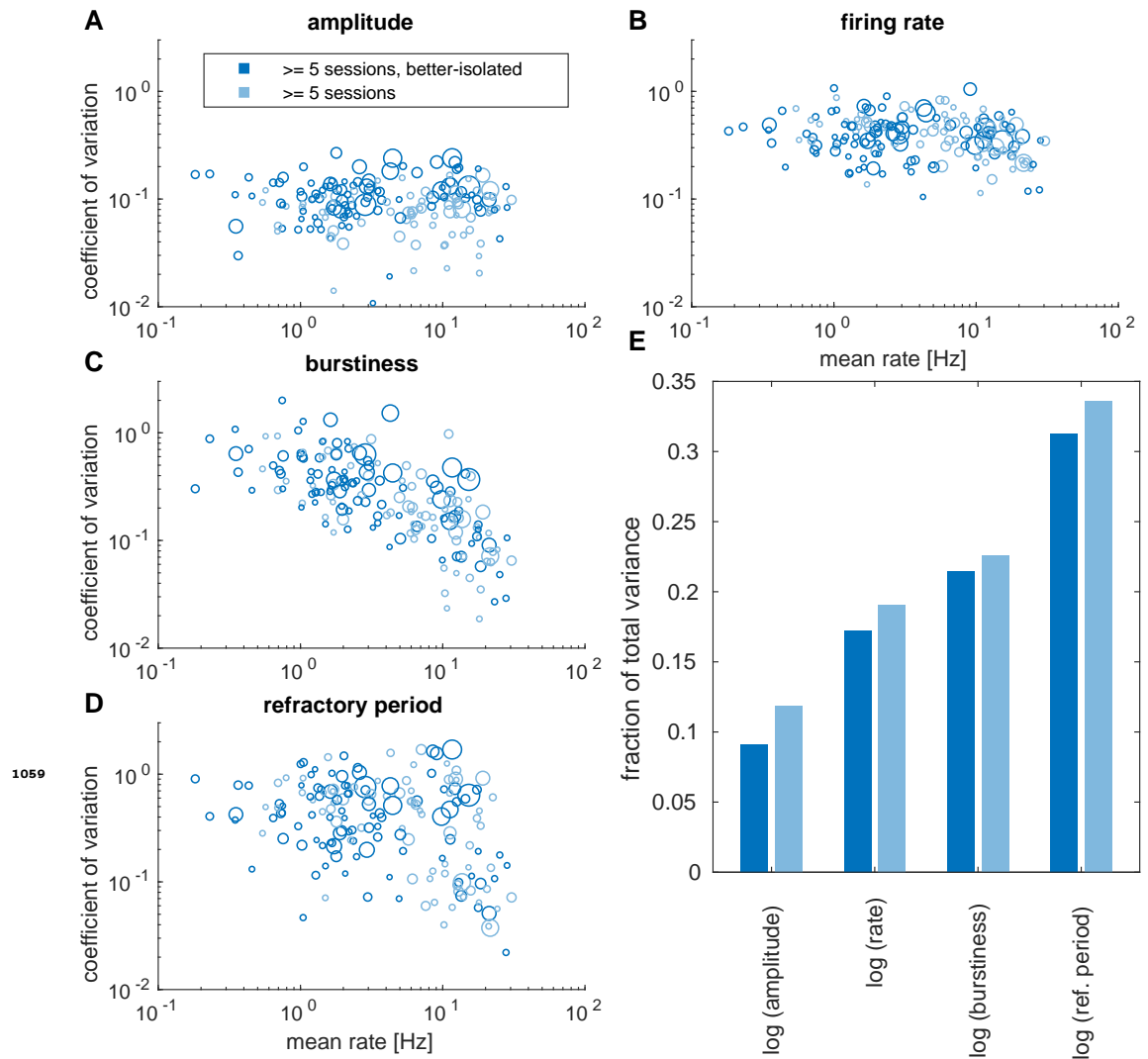


Figure 4-Figure supplement 4. Long-term statistics for marmoset B. (A) Relative amplitude variations of long-term clusters. Larger symbols represent clusters observed in more experimental sessions, darker shades correspond to better-isolated units (as in Figure 4). (B) Relative firing rate variations. (C-D) Averages and variability of relative spike triggered averaged firing rates. To quantify the propensity of spiking in a short window after a spike, we computed spike triggered spike count histograms in an interval from 0.2 - 50 ms after a spike. These were converted into firing rates, smoothed using a 2 ms Hanning window, and normalized by the estimated firing rate of a given session. The maximum relative spike triggered firing rate was termed 'burstiness', and its variability for individual units is shown in (C). A high value would correspond to an increased chance of firing shortly after a spike, and a value around one would reflect no burst firing. As an estimate for a relative refractory period (variability shown in (D)), we computed the temporal lag after a spike required to reach 3/4 of this maximum firing rate. (E) Fraction of the total variance explained by within unit and across session variability. In order to more equally weight clusters with lower averages, this analysis was performed on a logarithmic scale.

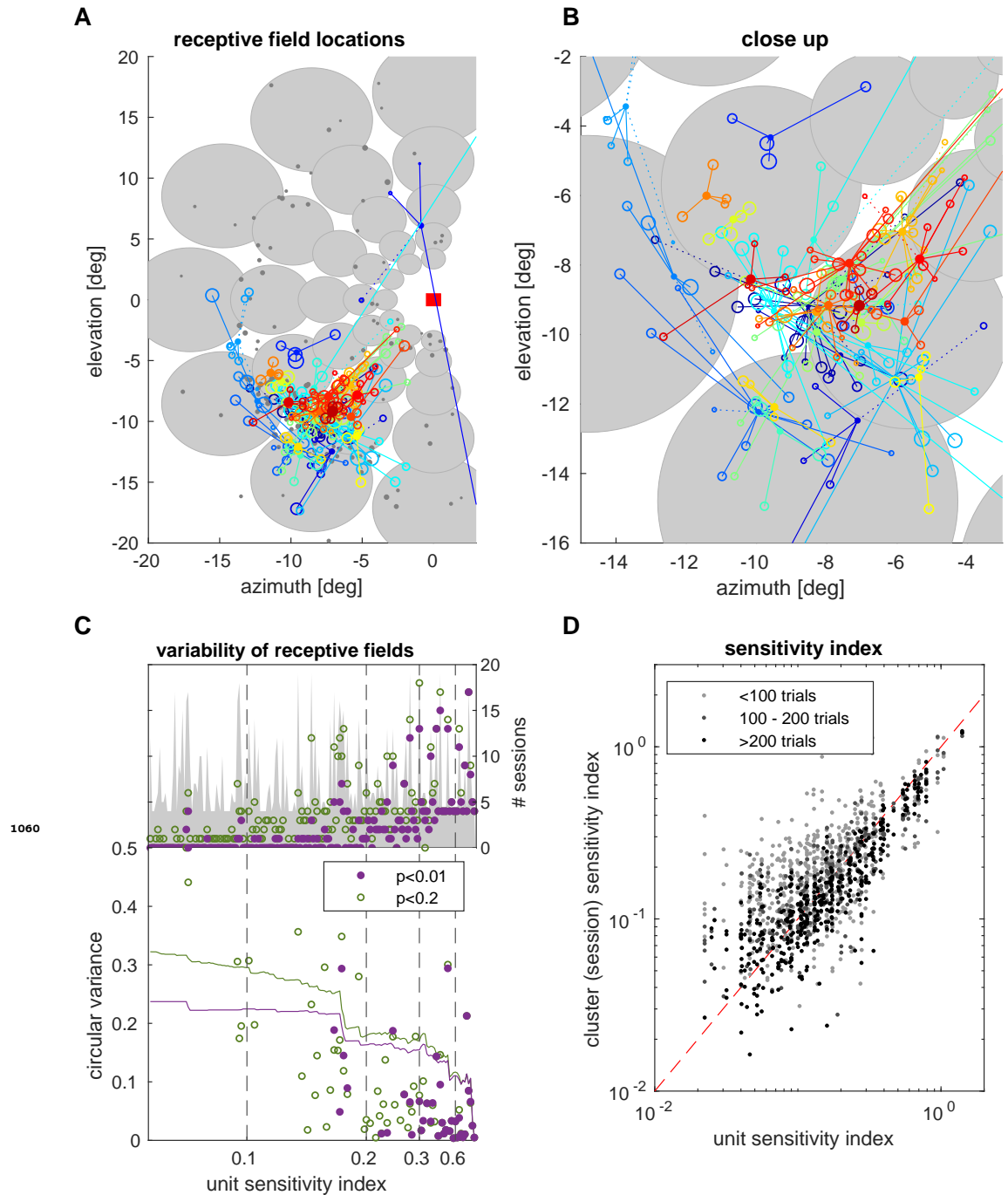


Figure 7-Figure supplement 1. Statistics for aggregate data. Receptive field locations were estimated by mapping the 5x7 grid of stimulus eccentricities and directions to circular variables equally spaced on unit circles. Summing up response vectors for different stimuli allowed forming a resultant vector with approximate Gaussian distribution for uniform responses (as null hypothesis), and mapping the preferred stimulus location back to world coordinates. (A,B) Receptive field locations of units observed for at least 4 sessions with receptive field mapping (filled circles). Size/color relates to sensitivity indices (red: high, blue:low, gray:<0.3). Open circles denote estimated receptive field locations in individual sessions, linked to the corresponding unit with a solid line for sessions with a significant ($p < 0.01$) spatial modulation of firing rates and and dotted line for a tendency ($p < 0.2$) of a spatial modulation. (C) Variation of receptive field locations across at least 4 sessions from the same unit with good ($p < 0.01$, purple) and weak ($p < 0.2$, green) spatial modulation, normalizing individual session resultant vectors and computing the circular variance across sessions. The circular variance of a population of clusters from units with a given minimum sensitivity index is shown as a reference (colored lines). (D) Scatterplot comparing sensitivity indices of units computed across sessions and the corresponding single session estimates.