

1 Reproducibility of in-vivo 2 electrophysiological measurements 3 in mice

4 **International Brain Laboratory***, Kush Banga⁷, Julius Benson¹¹, Niccolò Bonacchi²,
5 Sebastian A Bruijns¹, Rob Campbell¹³, Gaëlle A Chapuis⁵, Anne K Churchland⁶, M
6 Felicia Davatolhagh⁶, Hyun Dong Lee³, Mayo Faulkner⁷, Fei Hu⁹, Julia
7 Hunterberg², Anup Khanal⁶, Christopher Krasniak¹⁰, Guido T Meijer², Nathaniel J
8 Miska⁷, Zeinab Mohammadi¹², Jean-Paul Noel¹¹, Liam Paninski³, Alejandro
9 Pan-Vazquez¹², Noam Roth⁴, Michael Schartner², Karolina Socha⁷, Nicholas A
10 Steinmetz⁴, Karel Svoboda¹⁴, Marsa Taheri⁶, Anne E Urai⁸, Miles Wells⁷, Steven J
11 West⁷, Matthew R Whiteway³, Olivier Winter², Ilana B Witten¹²

***For correspondence:**

anne.churchland@internationalbrainlab.org (AKC); liam.paninski@internationalbrainlab.org (LP); nicholas.steinmetz@internationalbrainlab.org (NAS)

12 ¹Max-Planck-Institute, Tübingen, Germany; ²Champalimaud Foundation, Lisbon,
13 Portugal; ³Columbia University, NY, USA; ⁴University of Washington, WA, USA;
14 ⁵University of Geneva, Switzerland; ⁶University of California Los Angeles, USA;
15 ⁷University College London, UK; ⁸Leiden University, The Netherlands; ⁹University of
16 California, Berkeley, USA; ¹⁰Cold Spring Harbor Laboratory, NY, USA; ¹¹New York
17 University, NY, USA; ¹²Princeton University, NJ, USA; ¹³Sainsbury Wellcome Center,
18 London, UK; ¹⁴Allen Institute for Neural Dynamics WA, USA

19 May 12, 2022
20

21 Abstract

22 Understanding whole-brain-scale electrophysiological recordings will rely on the collective work
23 of multiple labs. Because two labs recording from the same brain area often reach different
24 conclusions, it is critical to quantify and control for features that decrease reproducibility. To
25 address these issues, we formed a multi-lab collaboration using a shared, open-source
26 behavioral task and experimental apparatus. We repeatedly inserted Neuropixels multi-electrode
27 probes targeting the same brain locations (including posterior parietal cortex, hippocampus, and
28 thalamus) in mice performing the behavioral task. We gathered data across 9 labs and developed
29 a common histological and data processing pipeline to analyze the resulting large datasets. After
30 applying stringent behavioral, histological, and electrophysiological quality-control criteria, we
31 found that neuronal yield, firing rates, spike amplitudes, and task-modulated neuronal activity
32 were reproducible across laboratories. To quantify variance in neural activity explained by task
33 variables (e.g., stimulus onset time), behavioral variables (timing of licks/paw movements), and
34 other variables (e.g., spatial location in the brain or the lab ID), we developed a multi-task neural
35 network encoding model that extends common, simpler regression approaches by allowing
36 nonlinear interactions between variables. We found that within-lab random effects captured by
37 this model were comparable to between-lab random effects. Taken together, these results
38 demonstrate that across-lab standardization of electrophysiological procedures can lead to
39 reproducible results across labs. Moreover, our protocols to achieve reproducibility, along with
40 our analyses to evaluate it are openly accessible to the scientific community, along with our
41 extensive electrophysiological dataset with corresponding behavior and open-source analysis

42 code.

43

44 Introduction

45 *Reproducibility* is a cornerstone of the scientific method: a given sequence of experimental meth-
46 ods should lead to comparable results if applied in different laboratories. In some areas of bi-
47 ological and psychological science, however, the reliable generation of reproducible results is a
48 well-known challenge (*Baker, 2016; Voelkl et al., 2020; Li et al., 2021; Errington et al., 2021*). In
49 systems neuroscience at the level of single-cell-resolution recordings, evaluating reproducibility
50 is difficult: experimental methods are sufficiently complex that replicating experiments is techni-
51 cally challenging, and many experimenters feel little incentive to do such experiments since nega-
52 tive results can be difficult to publish. Variability in experimental outcomes has nonetheless been
53 well-documented on a number of occasions. These include the existence and nature of “preplay”
54 (*Dragoi and Tonegawa, 2011; Silva et al., 2015; Ólafsdóttir et al., 2015; Grosmark and Buzsáki,*
55 *2016; Liu et al., 2019*), the persistence of place fields in the absence of visual inputs (*Hafting et al.,*
56 *2005; Barry et al., 2012; Chen et al., 2016; Waaga et al., 2022*), and the existence of spike-timing de-
57 pendent plasticity (STDP) in nematodes (*Zhang et al., 1998; Tsui et al., 2010*). In the latter example,
58 variability in experimental results arose from whether the nematode being studied was pigmented
59 or albino, an experimental feature that was not originally known to be relevant to STDP. This high-
60 lights that understanding the source of experimental variability can facilitate efforts to improve
61 reproducibility.

62 For electrophysiological recordings, several efforts are currently underway to document this
63 variability and reduce it through standardization of methods (*de Vries et al., 2020; Siegle et al.,*
64 *2021*). These efforts are promising, in that they suggest that when approaches are standardized
65 and results undergo quality control, observations conducted within a single organization can be
66 reassuringly reproducible. However, this leaves unanswered whether observations made in sepa-
67 rate, individual laboratories are reproducible when they likewise use standardization and quality
68 control. Answering this question is critical since most neuroscience data is collected within small,
69 individual laboratories rather than large-scale organizations.

70 We have previously addressed the issue of reproducibility in the context of mouse psychophys-
71 ical behavior, by training 140 mice in 7 laboratories and comparing their learning rates, speed, and
72 accuracy in a simple binary visually-driven decision task. We demonstrated that standardized pro-
73 tocols can lead to highly reproducible behavior (*The International Brain Laboratory et al., 2021*).
74 Here, we build on those results by measuring within- and across-lab variability in the context of
75 intra-cerebral electrophysiological recordings. We repeatedly inserted Neuropixels multi-electrode
76 probes (*Jun et al., 2017*) targeting the same brain regions (including posterior parietal cortex, hip-
77 pocampus, and thalamus) in mice performing the behavioral task from (*The International Brain*
78 *Laboratory et al., 2021*). We gathered data across 9 different labs and developed a common histo-
79 logical and data processing pipeline to analyze the resulting large datasets.

80 After applying stringent behavioral, histological, and electrophysiological quality-control crite-
81 ria, features such as neuronal yield, firing rate, and normalized LFP power were reproducible across
82 laboratories; their within-lab averages did not significantly deviate from the mean across labs. Sim-
83 ilarly, the proportions of cells modulated by task events was largely reproducible across labs, as
84 was the Fano Factor, a measure of neural variability. Finally, to quantify variance in neural activ-
85 ity explained by task variables (e.g., stimulus onset time), behavioral variables (timing of licks/paw
86 movements), and other variables (e.g., spatial location in the brain or the lab ID), we developed a
87 multi-task neural network encoding model that extends common, simpler regression approaches
88 by allowing nonlinear interactions between variables. Again, we found that within-lab random ef-
89 fects captured by this model were comparable to between-lab random effects. Taken together,
90 these results suggest that across-lab standardization of electrophysiological procedures can lead

91 to reproducible results across laboratories.

92 Results

93 Repeated-site recordings in the same task across multiple labs

94 To quantify reproducibility across electrophysiological recordings, we set out to establish standard-
95 ized procedures across the International Brain Laboratory (IBL) and to test whether this standard-
96 ization was successful. Nine IBL labs collected Neuropixels recordings from one repeated site,
97 targeting the same stereotaxic coordinates, during a standardized decision-making task in which
98 head-fixed mice reported the perceived position of a visual grating (*The International Brain Lab-*
99 *oratory et al., 2021*). The experimental pipeline was standardized across labs, including surgical
100 methods, behavioral training, recording procedures, histology, and data processing (Figure 1a, b);
101 see Methods for full details. In each experiment, Neuropixels 1.0 probes were inserted, targeted
102 at -2.0 mm AP, -2.24 mm ML, 4.0 mm DV relative to bregma; 15° angle (Figure 1c). This site was
103 selected because it encompasses brain regions implicated in visual decision-making, including vi-
104 sual area A (*Najafi et al., 2020; Harvey et al., 2012*), dentate gyrus, CA1, (*Turk-Browne, 2019*), and
105 thalamic nuclei LP and PO (*Saalmann and Kastner, 2011; Roth et al., 2016*).

106 Probe placement contributes to experimental variability

107 As a first test of experimental reproducibility, we assessed variability in Neuropixels probe place-
108 ment around the planned repeated site location. Brains were perfusion-fixed, dissected, and im-
109 aged using serial section 2-photon microscopy for 3D reconstruction of probes (Figure 2a). Whole
110 brain auto-fluorescence data was aligned to the Allen Common Coordinate Framework (CCF) (*Wang*
111 *et al., 2020*) using an elastix-based pipeline (*Klein et al., 2010*) adapted for mouse brain registra-
112 tion (*West, 2021*). CM-Dil labelled probe tracks were manually traced in the 3D volume. Trajectories
113 obtained from our stereotaxic system and traced histology were then compared to the planned
114 trajectory (Figure 2a,b, Figure 2b; supp. 1). To measure probe track variability, traced probe tracks
115 were linearly interpolated (Figure 2c).

116 Variability in brain insertions can be assessed by probe placement at the brain surface, and by
117 probe angle. Probe placement at the brain surface comprises two components. The first, 'target-
118 ing variability,' was obtained by calculating the difference between the planned and actual probe
119 placement, measured with the micro-manipulator at the time of recording (Figure 2d). Targeting
120 variability is expected to be non-zero because experimenters sometimes move probes slightly from
121 the planned location to avoid blood vessels or irregularities (Figure 2d, top, total mean displace-
122 ment = $115 \mu\text{m}$, exclusion criteria passed mean displacement = $72 \mu\text{m}$). Reproducibility of targeting
123 variability across labs was evaluated via a permutation test: values were shuffled between the lab
124 identities 10,000 times, and the original targeting variability mean per lab distribution was com-
125 pared to all permuted distributions to compute a p-value. Targeting variability shows no signifi-
126 cant effect across laboratories across all probes (Figure 2d, bottom), permutation test p-value for
127 all probes $p=0.2118$). When applying our exclusion criteria, including the anatomical requirement
128 that the probe must record from three of our five repeated site brain regions, the computed p-
129 value increased (Figure 2d, bottom), permutation test p-value for exclusion criteria passed probes
130 $p=0.2295$), indicating the data are more likely from the same distribution. Thus, targeting repro-
131 ducibility is enhanced with appropriate anatomical exclusion criteria.

132 The second component of probe placement variability in brain insertions is 'geometrical vari-
133 ability.' Geometrical variability was obtained by calculating the difference between our planned
134 position and the final identified probe position obtained from the reconstructed histology. This
135 encompasses the targeting variance above, plus anatomical differences and errors in defining
136 the stereotaxic coordinate system, including residual errors from a mismatch in skull landmarks
137 and underlying brain structure. Geometrical variability was likewise non-zero (Figure 2e, top, total
138 mean displacement = $392 \mu\text{m}$, exclusion criteria passed mean displacement = $253 \mu\text{m}$) with some

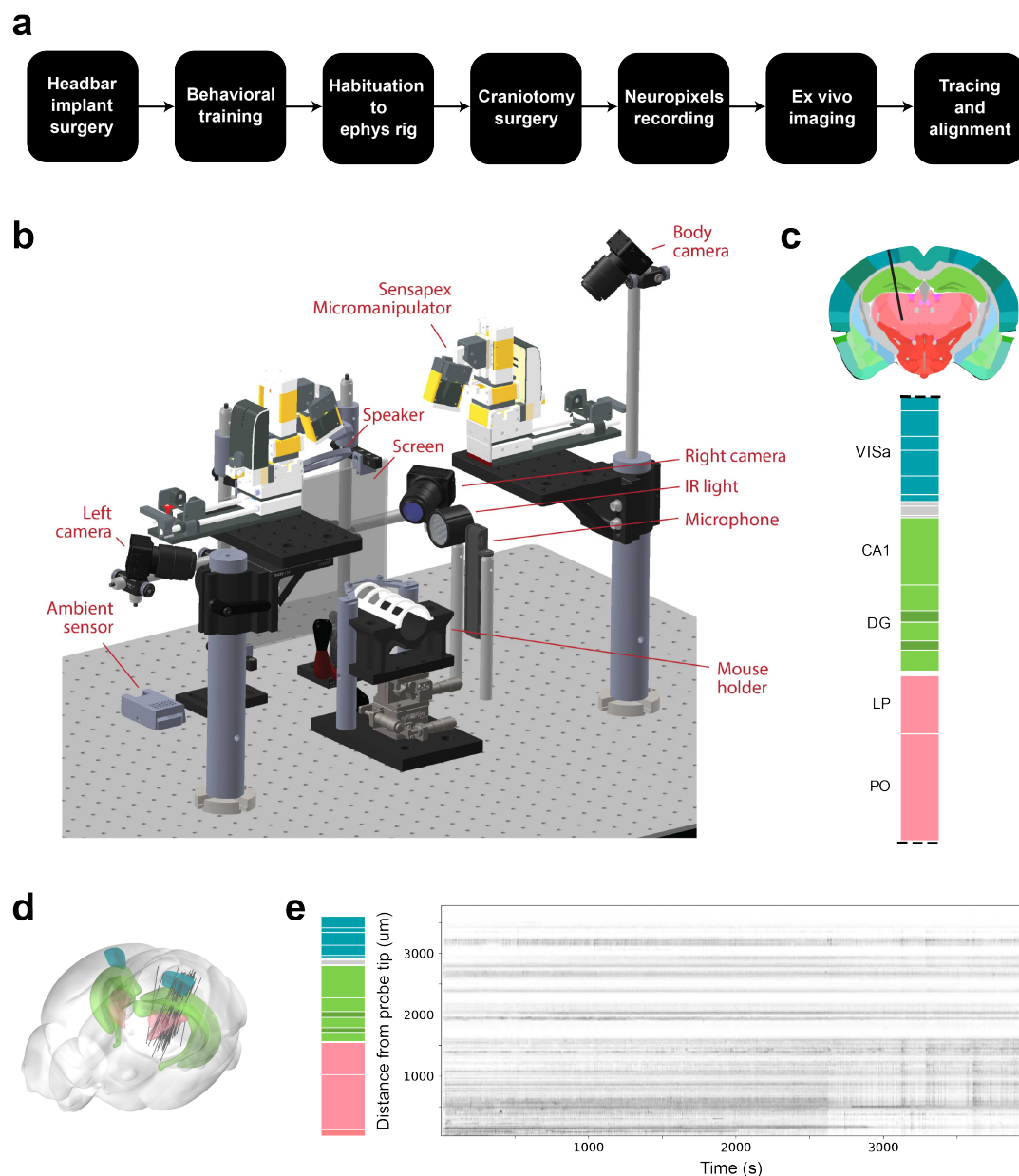


Figure 1. Standardized experimental pipeline and apparatus; location of the repeated site. **a**, The pipeline for electrophysiology experiments. **b**, Drawing of the experimental apparatus. **c**, Location and brain regions of the repeated site. VISa: Visual Area A; CA1: Hippocampal Field CA1; DG: Dentate Gyrus; LP: Lateral Posterior nucleus of the thalamus; PO: Posterior Nucleus of the Thalamus. **d**, Acquired repeated site trajectories shown within a 3D brain schematic. **e**, Raster plot from one example session.

Figure 1-Figure supplement 1. Detailed experimental pipeline for the Neuropixels experiment.

Figure 1-Figure supplement 2. Spiking activity qualitatively appears heterogeneous across recordings.

139 individual insertion locations up to 1500 μm from the planned coordinate. Assessing geometrical
 140 variability for all probes with permutation testing revealed no significant effect across laboratories
 141 (Figure 2e, bottom, permutation test p-value for all probes $p=0.1974$), which produced a higher
 142 p-value after the application of our exclusion criteria (Figure 2e, bottom, permutation test p-value
 143 for exclusion criteria passed probes $p=0.0.5499$). This demonstrates that after histology recon-
 144 struction, the reproducibility of probe placement is enhanced across labs for the brain insertion
 145 coordinate with the application of anatomical exclusion criteria.

146 The final way to assess variability in brain insertions is via 'angle variability,' also calculated from
147 the histological reconstructions. We observed a consistent mean displacement from the planned
148 angle in both medio-lateral (ML) and anterior-posterior (AP) angles (mean difference in angle from
149 planned: 7 degrees, Figure 2f, top). AP angle differences can be explained by the different ori-
150 entation of the CCF and the stereotaxic coordinate system; ML differences may result from the
151 histological asmples being compressed in the DV direction compared to the CCF. The difference
152 in histology angle to planned probe placement was assessed with permutation testing across labs,
153 and shows a significant difference with our exclusion criteria applied (Figure 2f, bottom, permu-
154 tation test p-value for all probes p=0.1993; permutation test p-value for exclusion criteria passed
155 probes p=0.0491). This significant result can be explained by the repeated use of the same rig
156 and micromanipulator angle within each laboratory, resulting in reduced variability in probe angle
157 within labs versus across labs.

158 To determine the extent that anatomical differences drive geometrical variability, we used the
159 micro-manipulator to histology distance at the brain surface and regressed this measurement
160 against animal weight. This easily measured parameter should correlate with mouse brain size
161 and provide a quantifiable predictor of anatomical differences. No such correlation was identified
162 ($R^2 < 0.01$), indicating differences between CCF and mouse brain sizes are not the major cause of
163 variance. We therefore surmise that geometrical variance in probe placement at the brain surface
164 is driven by inaccuracies in defining the stereotaxic coordinate system, including discrepancies
165 between skull landmarks and the underlying brain structures.

166 In conclusion, targeting, geometrical and angle variability revealed lab-to-lab differences that
167 can hinder reproducibility. To control this variability we applied a "targeting" exclusion criterion,
168 which discarded insertions from further analysis when they failed to include sites from at least 3
169 of the 5 selected areas. This exclusion criterion improved the reproducibility of probe placement
170 at the brain surface, and was used in all subsequent analyses. Probe angle reproducibility was not
171 improved with the exclusion criterion, and this appears to be driven by variance between recording
172 rigs repeatedly used for probe placement within labs. We were unable to identify a prescriptive
173 analysis to predict probe placement accuracy, which may reflect that the major driver of probe
174 placement variance derives from differences in skull landmarks used for establishing the coordi-
175 nate system, and the underlying brain structures.

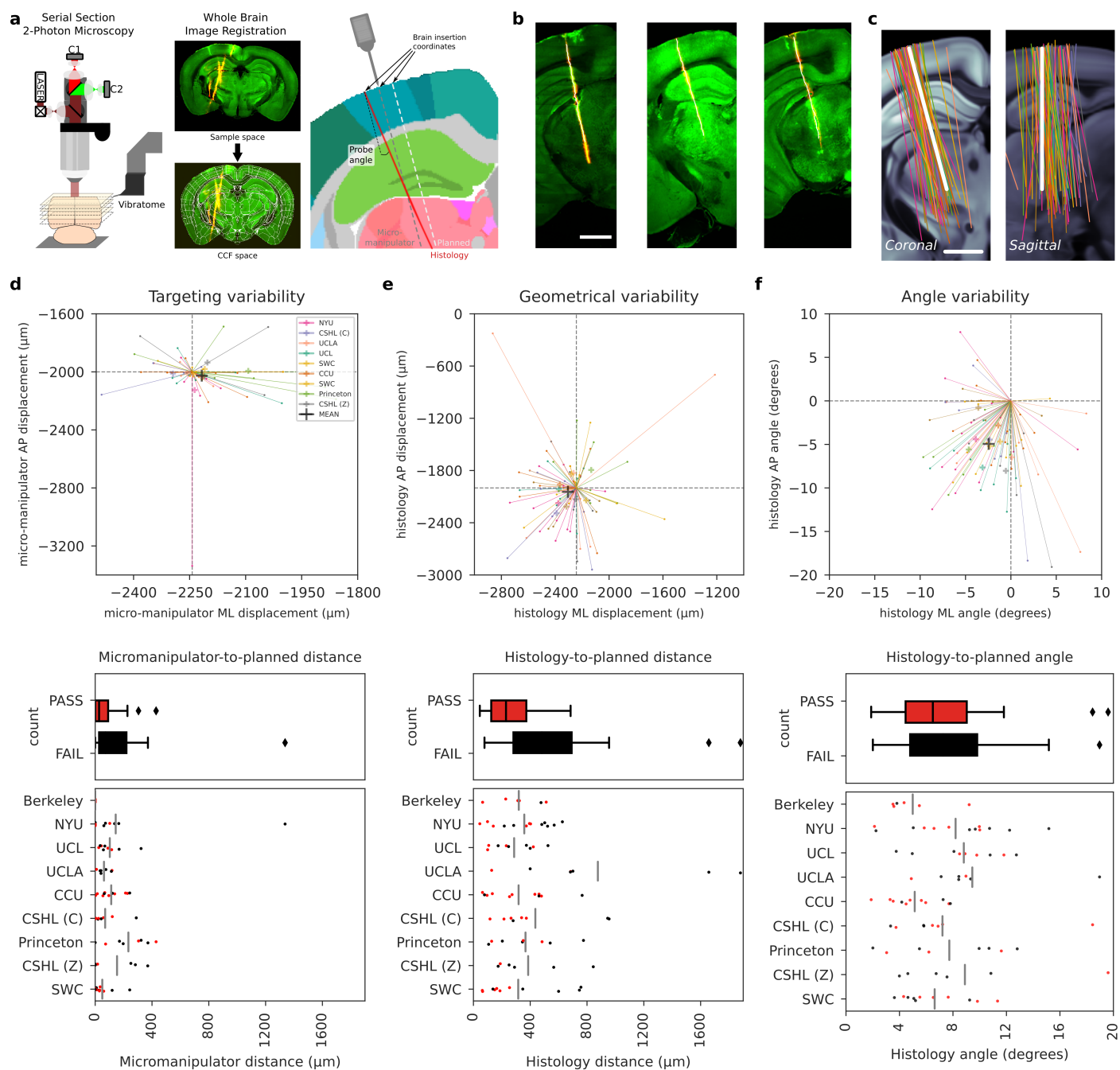


Figure 2. Probe placement shows variance that is reduced with exclusion criteria. **a**, The histology pipeline for electrode probe track reconstruction and its assessment, consisting of serial section 2-photon microscopy, and manual probe tracing. Three separate trajectories can be defined per probe: the planned trajectory; the micro-manipulator trajectory, based on the experimenter's stereotaxic coordinates; and the histology trajectory, interpolated from tracks traced in the histology data. **b**, Examples of tilted slices through the histology reconstructions showing the repeated site probe track. Plots show the green auto-fluorescence data used for CCF registration; and red cm-Dil signal, used to mark the probe track. White dots show the projections of channel positions onto each tilted slice. Scale bar: 1mm. **c**, Histology probe trajectories are interpolated from traced probe tracks and plotted as 2D projections in coronal and sagittal planes, tilted along the repeated site trajectory over the allen CCF, color coded by laboratory. Scale bar: 1mm. **d**, Targeting variability of probe placement on the brain surface: scatterplot showing the planned insertion coordinate on the brain surface in ML-AP dimensions, with the position of each subjects' insertion plotted according to the experimenter's stereotaxic coordinates of the probe, color coded by laboratory. Below, boxplots of the distances from planned to stereotaxic coordinates grouped by exclusion criteria, and dotplots by laboratory of stereotaxic-to-planned distances, colour coded by passing our exclusion criteria. **e**, Geometrical variability of probe placement on the brain surface: scatterplot of the planned insertion coordinate on the brain surface in ML-AP dimensions, with the position of each subjects' insertion plotted according to the histology-derived coordinates of the probe, color coded by laboratory. Below, boxplots of the distances from planned to histology coordinates grouped by exclusion criteria, and dotplots by laboratory of histology-to-planned distances, colour coded by passing our exclusion criteria. **f**, Angle variability of probe insertion angle: scatterplot showing the magnitude and direction of the probe angle in ML-AP dimensions, derived from histological reconstructions. Below, boxplots of the relative angles from histology to planned trajectories grouped by exclusion criteria, and dotplots by laboratory of histology-to-planned angle, colour coded by passing our exclusion criteria.

Figure 2-Figure supplement 1. Tilted slices along the histology insertion for all insertions from all labs used in assessing probe placement.

Criterion	Definition
Targeting criterion	At least 4 electrode channels in at least 3 of the 5 target brain regions
Behavior criterion	Mouse completed at least 400 trials
Yield criterion	At least 0.1 neurons (that pass single unit criteria*) per electrode channel in each region
Noise criterion	Median action-potential band RMS (AP RMS) less than 40 μ V and Median LFP power less than -140 dB
Session number criterion	For analyses that directly compared between labs (permutation tests: Fig 3d-f, Fig 4c, Fig 6), only labs with at least 3 passing sessions per brain region were included.
*Single unit metrics	Each neuron was defined as passing single unit QC if it passed three metrics: a refractory period violation metric, a noise cutoff metric, and a median amplitude threshold. Described further in (<i>The International Brain Laboratory et al., 2022a</i>).

Table 1. Quality control criteria for sessions and neurons

176 **Electrophysiological features are reproducible across laboratories**

177 In addition to the "targeting" exclusion criterion, we implemented four other exclusion criteria
178 (see Table 1). We recorded a total of 74 sessions targeted at our planned repeated site (Figure 3a).
179 Of these, 13 were excluded due to unsuccessful data acquisition that could occur from session
180 interruptions (e.g. power outage). Three recordings did not pass our targeting criterion (at least 5
181 electrode channels in at least 3 of the target brain regions). Six did not pass our behavior criterion
182 (at least 400 trials completed). Eight did not pass our criteria for low yield recordings. Finally, three
183 recordings did not pass our criterion for noise or other electrical artifacts. In subsequent figures,
184 only recordings that passed these quality control criteria were included. In analyses that directly
185 compared across labs (permutation tests; Fig 3d-f, 4c, 5d, 6), only labs which performed three
186 or more successful sessions were included. Furthermore, single units had to pass three quality
187 control metrics to be included in single unit analyses (*The International Brain Laboratory et al.,*
188 *2022a*). When plotting all recordings, including those that failed to meet quality control criteria,
189 one can observe that discarded sessions were often clear outliers (Figure 3b-c, supp. 1). Overall,
190 we analyzed data recorded from the 40 remaining sessions recorded in 9 labs to determine the
191 reproducibility of our electrophysiological recordings.

192

193 We set out to answer the question whether electrophysiological features, such as firing rates
194 and LFP power, were reproducible across laboratories. In other words, is there consistent varia-
195 tion across laboratories in these features that is larger than expected by chance? We first visualized
196 LFP power, a feature used by experimenters to guide the alignment of the probe position to brain
197 regions, for all the repeated site recordings (Figure 3b). The dentate gyrus (DG) is characterized
198 by high power spectral density of the LFP (*Penttonen et al., 1997; Bragin et al., 1995; Senzai and*
199 *Buzsáki, 2017*) and this feature was used to guide physiology-to-histology alignment of probe po-
200 sitions (Figure 3 supplementary 2). By plotting the LFP power of all recordings along the length of
201 the probe side-by-side, aligned to the boundary between the DG and thalamus, we confirmed that
202 this band of elevated LFP power was clearly visible in all recordings at the same depth. The probe
203 alignment allowed us to attribute the channels of each probe to their corresponding brain regions
204 to investigate the reproducibility of electrophysiological features for each of the target regions of
205 the repeated site. To visualize all the neuronal data, each neuron was plotted at the depth it was
206 recorded overlaid with the position of the target brain region locations (Figure 3b).

207 The reproducibility of electrophysiological features over laboratories was investigated using
208 permutation testing. The tested features included neuronal yield, firing rate, spike amplitude, LFP
209 power, and action-potential band RMS (AP RMS). For each feature and each brain region, the within-
210 lab and across-lab means were calculated (example in Figure 3c). If the electrophysiological feature
211 is reproducible across laboratories, there should be a small deviation between the mean over an-
212 imals within a lab and the mean over all the lab means. To investigate whether the deviation was
213 significantly larger than expected by chance, we performed permutation testing in which the lab
214 labels were shuffled and a p-value was calculated by comparing the actual deviation from the shuf-
215 fled null-distribution. Because a test is performed per region-metric pair, the p-values were cor-
216 rected for multiple testing using the Benjamini-Hochberg procedure (*Seabold and Perktold, 2010;*
217 *Benjamini and Hochberg, 1995*). We found that all electrophysiological features were reproducible
218 across laboratories for all regions studied.

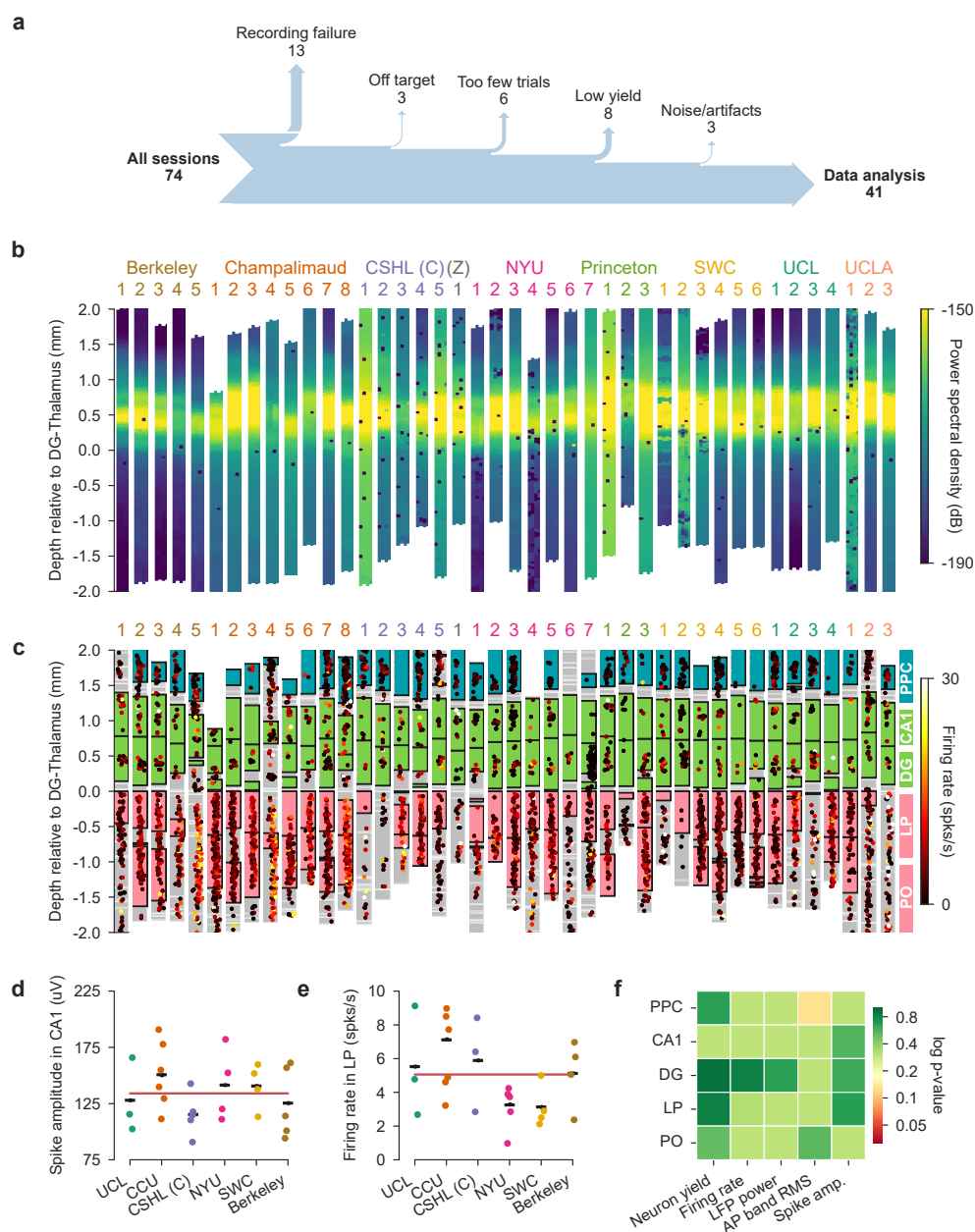


Figure 3. Electrophysiological features are reproducible across laboratories. **a**, Number of experimental sessions recorded; number of sessions used in analysis due to exclusion criteria. **b**, Power spectral density between 20 and 80 Hz of each channel of each probe insertion (vertical columns) shows reproducible alignment of electrophysiological features to histology. Insertions are aligned to the boundary between the dentate gyrus and the thalamus. CSHL: Cold Spring Harbor Laboratory [(C): Churchland lab, (Z): Zador lab], NYU: New York University, SWC: Sainsbury Wellcome Centre, UCL: University College London, UCLA: University of California, Los Angeles. **c**, Firing rates of individual neurons according to the depth at which they were recorded. Colored blocks indicate the target brain regions of the repeated site; if no block is plotted the neurons are in a region that is not one of the targets. Dots are neurons, colors indicate firing rate, displacement along the x-axis indicates spike amplitude. **d,e**, Examples of permutation testing to determine whether the deviation of lab means (short black lines) from the mean across labs (red line) was larger than expected by chance. For each region, only laboratories that had three or more recordings in that region were included in the permutation testing. Here the median spike amplitude in CA1 and median firing rate in LP is plotted per lab. A p-value was determined by shuffling the lab labels 10,000 times. **f**, P-values for five electrophysiological metrics, computed separately for all target regions. P-values are plotted on a log-scale to visually emphasize values close to significance.

Figure 3-Figure supplement 1. Electrophysiological features of *all* recordings, including recordings that failed quality criteria.

Figure 3-Figure supplement 2. High LFP power in dentate gyrus was used to align probe locations in the brain.

219 **Task-driven activity of brain regions is reproducible across laboratories**

220 Concerns about reproducibility include not only basic electrophysiological properties, but also
221 modulation of firing rates by task variables. To address this, we analysed the reproducibility of
222 the relationship between neural activity and task variables across laboratories. In particular, we
223 were interested in whether the brain regions targeted here have comparable neural responses
224 to task events, such as stimulus onset, movement onset, and reward delivery. An inspection of
225 individual neurons revealed clear modulation by, for instance, the onset of movement (Fig. 4a).
226 When considering all neurons within a single region of a given session however, it becomes clear
227 that, while a number of neurons are modulated, there is also a proportion of neurons that do
228 not change their firing in relation to task events (Fig. 4b) (Urai *et al.*, 2022). Plotting the session-
229 averaged response for each experiment in a given area reveals that despite variability, many key
230 features are reproduced, such as the general response time course and timing (Fig. 4c; also Fig.
231 6d).

232 Having observed that many individual neurons are modulated during the task, we then wanted
233 to compare how the proportion of modulated neurons differed across labs. This is especially im-
234 portant, as we are often interested in determining which regions are involved in the neural compu-
235 tations underlying task performance. Therefore, within each brain region, we compared the pro-
236 portion of the neural population that was sensitive to specific elements of the task. Using Wilcoxon
237 sign-rank tests and Wilcoxon rank-sum tests (Steinmetz *et al.*, 2019), we used seven tests to iden-
238 tify neurons with significantly modulated firing rates during specific time-periods of the task. The
239 general logic of these tests is displayed in Fig. 5b and Fig. 5-supplemental 1. The neurons that
240 were found by these tests showed a clear modulation to the tested events, as expected (Fig. 5a-b).
241 For most tests, the proportions of modulated neurons across sessions and across brain regions
242 were quite variable (Fig. 5c and Fig. 5-supplemental 1). However, when applying a permutation
243 test as used in our previous analyses, we found no significant differences across labs regarding
244 the proportion of task-modulated units (Fig. 5d). We can therefore conclude that task-modulated
245 activity is reproducible across labs.

246 To further investigate neuronal task-modulation, we also measured the Fano Factor of single
247 units. The Fano Factor is a useful measurement of firing rate variability and is defined as the spike
248 count variance over trials divided by spike count mean. The Fano Factor enables the comparison
249 of the fidelity of signals across neurons and regions, despite differences in firing rates (Tolhurst
250 *et al.*, 1983). Further, the temporal dynamics of the Fano Factor can be informative about under-
251 lying neural computations (Churchland *et al.*, 2010, 2011). We calculated the Fano Factor using a
252 sliding window over each trial. In most brain regions, the Fano Factor, averaged over all neurons,
253 decreased around the time of movement onset (Fig. 7-supplemental 4, left column). Based on the
254 Fano Factor time course, we selected the period between 40-200 ms after movement onset (for cor-
255 rect trials with full-contrast stimuli on the right side) to calculate an average Fano Factor per neuron
256 and quantify differences in Fano Factor across labs. While Fano Factor values varied between neu-
257 rons and across sessions, we found no difference across labs after applying a permutation test
258 (Fig. 5d). This argues that the decrease in neural variability around the time of movements is re-
259 producible and is present not only in cortical structures, as previously reported (Churchland *et al.*,
260 2010), but is also reliably present in the hippocampus and thalamus.

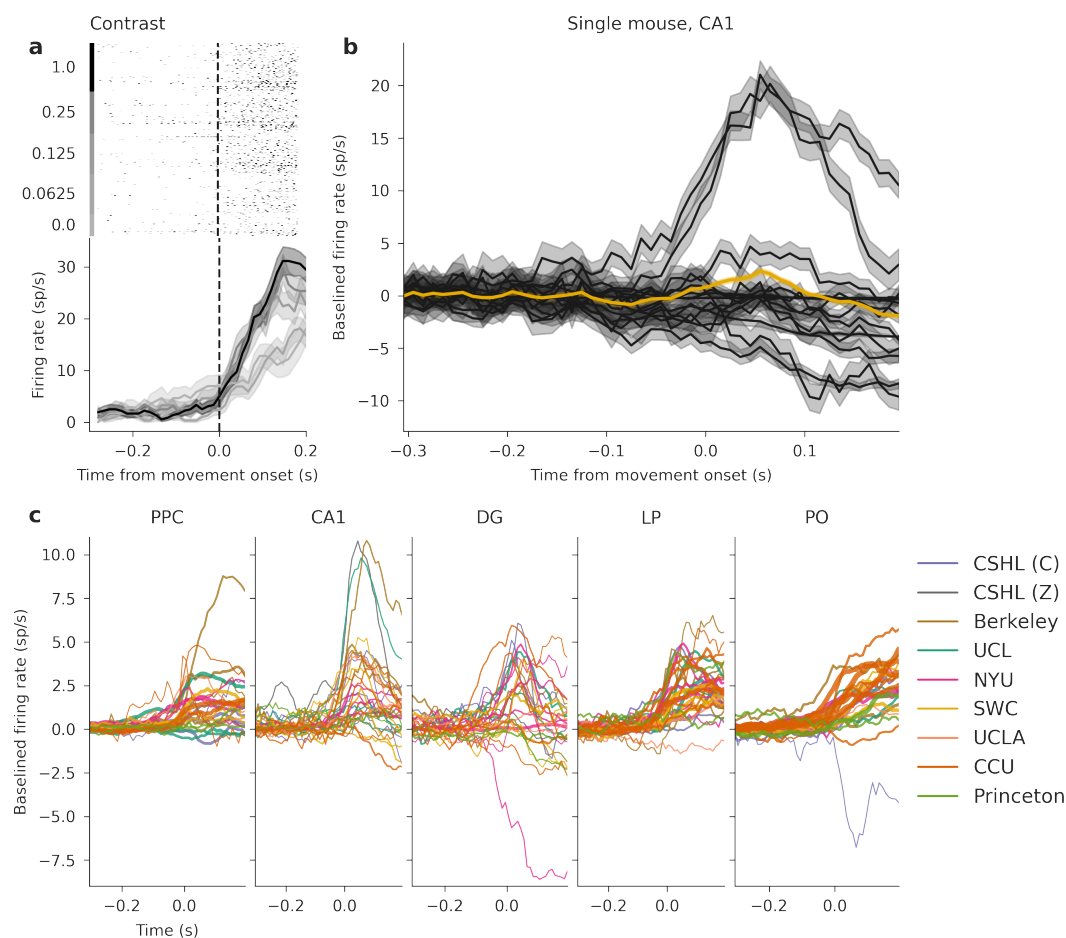


Figure 4. Neural activity is modulated during decision-making in 5 neural structures. **(a)** Neural activity in relation to movement onset towards the left for different contrasts, raster plot (top), peristimulus time histogram (bottom). **(b)** Peri-event time histograms (PETs) for correct left choices of all neurons from CA1 of a single mouse, aligned to movement onset. These PETs are baseline-subtracted by a pre-stimulus baseline. Shaded areas show standard error of mean (and propagated error for the overall mean). The thicker line shows the average over this entire population, coloured by the lab from which the recording originates. **(c)** Average PETs for correct left choices across regions within individual mice (same as thick line in (b)). Line thickness indicates how many neurons went into the average (min=4, max=86). (As we do not compare across labs, we do not subset to labs with sufficient recordings here).

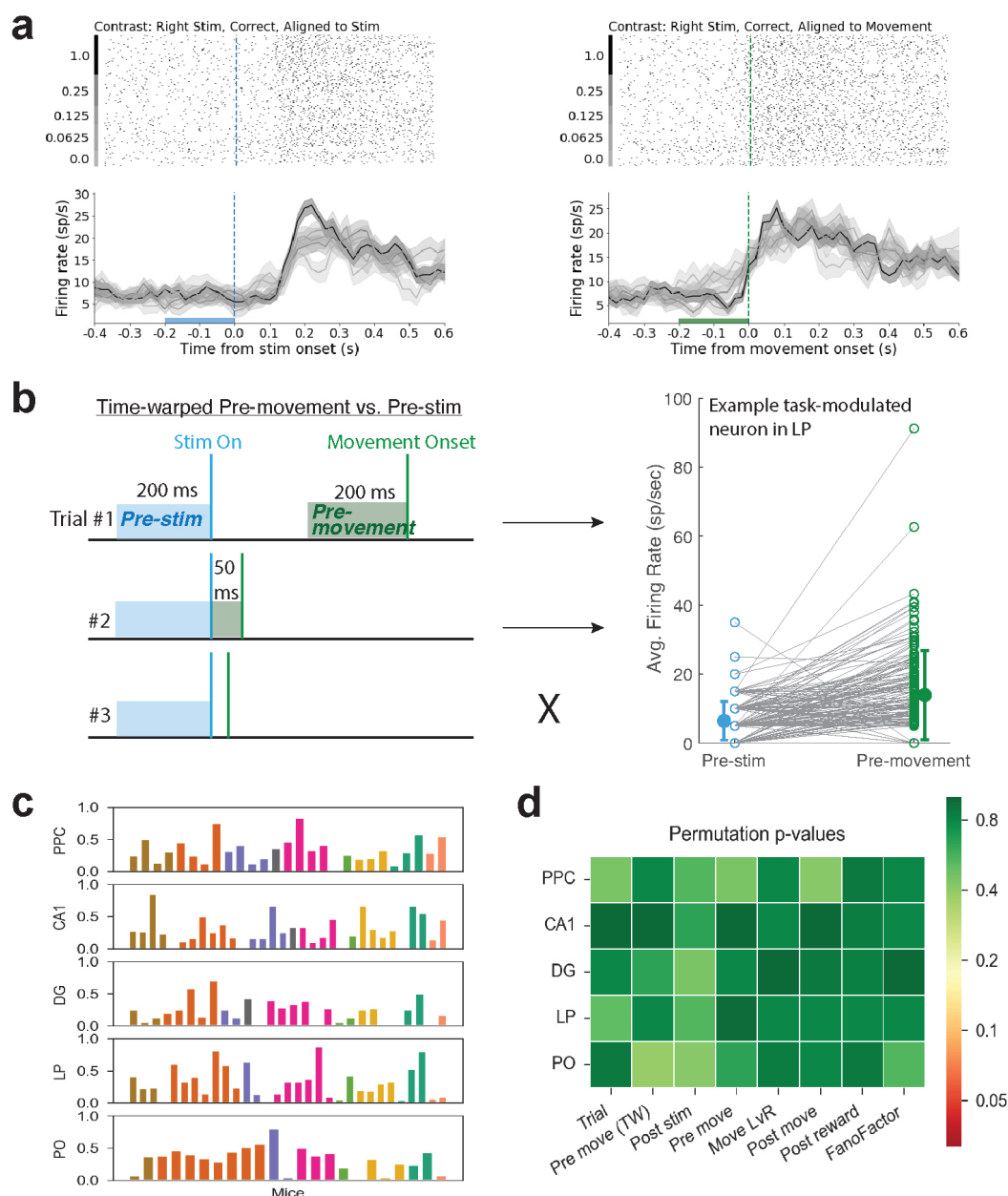


Figure 5. Task-modulated neurons are not significantly different between laboratories. (a) Raster plots and firing rate time courses of an example neuron in LP, aligned to either stimulus onset or movement onset; plotted only for right visual stimuli and correct movements. (The firing rates are calculated using a sliding window and are causal, such that each time point includes a 40 ms window prior to the indicated point.) (b) Schematic of the time-warped (TW) pre-movement vs. pre-stimulus test for finding task-modulated neurons (*left*), where the firing rate prior to movement onset is compared against the firing rate during 200 ms before the stimulus. This is only calculated for trials where the time between pre-movement time and stimulus is at least 50 ms (third example trial is excluded). Also, the pre-movement time is considered only up to 200 ms prior to the movement onset, i.e., the pre-movement period can range anywhere from 50 ms to 200 ms prior to the onset of the stimulus (resulting in continuous firing rates in the right panel), unlike the pre-stimulus period which is always set to 200 ms (thus, firing rates in the right panel change in increments of 5 sp/sec). (*right*) The change in firing rate of the example neuron in a, which is considered a task-modulated neuron using the TW pre-movement test; each gray line indicates one trial. Mean pre-stimulus and pre-movement firing rates across all trials are shown with filled circles (error bar: standard deviation). (c) Proportion of task-modulated neurons for each mouse in each of the five brain regions using the TW pre-movement test. Each column or color indicate, in order, a different recording session or lab. (Note that there is no correspondence here between columns across different brain regions.) (d) Permutation test results comparing across-lab variation in the proportion of task-modulated neurons found using each of the seven tests examined (the TW pre-movement test in b and c and six other tests described in Fig 5-Figure Supplement 1), as well as variation in the neuronal Fano Factors. All task-modulated comparisons were performed for correct trials with non-zero contrast stimuli.

Figure 5-Figure supplement 1. Proportion of task-modulated neurons, defined by six additional tests, across mice, labs, and brain regions.

261 **Principal component embedding analysis reveals little functional separation be-** 262 **tween labs**

263 In the previous section, we tested specific hypotheses about modulations in task-driven activity at
264 different times within the behavioral trial. We wondered if our conclusions about reproducibility
265 would remain consistent if we perform comparisons across labs and brain regions at the level of
266 the trial-averaged firing rate vectors computed over the entire trial.

267 The first step is to choose a summary of each cell's neural activity that can be directly compared
268 across experimental sessions and labs. The peri-event time histogram (PETH) is one such summary
269 that is commonly used. The PETH depends on the event used to align trials, and also discards in-
270 formation about behavioral variability across trials. To retain more of this information, we coarsely
271 split trials into two sets, one with fast reaction times (< 0.15 s) and one with slower reaction times
272 (> 0.15 s). Then we computed PETHs within each of these subsets and concatenated the resulting
273 vectors to obtain a more informative summary of each cell's average activity within these different
274 types of behavioral trials. (The results described below did not depend strongly on the details of
275 the trial-splitting we chose; for example, splitting trials by "left" vs "right" behavioral choice led to
276 similar results.) See Figure 6a for two example cells' PETHs, showing only the PETH obtained by
277 averaging fast reaction time trials.

278 Next, we project these high-dimensional summary vectors into a low-dimensional "embedding"
279 space that captures the variability of the neuronal population but at the same time allows for
280 easy visualization and further analysis. We found that a simple principal component analysis (PCA)
281 provided a useful embedding. Specifically, we stack each cell's summary double-PETH vector (de-
282 scribed above) into a matrix (containing the summary vectors for all cells across all sessions) and
283 run PCA to obtain a low-rank approximation of this matrix (see Methods). Figure 6a shows two
284 cells and the corresponding two-dimensional PCA approximation, with one high-accuracy recon-
285 struction example and one low-accuracy example shown here. Figure 6b displays the goodness of
286 this PCA approximation over the full population as a function of the number of PCs used, showing
287 that the PETHs of the majority of cells can be well reconstructed even with just 2 PCs.

288 Now we have obtained a simple two-dimensional summary of each cell's activity that we can
289 visualize easily; see Figure 6c. This simple embedding is already sufficiently powerful to distinguish
290 different brain regions: in Figure 6c we have colored cells by region, and we see that e.g. regions
291 PO and CA1 show displaced clusters, illustrating clear regional differences in cell activities. These
292 per-region differences are also visible in the region-averaged PETHs (Figure 6d). We quantified this
293 separation via a permutation test, computing the sum across each region's distance between its
294 mean embedded activity and the mean across all regions and comparing that to the null distribu-
295 tion of values obtained in the same way after shuffling the region labels. The p-value is < 0.0001 ,
296 indicating a significant difference between regional PCA-reduced PETHs.

297 To test for activity differences between labs, we subdivided the embedded point clouds (Figure
298 6c) by lab (Figure 6e and supp. Figure 1). The standard deviation of these activity point clouds
299 show large overlap across most labs, indicating similar activity. For each region separately, we
300 determined whether the sum across each lab's distance between its mean embedded activity and
301 the mean across all labs is significantly different, using the same permutation test as described in
302 the previous paragraph, this time shuffling lab labels. We obtain one false discovery rate corrected
303 p-value for this lab-permutation test per region - PO 0.706, LP 0.065, DG 0.706, CA1 0.168, PPC
304 $p < 0.0001$ - finding that for all regions except PPC the sum of mean lab embedded activities is
305 not significantly different than the mean over all labs. We thus see that embedded activity differs
306 clearly across regions but much less so across labs.

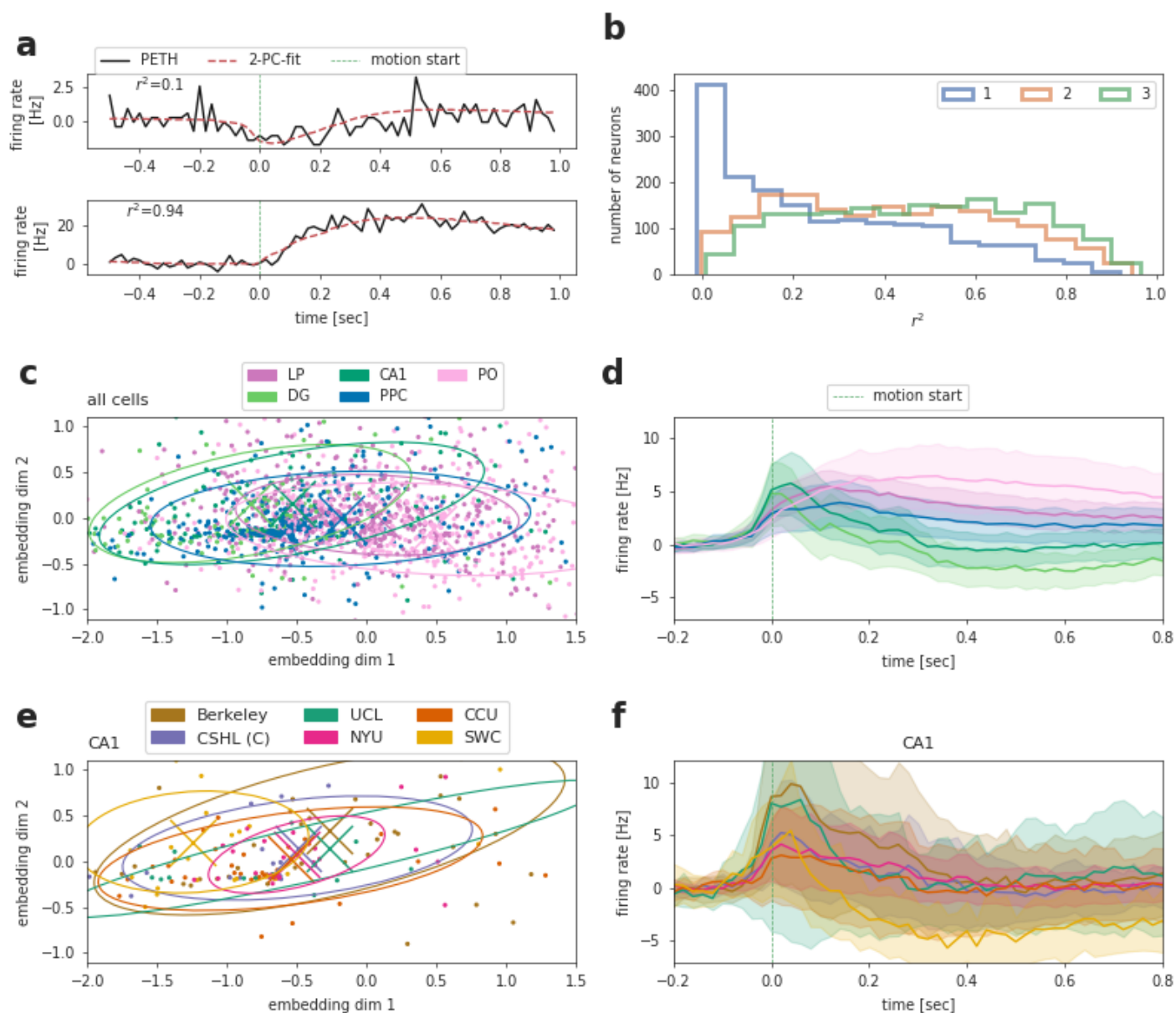


Figure 6. Principal component embedding of peri-event time histograms separates cells from different brain regions but not cells from different labs. (a) Two example cells' PETHs in black and 2-PC-based reconstruction; poor (top), good (bottom) fit with goodness of fit r^2 indicated on top. (b) Histograms of reconstruction goodness of fit across all cells based on reconstruction by 1-3 PCs. With only the first 2 PCs most PETHs are well approximated, justifying the subsequent two-dimensional embedding analysis. (c) Two-dimensional embedding of PETHs of all cells, colored by region (each dot corresponds to a single cell). X's and ellipses indicate the mean and standard deviation for each region. (d) Mean firing rates of all cells in each of the studied regions. As in the 2D embedding, mean values for PO and CA1 clearly separate. (Error bars are standard deviation across cells divided by square root of number or recordings per region). (e) Embedded activity of CA1 neurons plotted separately for each lab (colors). (f) Mean activity for all labs in CA1 (color conventions the same as in (e)). See supp. Figure 1 for the other regions. (Error bars are standard deviation across cells divided by square root of number of recordings per lab). Note that only 6 labs are included in this analysis, as we only include labs that have at least 3 recordings per region (see exclusion criterion Table 1).

Figure 6-Figure supplement 1. Regional 2-PC embedding and average PETH per lab

307 **Differences in neuronal spatial position and spike characteristics are a minor source**
308 **of variability across sessions**

309 While we found little variability between laboratories in terms of electrophysiological features and
310 task variables, we observed large variability between recording sessions and mice (Fig. 3, Fig. 5,
311 and Fig. 5-supplemental 1). Since the spatial position of the Neuropixels probe was variable be-
312 tween sessions (Fig. 2), we examined variability in targeting as a potential source of differences in
313 neuronal activity for each of the five repeated site brain regions. We also considered single-unit
314 spike waveform characteristics as a source of variability. In the next section, we examine other
315 potential sources of variability (e.g., mouse movements).

316 To investigate variability in session-averaged firing rates, we identified neurons which had fir-
317 ing rates different from the majority of neurons within each brain region (absolute deviation from
318 the median firing rate being >15% of the firing rate range). These outlier neurons, which mostly
319 turned out to be high-firing (except in PO), were compared against regular neurons in terms of five
320 features: spatial position (x, y, z, computed as the center-of-mass of each unit's spike template on
321 the probe, localized to CCF coordinates in the histology pipeline) and spike waveform character-
322 istics (amplitude, peak-to-trough duration). We observed that recordings in all areas, such as LP
323 (Fig. 7a), indeed spanned a wide space within that area. Interestingly, in areas other than DG, the
324 highest firing neurons were not entirely uniformly distributed in space. For instance, in LP, high
325 firing neurons tend to be positioned more laterally and centered on the anterior-posterior axis
326 (Fig. 7b). In PPC and PO, the spatial position of neurons, but not differences in spike character-
327 istics, contributes to differences in session-averaged firing rates (Fig. 7-supplemental 1b and 3c).
328 In contrast, high-firing LP, CA1, and DG neurons have different spike characteristics compared to
329 other neurons in their respective regions (7b and Fig. 7-supplemental 2b and 3a).

330 To quantify the amount of variability in average firing rates that can be explained by spatial
331 position or spike characteristics, we fit a linear regression model with these five features (x, y, z,
332 spike amplitude, and duration) as the inputs. We found similar results: In PPC, z position, or neuron
333 depth, explained part of the variance (had a significant weight); in CA1 and DG, spike amplitude, not
334 spatial position, explained part of the variance; in LP, x and y positions as well as spike amplitude
335 explained some of the variance; in PO, x and y position captured more variance than the other
336 features. In LP, where the most amount of variability can be explained by this regression model,
337 these features account for a total of ~12% of the firing variability. In PPC, CA1, DG, and PO, they
338 account for approximately 3%, 6%, 6%, and 5% of the variability, respectively.

339 Next, we examined whether neuronal spatial position and spike features contributed to vari-
340 ability in task-modulated activity. We found that all brain regions, except CA1, had minor, yet sig-
341 nificant, differences in spatial positions of task-modulated and non-modulated neurons (using the
342 definition of at least of one of the seven tests in Fig. 5d). For instance, task-modulated LP neu-
343 rons defined by the time-warped pre-movement test, were positioned more ventrally and centered
344 along the anterior-posterior axis (Fig. 7c), while task-modulated LP neurons defined by the left ver-
345 sus right pre-movement test, tended to be more ventral (Fig. 7d). Other brain regions had less
346 spatial differences than LP (Fig. 7- supplemental 1, 2, 3). Spike characteristics were significantly
347 different between task-modulated and non-modulated neurons only for some tests and only in
348 PPC, DG, and PO (Fig. 6-supplemental 1c-d and 3)b-d. On the other hand, the task-aligned Fano
349 Factor of neurons did not have any differences in spatial position except for in PPC, where lower
350 Fano Factors (<1) tended to be located ventrally (Fig. 7- supplemental 4a). Spike characteristics of
351 neurons with lower vs. higher Fano Factors were only different in the LP and PO (Fig. 7- supple-
352 mental 4). Lastly, we trained a linear regression model to predict the 2D embedding of PETHs of
353 each cell shown in Fig 6c from the x, y, z coordinates and found that spatial position contains little
354 information ($r^2 \sim 4\%$) about the embedded PETHs of cells.

355 In summary, our results suggest that spatial position is a small contributor to variability for
356 session-averaged firing rates in all brain regions except DG, and to a lesser degree for task-modulated

357 neuronal activity in all brain regions except CA1. In all regions, spike characteristics also have a mi-
358 nor contribution to the observed variability. Since, overall, the contributions of spatial position and
359 spike features were small, despite being significant, we examine other sources of variability in the
360 next section.

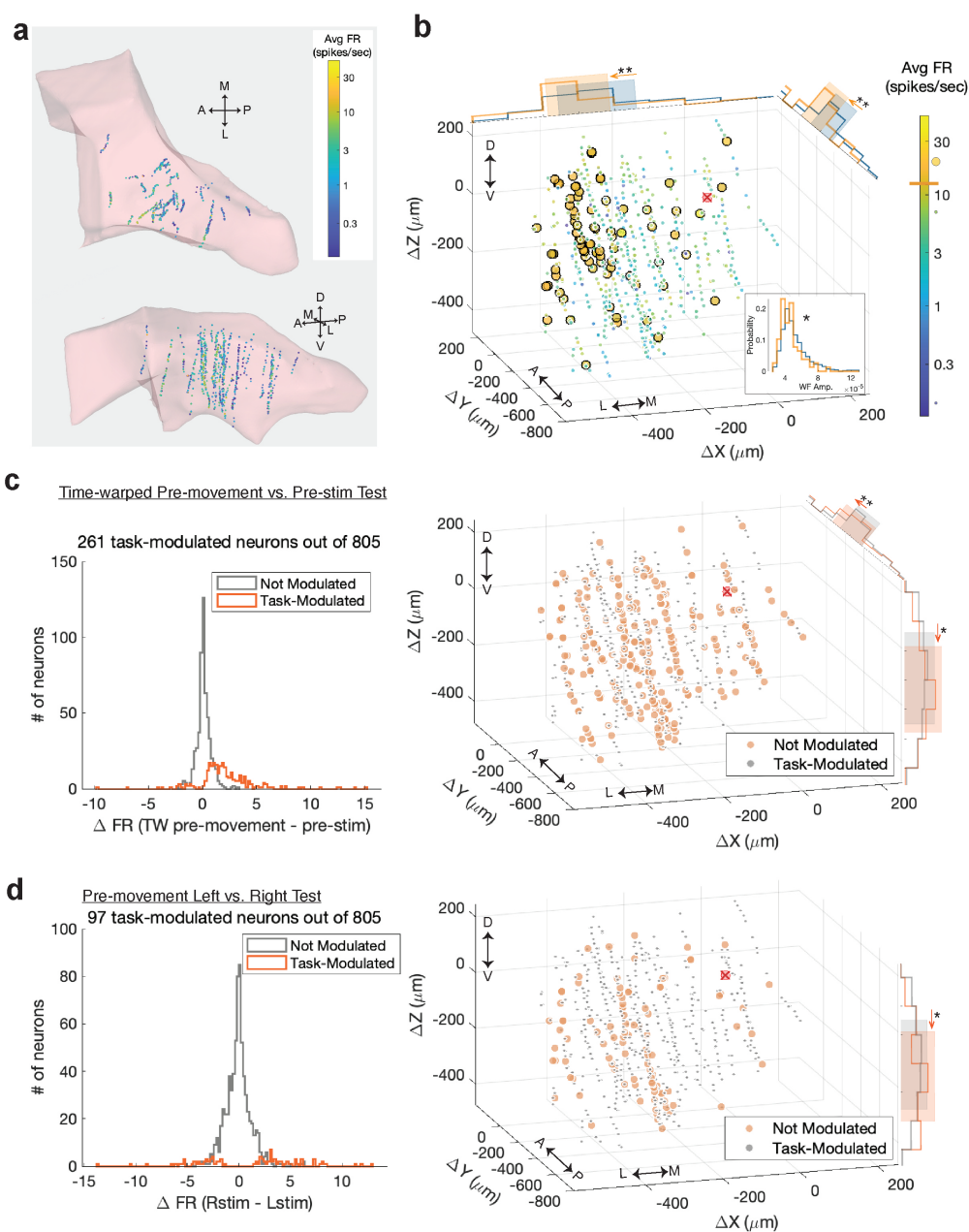


Figure 7. High-firing and task-modulated LP neurons have slightly different spatial positions than other LP neurons, potentially contributing to variability between sessions. (a) Spatial positions of recorded neurons in LP, color-coded with their firing rates averaged over the recording session. (b) Spatial positions of LP neurons plotted as distance from the planned target center of mass, indicated with the red x. (To enable visualization of overlapping data points, jitter was added to the unit locations.) Larger circles indicate outlier neurons (defined by a normalized firing rate deviation > 15%, resulting in a threshold of 12 sp/sec for LP, shown on the colorbar; here, 78 out of 805 neurons were outliers). Only histograms of the spatial positions and waveform features that were significantly different between the outlier (yellow) and regular (blue) units are shown (two-sample Kolmogorov-Smirnov test with Bonferroni correction for multiple comparisons; * and ** indicate corrected p-values of <0.05 and <0.01, in order). Shaded areas indicate the area between 20th and 80th percentiles of the neurons' locations. (c) (Left) Histogram of firing rate changes during the pre-movement period from the pre-stimulus period (using the time-warped test, Fig. 5b-c) for task-modulated (orange) and non-modulated (gray) neurons. (Right) Spatial positions of task-modulated and non-modulated LP neurons, with histograms of significant features (here, y and z positions) shown. (d) Same as c but using the pre-movement left vs. right test to identify task-modulated units.

Figure 7-Figure supplement 1. High-firing and task-modulated PPC neurons.

Figure 7-Figure supplement 2. High-firing and task-modulated CA1 neurons.

Figure 7-Figure supplement 3. High-firing and task-modulated DG and PO neurons.

Figure 7-Figure supplement 4. Time-course and spatial position of neuronal Fano Factors.

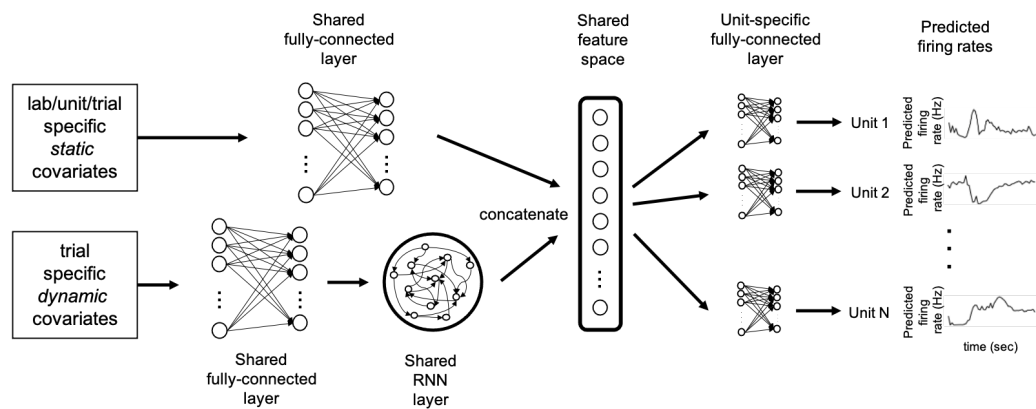


Figure 8. Schematic of the multi-task neural network model architecture: We adapt a multi-task neural network approach for unit-specific firing rate prediction. The model takes in a set of covariates, and outputs time-varying firing rates for each neuron for each trial. The covariates include the lab ID, 3-D unit location, and trial event times (e.g., stimulus onset); see Table 2 for a full list. The initial embedding layer of the network is shared across all units, and serves to learn a useful (nonlinear) shared set of features that all the individual units can regress onto for their predictions.

361 **A multi-task neural network accurately predicts activity and quantifies sources of** 362 **neural variability**

363 As discussed above, variability in neural activity between labs or between sessions can be due to
364 many factors. These include differences in behavior between animals, differences in probe place-
365 ment between sessions, and uncontrolled differences in experimental setups between labs. How
366 can we quantify and distinguish these different sources of variability? Simple linear regression
367 models or generalized linear models (GLMs) are likely too inflexible to capture the nonlinear con-
368 tributions that many of these variables, including lab IDs and spatial positions of neurons, might
369 make to neural activity. On the other hand, fitting a different nonlinear regression model (involv-
370 ing many covariates) individually to each recorded unit would be computationally expensive and
371 could lead to poor predictive performance due to overfitting.

372 To estimate a flexible nonlinear model given constraints on available data and computation
373 time, we adapt an approach that has proven useful in the context of sensory neuroscience (*Mcln-*
374 *tosh et al., 2016; Batty et al., 2016; Cadena et al., 2019*). We use a "multi-task" neural network
375 (MTNN; Figure 8) that takes as input a set of covariates (including the lab ID, the unit's 3D spatial
376 position in standardized CCF coordinates, the animal's estimated pose extracted from behavioral
377 video monitoring, feedback times, and others; see Table 2 for a full list). The model learns a shared
378 set of nonlinear features (shared over all recorded units) and fits a Poisson regression model on
379 this shared feature space for each unit. (With this approach we effectively solve multiple nonlin-
380 ear regression tasks simultaneously; hence the "multi-task" nomenclature.) The model extends
381 simpler regression approaches by allowing nonlinear interactions between variables. In particular,
382 previous reduced-rank regression approaches (*Kobak et al., 2016; Izenman, 1975*) can be seen as a
383 special case of the multi-task neural network, with a single hidden layer and linear weights in each
384 layer.

385 Figure 9a shows model predictions on held-out trials for a single CA1 unit. We plot the observed
386 and predicted peri-event time histograms and raster plots, split into left vs. right trials. As a visual
387 overview of which behavioral covariates are highly correlated with this cell's activity on each trial,
388 various behavioral covariates that are input into the MTNN are shown in Figure 9b. Overall, the
389 MTNN approach accurately predicts the observed firing rates. When the MTNN and GLMs are
390 trained on a reduced set of covariates, consisting of stimulus onset timing, stimulus contrast and
391 side, feedback type and timing, first movement onset timing, wheel velocity, and mouse's prior, the

Variable Name	Type	Group	Note
Lab ID	Categorical / Static		
Session ID	Categorical / Static		
Unit 3D spatial position	Real / Static	Electrophysiological	In standardized CCF coordinates
Unit amplitude	Real / Static	Electrophysiological	Template amplitude
Unit waveform width	Real / Static	Electrophysiological	Template width
Paw speed	Real / Dynamic	Behavioral	Inferred from DLC
Nose speed	Real / Dynamic	Behavioral	Inferred from DLC
Pupil diameter	Real / Dynamic	Behavioral	Inferred from DLC
Motion energy	Real / Dynamic	Behavioral	
Stimulus	Real / Dynamic	Task-related	Stimulus side, contrast and onset timing
Go cue	Binary / Dynamic	Task-related	
First movement	Binary / Dynamic	Task-related	
Choice	Binary / Dynamic	Task-related	
Feedback	Binary / Dynamic	Task-related	
Wheel velocity	Real / Dynamic	Behavioral	
Mouse Prior	Real / Static		Mouse's prior belief
Last Mouse Prior	Real / Static		Mouse's prior belief in previous trial
Lick	Binary / Dynamic	Behavioral	
Decision Strategy	Real / Static		Decision-making strategy (<i>Ashwood et al., 2022</i>)
Brain region	Categorical / Static	Electrophysiological	5 repeated site regions

Table 2. List of covariates input to the multi-task neural network. See Appendix for additional details.

392 MTNN and GLMs perform similarly on predicting the firing rates of held-out test trials. Furthermore,
 393 the MTNN trained on the full set of covariates in Table 2 outperforms the MTNN and GLMs trained
 394 on the reduced covariate set (See Figure 9 supplemental 2).

395 Next we use the predictive model performance to quantify the contribution of each covariate
 396 to the fraction of variance explained by the model. Following *Musall et al. (2019)*, we run two com-
 397plementary analyses to quantify these effect sizes: *single-covariate fits*, in which we fit the model
 398 using just one of the covariates, and *leave-one-out fits*, in which we train the model with one of the
 399 covariates left out and compare the predictive explained to that of the full model. As an exten-
 400sion of the leave-one-out analysis, we run the *leave-group-out analysis*, in which we quantify the
 401 contribution of each group of covariates (electrophysiological, task-related, and behavioral) to the
 402 model performance. Using data simulated from GLMs, we first validate that the MTNN leave-one-
 403 out analysis is able to partition and explain different sources of neural variability (See Figure 10
 404 supplemental 1).

405 We then run single-covariate, leave-one-out, and leave-group-out analyses to quantify the con-
 406 tributions of the covariates listed in Table 2 to the predictive performance of the model on held-out
 407 test trials. The results are summarized in Figure 10. According to the single-covariate analysis (Fig-
 408 ure 10a), face motion energy (derived from behavioral video), wheel velocity, and some task vari-
 409 ables (e.g., stimulus information and first movement onset timing) can individually explain about
 410 5-10% of variance of the units on average. The leave-one-out analysis (Figure 10b left) shows that
 411 most covariates have low unique contribution to the predictive power. This is because many vari-
 412 ables are correlated and are capable of capturing variance in the neural activity even if one of the
 413 covariates is dropped (See behavioral raster plots in Figure 9b). According to the leave-group-out

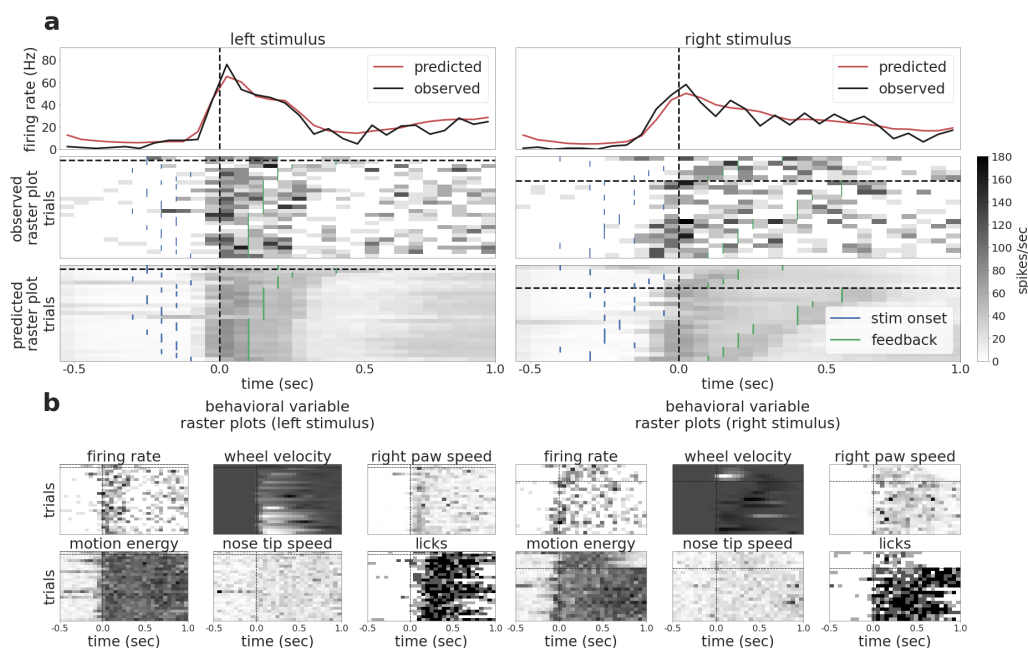


Figure 9. The MTNN model accurately estimates firing rates on held-out test trials from a CA1 neuron: (a) MTNN model estimates of firing rates (50 ms bin size) of a CA1 neuron from an example subject during held-out test trials. The trials are split into those that had stimulus on the left/right and are aligned to the first movement onset time (vertical dashed lines). We plot the observed and predicted peri-event time histograms (1st row) and the observed and predicted raster plots (2nd and 3rd rows). The blue ticks in the raster plots indicate stimulus onset, and the green ticks indicate feedback times. The black horizontal dashed line separates the incorrect/correct trials (i.e., the trials above the dashed line are incorrect trials), and the trials are ordered by reaction time. The trained model does well in predicting the (normalized) firing rates. The MTNN prediction quality measured in R^2 is 0.32 on held-out test trials and 0.90 on PETHs of held-out test trials. (b) We plot the raster plots of behavioral variables (wheel velocity, paw speed, motion energy, nose speed, and licks), ordering the trials in the same manner as in (a). We see that the MTNN firing rate predictions are modulated synchronously with the behavioral variables.

Figure 9-Figure supplement 1. Scatter plot of MTNN prediction quality (R^2) vs. mean firing rate (spikes/sec)

Figure 9-Figure supplement 2. MTNN slightly outperforms GLMs on predicting the firing rates of held-out test trials.

Figure 9-Figure supplement 3. PETHs and MTNN predictions for held-out test trials

414 analysis, the behavioral covariates as a group have the highest unique contribution to the model's
 415 performance while the task-related and electrophysiological variables have close-to-zero unique
 416 contribution. Most importantly, the leave-one-out analysis shows that lab and session IDs, condi-
 417 tioning on the covariates listed in Table 2, have close to zero effect sizes, indicating that within-lab
 418 and between-lab random effects are small and comparable.

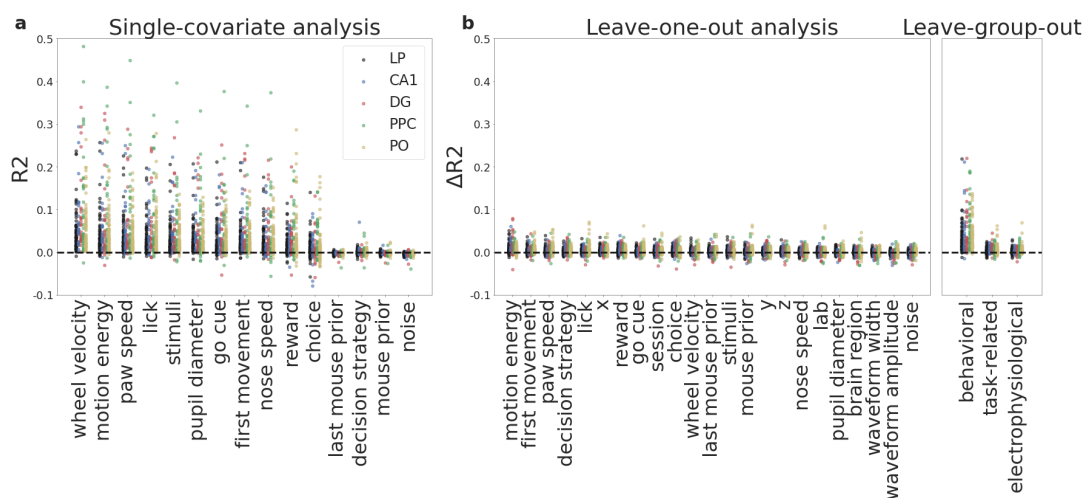


Figure 10. Single-covariate, leave-one-out, and leave-group-out analyses show the contribution of each (group of) covariate(s) to the model. Lab and session IDs have low contributions to the model. (a) Single-covariate analysis, colored by the brain region. Each dot corresponds to a single cell in each plot. **(b)** Leave-one-out analysis, colored by the brain region. The analyses are run on 246 responsive units across 20 sessions. The leave-one-out analysis shows the unique contribution of each covariate to the model, and the single-covariate analysis shows the upper limit of the contribution of each covariate to the model. The leave-group-out analysis shows how groups of electrophysiological, task-related, and behavioral covariates contribute to the model. The leave-one-out analysis shows that lab/session IDs have low effect sizes on average, indicating that within-lab and between-lab random effects are small and comparable. The “noise” covariate is a dynamic covariate (white noise randomly sampled from a Gaussian distribution) and is included as a negative control: the model correctly assigns zero effect size to this covariate. Covariates that are constant across trials (i.e., lab and session IDs, unit’s 3D spatial location) are left out from the single-covariate analysis.

Figure 10–Figure supplement 1. MTNN prediction quality on the data simulated from GLMs is comparable to the GLMs’ prediction quality. The effect sizes computed by the MTNN leave-one-out analysis are similar to the effect sizes computed by the GLMs’ leave-one-out analysis

Figure 10–Figure supplement 2. Pairwise scatterplots of MTNN single-covariate effect sizes.

419 Discussion

420 We set out to test whether electrophysiological responses, notoriously variable across labs, could
421 be reproducible across geographically separated laboratories after appropriate standardization
422 of experiments, data processing, and analyses. After applying stringent behavioral, histological,
423 and electrophysiological quality-control criteria, we found that electrophysiological features such
424 as neuronal yield, firing rate, and normalized LFP power were reproducible across laboratories;
425 their within-lab averages did not significantly deviate from the mean across labs. Similarly, the pro-
426 portion of cells whose responses are tuned to behaviorally-relevant task events is reproducible
427 across labs. Finally, a multi-task neural network approach can predict the firing rates of differ-
428 ent units across sessions, and again, the within-lab random effects estimated by this model were
429 comparable to between-lab random effects. Taken together, our results suggest that careful stan-
430 dardization can lead to reproducible electrophysiological results across laboratories.

431 Reproducibility in our electrophysiology studies depended on rigorous metrics of quality. We
432 found that it was necessary to exclude a significant fraction of datasets to reach a desired level
433 of reproducibility. Quality control was enforced for diverse aspects of the experiments, including
434 histology, behavior, targeting, neuronal yield, and the total number of completed sessions. Among
435 these measures, recordings with high noise and low neuronal yield were significantly represented
436 in sessions that were excluded (40/74 sessions). A number of issues contributed here, including
437 artifacts present in the recordings, inadequate grounding, and a decline in craniotomy health; all
438 of these can potentially be improved with experimenter experience. Sub-standard behavior (for in-
439 stance, too few trials in a session) led to the elimination of another substantial fraction of datasets.
440 Trial counts are likely to be highly variable across labs, as there is currently no agreed upon stan-
441 dard for what constitutes suitable behavior for an electrophysiology experiment. This has already
442 been shown to cause differences in the internal states visited by animals as they make decisions
443 (*Ashwood et al., 2022*).

444 These observations suggest that future experiments might enjoy greater reproducibility if re-
445 searchers followed, or at least reported, a number of agreed upon criteria, such as those we de-
446 fine in Table 1. This approach has been successful in other fields: for instance, the neuroimaging
447 field has agreed upon a set of guidelines for “best practices,” and has identified factors that can
448 impede those practices (*Nichols et al., 2017*). The genomics field likewise adopted the Minimum
449 Information about a Microarray Experiment (MIAME) standard, designed to ensure that data from
450 microarrays could be meaningfully interpreted and experimentally verified (*Brazma et al., 2001*).
451 Our work here suggests the creation of a similar set of standards for electrophysiology and be-
452 havioral experiments would be beneficial. These could include expectations for reporting (such
453 as histological information and behavioral trial numbers) as well as suggestions for minimizing
454 variability (e.g., agreed upon standards for the noise level that would exclude a recording).

455 We found probe targeting to be a large source of variability, driven by micro-manipulator po-
456 sitioning and anatomical discrepancies. The majority of the variance in targeting was due to the
457 probe entry positions at the brain surface, which showed no bias in placement across the dataset.
458 The source of this variance could be due to a discrepancy in skull landmarks compared to the un-
459 derlying brain anatomy. Accuracy in placing probes along a planned trajectory is therefore limited
460 by this variability (about 400 μ m). Probe angle also showed a small degree of variance, and a bias in
461 both anterior-posterior and medio-lateral directions; indicating that the Allen Common Coordinate
462 Framework (CCF) (*Wang et al., 2020*) and stereotaxic coordinate systems are slightly offset. Mini-
463 mizing variance in probe targeting is an important element in increasing reproducibility, as slight
464 deviations in probe entry position and angle can lead to samples from different populations of
465 neurons. Our approach suggests a path forward to minimize these biases: probe angles must be
466 carefully computed from the CCF, as the CCF and stereotaxic coordinate systems do not define the
467 same coronal plane angle. Small differences in probe location may be responsible for other stud-
468 ies arriving at different conclusions, highlighting the need for agreed upon methods for targeting

469 specific areas (*Rajasethupathy et al., 2015; Andrianova et al., 2022*).

470 Our results also highlight the critical importance of reproducible histological processing and
471 subsequent probe alignment. Specifically, we used a centralized histology and registration pipeline
472 to assign each recording site on each probe to a particular anatomical location, based on registra-
473 tion of the histological probe trajectories to the CCF and the electrophysiological features recorded
474 at each site. This differs from previous approaches, in which stereotaxic coordinates alone were
475 used to target an area of interest and exclusion criteria were not specified; see e.g. (*Najafi et al.,*
476 *2020; Harvey et al., 2012; Goard et al., 2016; Raposo et al., 2014; Erlich et al., 2015*). The reliance on
477 stereotaxic coordinates for localization, instead of standardized histological registration, is a possi-
478 ble explanation for conflicting results across labs. Our results speak to the importance of adopting
479 standardized procedures more broadly across laboratories.

480 A major contribution of our work is open-source data and code: we share our full dataset ([link](#)
481 [to data portal](#)) and suite of analysis tools for quantifying reproducibility ([link to code repository](#)).
482 The analyses here required significant improvements in data architecture, visualization, spike sort-
483 ing, histology image analysis, and video analysis. Our analyses uncovered major gaps and issues
484 in the existing toolsets that required improvements (see Methods and *The International Brain*
485 *Laboratory (2021a,b)* for full details); the large-scale dataset analyzed here proved to be a use-
486 ful stress test pointing to improved analysis pipelines. For example, we improved existing spike
487 sorting pipelines with regard to scalability, reproducibility, and stability. These improvements con-
488 tribute towards advancing automated spike sorting, and move beyond subjective manual curation,
489 which scales poorly and limits reproducibility. We anticipate that our open-source dataset will play
490 an important role in further improvements to these pipelines and also the development of further
491 methods for modeling the spike trains of many simultaneously recorded neurons across multiple
492 brain areas and experimental sessions.

493 Scientific advances rely on the reproducibility of scientific findings. The current study demon-
494 strates that reproducibility is attainable for large-scale neural recordings during a standardized
495 perceptual detection task across 9 laboratories. We offer several recommendations to increase
496 reproducibility, including (1) standardized protocols for data collection, (2) data processing, and
497 (3) rigorous data quality metrics. Furthermore, we have made improvements in data architecture
498 and processing, now available to the public. Our study provides a framework for the collection and
499 analysis of large neural datasets in a reproducible manner that will play a key role as neuroscience
500 continues to move towards increasingly complex datasets.

501 Resources

502 Data access

503 Please visit https://int-brain-lab.github.io/iblenv/notebooks_external/data_release_repro_ephys.html
504 to access the data used in this article.

505 Code repository

506 Please visit <https://github.com/int-brain-lab/paper-reproducible-ephys/> to access the code used to
507 produce the results and figures presented in this article.

508 Protocols and pipelines

509 Please visit https://figshare.com/projects/Reproducible_Electrophysiology/138367 to access the proto-
510 cols and pipelines used in this article.

511 Methods and Materials

512 All procedures and experiments were carried out in accordance with local laws and following ap-
513 proval by the relevant institutions: the Animal Welfare Ethical Review Body of University College
514 London; the Institutional Animal Care and Use Committees of Cold Spring Harbor Laboratory,
515 Princeton University, and University of California at Berkeley; the University Animal Welfare Com-
516 mittee of New York University; and the Portuguese Veterinary General Board.

517 Animals

518 Mice were housed under a 12/12 h light/dark cycle (normal or inverted depending on the labora-
519 tory) with food and water available ad libitum, except during behavioural training days. Electro-
520 physiological recordings and behavioural training were performed during either the dark or light
521 phase of the cycle depending on the laboratory. N=48 adult mice (C57BL/6, male and female, ob-
522 tained from either Jackson Laboratory or Charles River) were used in this study. Mice were aged
523 17-41 weeks and weighed 16.4-34.5 g on the day of the headbar implant surgery.

524 Materials and apparatus

525 For detailed parts lists and installation instructions, see Appendix 1 (*The International Brain Lab-*
526 *oratory, 2022a*).

527 Briefly, each lab installed a standardized electrophysiological rig (named ‘ephys rig’ throughout
528 this text), which differed slightly from the apparatus used during behavioral training (*The Interna-*
529 *tional Brain Laboratory et al., 2021*). The general structure of the rig was constructed from Thor-
530 labs parts and was placed inside a custom acoustical cabinet clamped on an air table (Newport,
531 M-VIS3036-SG2-325A). A static head bar fixation clamp and a 3D-printed mouse holder were used
532 to hold a mouse such that its forepaws rest on the steering wheel (86652 and 32019, LEGO) (*The*
533 *International Brain Laboratory et al., 2021*). Silicone tubing controlled by a pinch valve (225P011-
534 21, NResearch) was used to deliver water rewards to the mouse. The display of the visual stimuli
535 occurred on a LCD screen (LP097Q × 1, LG). To measure the precise times of changes in the visual
536 stimulus, a patch of pixels on the LCD screen flipped between white and black at every stimulus
537 change, and this flip was captured with a photodiode (Bpod Frame2TTL, Sanworks). Ambient tem-
538 perature, humidity, and barometric air pressure were measured with the Bpod Ambient module
539 (Sanworks), wheel position was monitored with a rotary encoder (05.2400.1122.1024, Kubler).

540 Videos of the mouse were recorded from 3 angles (left, right and body) with USB cameras (CM3-
541 U3-13Y3M-CS, Point Grey) sampling at 60, 150, 30 Hz respectively (for details see Appendix 1 (*The*
542 *International Brain Laboratory, 2022a*)). A custom speaker (Hardware Team of the Champalimaud
543 Foundation for the Unknown, V1.1) was used to play task-related sounds, and an ultrasonic mi-
544 crophone (Ultramic UM200K, Dodotronic) was used to record ambient noise from the rig. All task-
545 related data was coordinated by a Bpod State Machine (Sanworks). The task logic was programmed

546 in Python and the visual stimulus presentation and video capture was handled by Bonsai (*Lopes*
547 *et al., 2015*) and the Bonsai package BonVision (*Lopes et al., 2021*).

548 All recordings were made using Neuropixels probes (Imec, 3A and 3B models), advanced in the
549 brain using a micromanipulator (Sensapex, uMp-4) tilted by a 15 degree angle from the vertical
550 line. The aimed electrode penetration depth was 4.0 mm. Data were acquired via an FPGA (for 3A
551 probes) or PXI (for 3B probes, National Instrument) system and stored on a PC.

552 **Headbar implant surgery**

553 A detailed account of the surgical methods is in Appendix 1 (*The International Brain Laboratory*
554 *et al., 2021*).

555 Briefly, mice were anesthetized with isoflurane and head-fixed in a stereotaxic frame. The hair
556 was then removed from their scalp, much of the scalp and underlying periosteum was removed
557 and bregma and lambda were marked. Then the head was positioned such that there was a 0
558 degree angle between bregma and lambda in all directions. The head bar was then placed in
559 one of three stereotactically defined locations and cemented in place. The location of the future
560 craniotomies were measured using a pipette referenced to bregma, and marked on the skull using
561 either a surgical blade or pen. The exposed skull was then covered with cement and clear UV curing
562 glue, ensuring that the remaining scalp was unable to retract from the implant.

563 **Behavioral training and habituation to the ephys rig**

564 For a detailed protocol on animal training, see Appendix 2 (*The International Brain Laboratory*
565 *et al., 2021*).

566 Once the mouse is classified as having learned the biasedChoiceWorld task (criteria 'ready4ephyRig'
567 reached, cf Appendix 2 for definition (*The International Brain Laboratory et al., 2021*)), it is trans-
568 ferred onto the ephys rig.

569 The mouse is habituated to behave on the ephys rig in a series of steps that do not involve
570 any electrophysiology recording. First, the mouse needs to perform one session of biasedChoice-
571 World on the electrophysiology rig, with at least 400 trials and 90% correct on easy contrasts (col-
572 lapsing across block types). Once this criterion is reached, time delays are introduced prior to the
573 task; these delays would eventually serve to mimic the time it would take to insert electrodes in
574 the brain. The mouse has to maintain performance for 3 subsequent sessions (same criterion as
575 'ready4ephyRig'), but with a minimum of one session that has a 15 minutes delay and is a mock
576 recording.

577 **Electrophysiological recording using Neuropixels probes**

578 Data acquisition

579 For details, see Appendix 2 and 3 (*The International Brain Laboratory, 2022b,c*).

580 Briefly, upon the day of electrophysiological recording, the animal was anaesthetised using
581 isoflurane and surgically prepared. The cement and glue were removed, exposing the skull over
582 both hemispheres. A test was made to check whether the implant could hold liquid, and if suc-
583 cessful a grounding pin was implanted. One or two craniotomies (1 × 1 mm) were made over the
584 marked locations. The dura was left intact, and the brain was lubricated with ACSF. DuraGel was
585 applied over the dura as a moisturising sealant, and covered with a layer of Kwikcast. The mouse
586 was administered with analgesics subcutaneously, and left to recover in a heating chamber until
587 locomotor and grooming activity were fully recovered.

588 Once the animal was recovered from the craniotomy, it was fixed in the apparatus. Once a
589 craniotomy was made, up to 4 subsequent recording sessions were made in that same craniotomy.
590 Up to two probes were implanted in the brain on a given session. The probes were labelled with
591 CM-Dil (see Appendix 4 (*The International Brain Laboratory, 2022d*) and (*Liu, 2019*)).

592 Spike sorting

593 The spike sorting pipeline used at IBL is described in details in (*The International Brain Laboratory*
594 *et al., 2022a*). Briefly, spike sorting was performed using a modified version of the Kilosort 2.5
595 algorithm (*Steinmetz et al., 2021*). We found it necessary to improve the original code in several
596 aspects (scalability, reproducibility, and stability, discussed below), and developed an open-source
597 Python port; the code repository is here: (*The International Brain Laboratory, 2021b*).

598 Regarding scalability: we found that the original code failed on recording sessions with a large
599 number of detected spikes. Therefore we improved the CPU memory usage of the code to better
600 handle these cases.

601 Regarding reproducibility: spike sorting algorithms are still in heavy development; we needed to
602 tag and validate code versions and parameter settings internally so we could release the algorithm
603 to our data-processing computers across multiple labs on our own schedule. We also defined
604 a set of integration tests on short (100 seconds) recordings, using hybrid ground-truth datasets
605 (*Pachitariu et al., 2016*) to validate algorithm changes before new version releases.

606 Regarding stability: we observed a number of clear artifacts in the raw Neuropixels output
607 ("dead" channels, simultaneous "glitch" artifacts across multiple channels, mis-alignment errors,
608 etc.) that were not handled properly by previous algorithms. We developed new methods to han-
609 dle each of these artifact types, resulting in significantly more stable sorting outputs. See (*The*
610 *International Brain Laboratory et al., 2022a*) for full details.

611 Local field potential (LFP)

612 Concurrently with the action potential band, each channel of the Neuropixel probe recorded a low-
613 pass filtered trace at a sampling rate of 2500 Hz. The power spectral density at different frequencies
614 was estimated per channel using the Welch's method with partly overlapping Hanning windows of
615 1024 samples. Power spectral density (PSD) was converted into dB as follows:

$$dB = 10 * \log(PSD) \quad (1)$$

616 Serial section two-photon imaging

617 Mice were given a terminal dose of pentobarbital and perfuse-fixed with PBS followed by 4%
618 formaldehyde solution (ThermoFisher 28908) in 0.1M PB pH 7.4. Whole mouse brain was dissected,
619 and post-fixed in the same fixative for a minimum of 24 hours at room temperature. Tissues were
620 washed and stored for up to 2-3 weeks in PBS at 4C, prior to shipment to the Sainsbury Wellcome
621 Centre for image acquisition. For full details, see Appendix 5 (*The International Brain Laboratory,*
622 *2022e*).

623 For imaging, brains were equilibrated with 50mM PB solution and embedded into 5% agarose
624 gel blocks. The brains were imaged using serial section two-photon microscopy (*Ragan et al., 2012;*
625 *Economo et al., 2016*). The microscope was controlled with ScanImage Basic (Vidrio Technologies,
626 USA), and BakingTray, a custom software wrapper for setting up the imaging parameters (*Camp-*
627 *bell, 2020*). Image tiles were assembled using StitchIt (*Campbell, 2021*). Whole brain coronal image
628 stacks were acquired with a resolution of 4.4 x 4.4 x 25.0 μm in XYZ, with a two-photon laser wave-
629 length of 920nm, and power of 35% of 1800mW from the source laser, yielding approximately
630 150mW at the block face. Serial section microscopy proceeded with 2 z slices taken for each 50 μm
631 tissue slice, at a depth of 30 μm and 55 μm from the tissue surface. Two channels of image data
632 was acquired on two PMTs for green (bandpass filter ET525/50m) and red (bandpass filter ET570lp)
633 fluorescence.

634 Whole brain images were downsampled to 25 μm XYZ pixels and registered to the adult mouse
635 Allen common coordinate framework (*Wang et al., 2020*) using BrainRegister (*West, 2021*), an elastix-
636 based (*Klein et al., 2010*) registration pipeline with optimised parameters for mouse brain registra-
637 tion. For full details, see Appendix 7 (*The International Brain Laboratory, 2022g*).

638 **Probe track tracing and alignment**

639 Neuropixels probe tracks were manually traced to yield a probe trajectory using Lasagna (*Camp-*
640 *bell et al., 2020*), a Python-based image image viewer equipped with a plugin tailored for this task.
641 Traced probe track data was uploaded to an Alyx server (*Rossant et al., 2021*); a database designed
642 for experimental neuroscience laboratories. Neuropixels channels were then manually aligned to
643 anatomical features along the trajectory using electrophysiological landmarks with [ephys align-
644 ment tool] (*Faulkner, 2020*) (*Liu et al., 2021*). For full details, see Appendix 6 (*The International*
645 *Brain Laboratory, 2022f*).

646 **Permutation tests**

647 We use permutation tests to study the reproducibility of neural features across laboratories. To
648 this end, we first defined a test statistic that is sensitive to systematic deviations between labora-
649 tories: the sum of the absolute differences between laboratory means and overall mean. The
650 null-hypothesis is that there is no difference between the different laboratory means, i.e. the
651 assignment of mice to laboratories is completely random. We constructed the corresponding
652 null-distribution by permuting these assignments between laboratories and mice randomly 10000
653 times (leaving the relative numbers of mice in laboratories intact) and computing the test statistic
654 on these randomised samples. Given this constructed null-distribution, the p-value of the permu-
655 tation test is the proportion of the null-distribution that has more extreme values than the test
656 statistic that was computed on the real data.

657 **Dimensionality reduction of peri-event time histograms via principal component** 658 **analysis**

659 In Figure 6 we use principal component analysis (PCA) to embed peri-event time histograms (PETHs)
660 into a two-dimensional feature space for visualization and further analysis. Our overall approach
661 is to compute PETHs, split into fast-reaction-time and slow-reaction-time trials, then concatenate
662 these PETH vectors for each cell to obtain an informative summary of each cell's activity. Next
663 we stack these double PETHs from all labs into a single matrix and use PCA to obtain a low-rank
664 approximation of this PETH matrix.

665 In detail, the two PETHs consist of one averaging fast reaction time ($< 0.15sec$) trials and the
666 other slow reaction time ($> 0.15sec$) trials, each of length T time steps. We used 20ms bins, from
667 $-0.5sec$ to $1.5sec$ relative to motion onset, so $T = 100$. We also performed a simple normalization
668 on each PETH, dividing the firing rates by the baseline firing rate (prior to motion onset) of each
669 cell plus a small positive offset term (to avoid amplifying noise in very low-firing cells), following
670 *Steinmetz et al. (2021)*.

671 Let the stack of these double PETH vectors be Y , being a $N \times 2T$ matrix, where N is the total num-
672 ber of neurons recorded across 5 brain regions and labs. Running principal components analysis
673 (PCA) on Y (singular value decomposition) is used to obtain the low-rank approximation $UV \approx Y$.
674 This provides a simple low-d embedding of each cell: U is $N \times k$, with each row of U representing
675 a k -dimensional embedding of a cell that can be visualized easily across labs and brain regions. V
676 is $k \times 2T$ and corresponds to the k temporal basis functions that PCA learns to best approximate
677 Y . Figure 6(a) shows two cells of Y and the corresponding PCA approximation from UV .

678 The scatter plots in Figure 6 show the embedding U across labs and brain regions, with the
679 embedding dimension $k = 2$. Each $k \times 1$ vector in U , corresponding to a single cell, is assigned to a
680 single dot in Figure 6c.

681 **Video analysis**

682 Some of the behavioral time series used in the neural network analysis are derived from video
683 recordings of the animals. Full details of the video analysis pipeline are here: (*The International*
684 *Brain Laboratory et al., 2022b*), and the code is available here: (*The International Brain Laboratory,*
685 *2021a*).

686 Briefly, in the recording rigs, there are three cameras, one called 'left' at full resolution 1280x1024
687 and 60 Hz filming the mouse from one side, one called 'right' at half resolution (640x512) and 150
688 Hz, filming the mouse symmetrically from the other side, and one called 'body' filming the trunk of
689 the mouse from above. Several quality control metrics were developed to detect video issues such
690 as poor illumination (as infra red light bulbs broke) or accidental misplacement of the cameras.

691 Marker-less tracking of body parts is achieved using Deeplabcut (*Mathis et al., 2018*), a deep-
692 learning-based tool that is used within a fully automated pipeline in IBL to track various body parts
693 such as the paws. The pipeline first detects 3 regions of interest (ROI) in each frame, crops these
694 ROIs using ffmpeg (*Tomar, 2006*) and applies a separate network for each ROI to track features.
695 For each side video we track the following points:

696 • ROI eye:

697 'pupil_top_r', 'pupil_right_r', 'pupil_bottom_r', 'pupil_left_r'

698 • ROI mouth:

699 'nose_tip', 'tongue_end_r', 'tongue_end_l'

700 • ROI paws:

701 'paw_r', 'paw_l'

702 The right side video was flipped and spatially up-sampled to look like the left side video, such that
703 we could apply the same Deeplabcut networks.

704 Extensive curating of the training set of images for each network was required to obtain reliable
705 tracking across animals and laboratories. We annotated in total more than 10K frames, across sev-
706 eral iterations, using a semi-automated tracking failure detection approach, which found frames
707 with temporal jumps, 3d re-projection errors when combining both side views, and heuristic mea-
708 sures of spatial violations. These selected 'bad' frames were then annotated and the network re-
709 trained. To find further raw video and Deeplabcut issues, we inspected trial-averaged behaviors
710 obtained from the tracked features, such as licking aligned to feedback time, paw speed aligned
711 to stimulus onset and scatter plots of animal body parts across a session superimposed onto ex-
712 ample video frames. See (*The International Brain Laboratory et al., 2022b*) for further details and
713 example quality control images.

714 **Multi-task neural network model to quantify sources of variability**

715 Data preprocessing

716 For the Multi-task neural network (MTNN) analysis, we used data from 20 sessions recorded in
717 CCU, CSHL (C), SWC, Berkeley, and NYU. We included various covariates in our feature set (e.g. go-
718 cue signals, stimulus/reward type, Deep Lab Cut behavioral outputs). For the "decision strategy"
719 covariate, we used the posterior estimated state probabilities of the 4-state GLM-HMMs trained
720 on the sessions used for the MTNN analysis (*Ashwood et al., 2022*). Both biased and unbiased
721 data were used when training the 4-state model. For each session, we first filtered out the trials
722 where no choice is made. We then selected the trials whose stimulus onset time is within 0.4
723 seconds before the first movement onset time and feedback time is within 0.9 seconds after the
724 first movement onset time. Finally, we selected responsive units whose mean firing rate is greater
725 than 5 spikes/second for further analyses. For sessions with more than 15 responsive units, we
726 randomly sampled 15 units.

727 Model Architecture

728 Given a set of covariates in Table 2, the MTNN predicts the target sequence of firing rates from
729 0.5 seconds before first movement onset to 1 second after, with bin width set to 50 ms (30 time
730 bins). More specifically, a sequence of feature vectors $x_{\text{dynamic}} \in \mathbb{R}^{D_{\text{dynamic}} \times T}$ that include dynamic
731 covariates, such as Deep Lab Cut (DLC) outputs, and wheel velocity, and a feature vector $x_{\text{static}} \in$

732 $\mathbb{R}^{D_{\text{static}}}$ that includes static covariates, such as the lab ID, unit's 3-D location, are input to the MTNN
 733 to compute the prediction $y^{\text{pred}} \in \mathbb{R}^T$, where D_{static} is the number of static features, D_{dynamic} is the
 734 number of dynamic features, and T is the number of time bins. The MTNN has initial layers that
 735 are shared by all units, and each unit has its designated final fully-connected layer.

Given the feature vectors x_{dynamic} and x_{static} for session s and unit u , the model predicts the firing rates y^{pred} by:

$$e_{\text{static}} = f(w_{\text{static}}^T x_{\text{static}} + b_{\text{static}}) \quad (2)$$

$$e_{\text{dynamic}} = f(w_{\text{dynamic}}^T x_{\text{dynamic}} + b_{\text{dynamic}}) \quad (3)$$

$$h_t^{(\text{forward})} = \max(0, U_1 e_{\text{dynamic}, t} + V_1 h_{t-1}^{(\text{forward})} + b_{\text{forward}}) \quad (4)$$

$$h_t^{(\text{backward})} = \max(0, U_2 e_{\text{dynamic}, t} + V_2 h_{t+1}^{(\text{backward})} + b_{\text{backward}}) \quad (5)$$

$$y_t^{\text{pred}} = f(w_{(s,u)}^T \text{concat}(e_{\text{static}}, h_t^{(\text{forward})}, h_t^{(\text{backward})}) + b_{(s,u)}) \quad (6)$$

736 where f is the activation function. Eqn. (2) and Eqn. (3) are the shared fully-connected layers
 737 for static and dynamic covariates, respectively. Eqn. (4) and Eqn. (5) are the shared one-layer
 738 bidirectional recurrent neural networks (RNNs) for dynamic covariates, and Eqn. (6) is the unit-
 739 specific fully-connected layer, indexed by (s, u) . Each part of the MTNN architecture can have an
 740 arbitrary number of layers. For our analysis, we used two fully-connected shared layers for static
 741 covariates (Eqn. (2)) and three-layer bidirectional RNNs for dynamic covariates, with the embedding
 742 size set to 64.

743 Model training

744 The model was implemented in PyTorch and trained on a single GPU. The training was performed
 745 using Stochastic Gradient Descent on the Poisson negative loglikelihood (Poisson NLL) loss with
 746 learning rate set to 0.1, momentum set to 0.9, and weight decay set to 10^{-15} . We used a learning
 747 rate scheduler such that the learning rate for the i -th epoch is 0.1×0.95^i , and the dropout rate was
 748 set to 0.2. We also experimented with mean squared error (MSE) loss instead of Poisson NLL loss,
 749 and the results were similar. The batch size was set to 512.

750 The dataset consists of 20 sessions, 246 units and 6878 active trials in total. For each session,
 751 20% of the trials are used as the test data and the remaining trials are split 20:80 for the validation
 752 and training sets. During training, the performance on the held-out validation set is checked after
 753 every 3 passes through the training data. The model is trained for 100 epochs, and the model
 754 parameters with the best performance on the held-out validation set are saved and used for pre-
 755 dictions on the test data.

756 Simulated experiments

757 For the simulated experiment in Figure 10 supplemental 1, we first trained GLMs on the same set
 758 of 246 responsive neural units from 20 sessions used for the analysis in Figure 10, with a reduced
 759 set of covariates consisting of stimulus timing, stimulus side and contrast, first movement onset
 760 timing, feedback type and timing, wheel velocity, and mouse's priors for the current and previous
 761 trials. The kernels of the trained GLMs show the contribution of each of the covariates to the firing
 762 rates of each unit. For each simulated unit, we used these kernels of the trained GLM to simulate
 763 its firing rates for 350 randomly initialized trials. The random trials were 1.5 seconds long with 50
 764 ms bin width. For all trials, the first movement onset timing was set to 0.5 second after the start
 765 of the trial, and the stimulus contrast, side, onset timing and feedback type, timing were randomly
 766 sampled. We used wheel velocity traces and mouse's priors from real data for simulation. We
 767 finally ran the leave-one-out analyses with GLMs/MTNN on the simulated data and compared the
 768 effect sizes estimated by GLMs and MTNN.

769

770 Acknowledgments

771 This work was supported by grants from the Wellcome Trust (209558 and 216324), National Insti-
772 tutes of Health (1U19NS123716) and the Simons Foundation. We thank T. Zador, P. Dayan, and C.
773 Hurwitz for helpful comments on the manuscript. The production of all IBL Platform Papers is led
774 by a Task Force, which defines the scope and composition of the paper, assigns and/or performs
775 the required work for the paper, and ensures that the paper is completed in a timely fashion. The
776 Task Force members for this platform paper include authors SAB, GC, AC, MFD, HDL, MF, GM, LP,
777 NR, MS, NS, MT, and SW.

778 References

- 779 **Andrianova L**, Yanakieva S, Margetts-Smith G, Kohli S, Brady ES, Aggleton JP, Craig MT. No evidence from com-
780plementary data sources of a direct projection from the mouse anterior cingulate cortex to the hippocampal
781formation. *bioRxiv*. 2022; .
- 782 **Ashwood ZC**, Roy NA, Stone IR, Urai AE, Churchland AK, Pouget A, Pillow JW. Mice alternate between discrete
783strategies during perceptual decision-making. *Nature Neuroscience*. 2022; 25(2):201–212.
- 784 **Baker M**. 1,500 scientists lift the lid on reproducibility. *Nature*. 2016; 533(7604). doi: 10.1038/533452a.
- 785 **Barry C**, Ginzberg LL, O’Keefe J, Burgess N. Grid cell firing patterns signal environmental novelty by expansion.
786 *Proceedings of the National Academy of Sciences*. 2012; 109(43):17687–17692.
- 787 **Batty E**, Merel J, Brackbill N, Heitman A, Sher A, Litke A, Chichilnisky E, Paninski L. Multilayer recurrent network
788models of primate retinal ganglion cell responses. *ICLR*. 2016; .
- 789 **Benjamini Y**, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple
790Testing. *Journal of the Royal Statistical Society Series B (Methodological)*. 1995; 57(1):289–300. <https://www.jstor.org/stable/2346101>, publisher: [Royal Statistical Society, Wiley].
- 791
- 792 **Bragin A**, Jando G, Nadasdy Z, van Landeghem M, Buzsáki G. Dentate EEG spikes and associated interneuronal
793population bursts in the hippocampal hilar region of the rat. *Journal of Neurophysiology*. 1995; 73(4):1691–
7941705. doi: 10.1152/jn.1995.73.4.1691.
- 795 **Brazma A**, Hingamp P, Quackenbush J, Sherlock G, Spellman P, Stoeckert C, Aach J, Ansorge W, Ball CA, Causton
796HC, et al. Minimum information about a microarray experiment (MIAME)—toward standards for microarray
797data. *Nature genetics*. 2001; 29(4):365–371.
- 798 **Cadena SA**, Denfield GH, Walker EY, Gatys LA, Tolias AS, Bethge M, Ecker AS. Deep convolutional models im-
799prove predictions of macaque V1 responses to natural images. *PLOS Computational Biology*. 2019; 15(4):1–
80027. doi: 10.1371/journal.pcbi.1006897.
- 801 **Campbell R**, BakingTray; 2020. <https://github.com/SainsburyWellcomeCentre/BakingTray>, doi:
802<https://doi.org/10.5281/zenodo.3631609>.
- 803 **Campbell R**, StitchIt; 2021. <https://github.com/SainsburyWellcomeCentre/StitchIt>, doi:
804<https://zenodo.org/badge/latestdoi/57851444>.
- 805 **Campbell R**, Blot A, Rousseau C, Winter O, Lasagna; 2020. <https://github.com/SainsburyWellcomeCentre/lasagna>,
806doi: 10.5281/zenodo.3941894.
- 807 **Chen G**, Manson D, Cacucci F, Wills TJ. Absence of visual input results in the disruption of grid cell firing in the
808mouse. *Current Biology*. 2016; 26(17):2335–2342.
- 809 **Churchland AK**, Kiani R, Chaudhuri R, Wang XJ, Pouget A, Shadlen MN. Variance as a signature of neural
810computations during decision making. *Neuron*. 2011; 69(4):818–31. <http://www.ncbi.nlm.nih.gov/pubmed/21338889>, doi: 10.1016/j.neuron.2010.12.037.
- 811
- 812 **Churchland MM**, Yu BM, Cunningham JP, Sugrue LP, Cohen MR, Corrado GS, Newsome WT, Clark AM, Hosseini
813P, Scott BB, Bradley DC, Smith MA, Kohn A, Movshon JA, Armstrong KM, Moore T, Chang SW, Snyder LH,
814Lisberger SG, Priebe NJ, et al. Stimulus onset quenches neural variability: a widespread cortical phenomenon.
815 *Nat Neurosci*. 2010; 13(3):369–78. <http://www.ncbi.nlm.nih.gov/pubmed/20173745>, doi: 10.1038/nn.2501.
- 816 **Dragoi G**, Tonegawa S. Preplay of future place cell sequences by hippocampal cellular assemblies. *Nature*.
8172011; 469(7330):397–401.

- 818 **Economo MN**, Clack NG, Lavis LD, Gerfen CR, Svoboda K, Myers EW, Chandrashekar J. A platform for brain-wide
819 imaging and reconstruction of individual neurons. *eLife*. 2016; 5(e10566). doi: [10.7554/eLife.10566](https://doi.org/10.7554/eLife.10566).
- 820 **Erllich JC**, Brunton BW, Duan CA, Hanks TD, Brody CD. Distinct effects of prefrontal and parietal cortex inacti-
821 vations on an accumulation of evidence task in the rat. *Elife*. 2015; 4:e05457.
- 822 **Errington TM**, Mathur M, Soderberg CK, Denis A, Perfito N, Iorns E, Nosek BA. Investigating the replicability
823 of preclinical cancer biology. *eLife*. 2021; 10:e71601. doi: [10.7554/eLife.71601](https://doi.org/10.7554/eLife.71601), publisher: eLife Sciences
824 Publications, Ltd.
- 825 **Faulkner M**, Ephys Atlas GUI; 2020. <https://github.com/int-brain-lab/iblapps/tree/master/atlas electrophysiology>.
- 826 **Goard MJ**, Pho GN, Woodson J, Sur M. Distinct roles of visual, parietal, and frontal motor cortices in memory-
827 guided sensorimotor decisions. *elife*. 2016; 5:e13764.
- 828 **Grosmark AD**, Buzsáki G. Diversity in neural firing dynamics supports both rigid and learned hippocampal
829 sequences. *Science*. 2016; 351(6280):1440–1443.
- 830 **Hafting T**, Fyhn M, Molden S, Moser MB, Moser EI. Microstructure of a spatial map in the entorhinal cortex.
831 *Nature*. 2005; 436(7052):801–806.
- 832 **Harvey CD**, Coen P, Tank DW. Choice-specific sequences in parietal cortex during a virtual-navigation decision
833 task. *Nature*. 2012; 484(7392):62–68.
- 834 **Izenman AJ**. Reduced-rank regression for the multivariate linear model. *Journal of multivariate analysis*. 1975;
835 5(2):248–264.
- 836 **Jun JJ**, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA, Andrei A, Aydın Ç, et al.
837 Fully integrated silicon probes for high-density recording of neural activity. *Nature*. 2017; 551(7679):232–236.
- 838 **Klein S**, Staring M, Murphy K, Viergever MA, Pluim JPW. elastix: a toolbox for intensity based medical image
839 registration. *IEEE Transactions on Medical Imaging*. 2010; 29(1):196–205. doi: [10.1109/TMI.2009.2035616](https://doi.org/10.1109/TMI.2009.2035616).
- 840 **Kobak D**, Brendel W, Constantinidis C, Feierstein CE, Kepecs A, Mainen ZF, Qi XL, Romo R, Uchida N, Machens
841 CK. Demixed principal component analysis of neural population data. *Elife*. 2016; 5:e10989.
- 842 **Li X**, Ai L, Giavasis S, Jin H, Feczko E, Xu T, Clucas J, Franco A, Sólón Heinsfeld A, Adebimpe A, Vogelstein JT,
843 Yan CG, Esteban O, Poldrack RA, Craddock C, Fair D, Satterthwaite T, Kiar G, Milham MP. Moving Beyond
844 Processing and Analysis-Related Variation in Neuroscience. *bioRxiv*. 2021; [https://www.biorxiv.org/content/
845 early/2021/12/03/2021.12.01.470790](https://www.biorxiv.org/content/early/2021/12/03/2021.12.01.470790), doi: [10.1101/2021.12.01.470790](https://doi.org/10.1101/2021.12.01.470790).
- 846 **Liu L**. Painting Neuropixels probes and other silicon probes for electrophysiological recordings. *protocolsio*.
847 2019; doi: [dx.doi.org/10.17504/protocols.io.wxqffmw](https://doi.org/10.17504/protocols.io.wxqffmw).
- 848 **Liu LD**, Chen S, Hou H, West SJ, Faulkner M, Economo MN, Li N, Svoboda K, the International Brain Labora-
849 tory. Accurate localization of linear probe electrode arrays across multiple brains. *eNeuro*. 2021; 8(6). doi:
850 [10.1523/ENEURO.0241-21.2021](https://doi.org/10.1523/ENEURO.0241-21.2021).
- 851 **Liu Y**, Dolan RJ, Kurth-Nelson Z, Behrens TE. Human replay spontaneously reorganizes experience. *Cell*. 2019;
852 178(3):640–652.
- 853 **Lopes G**, Bonacchi N, Frazão J, Neto JP, Atallah BV, Soares S, Moreira L, Matias S, Itskov PM, Correia PA, et al. Bon-
854 sai: an event-based framework for processing and controlling data streams. *Frontiers in neuroinformatics*.
855 2015; 9:7.
- 856 **Lopes G**, Farrell K, Horrocks EA, Lee CY, Morimoto MM, Muzzu T, Papanikolaou A, Rodrigues FR, Wheatcroft T,
857 Zucca S, et al. Creating and controlling visual environments using BonVision. *Elife*. 2021; 10:e65541.
- 858 **Mathis A**, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, Bethge M. DeepLabCut: markerless pose
859 estimation of user-defined body parts with deep learning. *Nature neuroscience*. 2018; 21(9):1281–1289.
- 860 **McIntosh LT**, Maheswaranathan N, Nayebi A, Ganguli S, Baccus SA. Deep learning models of the retinal re-
861 sponse to natural scenes. *Advances in neural information processing systems*. 2016; 29:1369.
- 862 **Musall S**, Kaufman MT, Juavinett AL, Gluf S, Churchland AK. Single-trial neural dynamics are dominated by
863 richly varied movements. *Nature neuroscience*. 2019; 22(10):1677–1686.

- 864 **Najafi F**, Elsayed GF, Cao R, Pnevmatikakis E, Latham PE, Cunningham JP, Churchland AK. Excitatory and in-
865 hibitory subnetworks are equally selective during decision-making and emerge simultaneously during learn-
866 ing. *Neuron*. 2020; 105(1):165–179.
- 867 **Nichols TE**, Das S, Eickhoff SB, Evans AC, Glatard T, Hanke M, Kriegeskorte N, Milham MP, Poldrack RA, Poline
868 JB, et al. Best practices in data analysis and sharing in neuroimaging using MRI. *Nature neuroscience*. 2017;
869 20(3):299–303.
- 870 **Ólafsdóttir HF**, Barry C, Saleem AB, Hassabis D, Spiers HJ. Hippocampal place cells construct reward related
871 sequences through unexplored space. *Elife*. 2015; 4:e06063.
- 872 **Pachitariu M**, Steinmetz NA, Kadir SN, Carandini M, Harris KD. Fast and accurate spike sorting of high-channel
873 count probes with KiloSort. In: Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R, editors. *Advances in Neural*
874 *Information Processing Systems*, vol. 29; 2016. .
- 875 **Penttonen M**, Kamondi A, Sik A, Acsády L, Buzsáki G. Feed-forward and feed-back activation of the dentate
876 gyrus in vivo during dentate spikes and sharp wave bursts. *Hippocampus*. 1997; 7(4):437–450.
- 877 **Ragan T**, Kadiri LR, Venkataraju KU, Bahlmann K, Sutin J, Taranda J, Arganda-Carreras I, Kim Y, Seung HS, Osten P.
878 Serial two-photon tomography for automated ex vivo mouse brain imaging. *Nat Methods*. 2012; 9(3):255–8.
879 doi: [10.1038/nmeth.1854](https://doi.org/10.1038/nmeth.1854).
- 880 **Rajaseethupathy P**, Sankaran S, Marshel JH, Kim CK, Ferenczi E, Lee SY, Berndt A, Ramakrishnan C, Jaffe A, Lo M,
881 Liston C, Deisseroth K. Projections from neocortex mediate top-down control of memory retrieval. *Nature*.
882 2015; 526(7575):653–659.
- 883 **Raposo D**, Kaufman MT, Churchland AK. A category-free neural population supports evolving demands during
884 decision-making. *Nature neuroscience*. 2014; 17(12):1784–1792.
- 885 **Rossant C**, Winter O, Hunter M, Huntenburg J, Faulkner M, Wells M, Steinmetz N, Harris K, Bonacchi N, Alyx;
886 2021. <https://github.com/cortex-lab/alyx>.
- 887 **Roth MM**, Dahmen JC, Muir DR, Imhof F, Martini FJ, Hofer SB. Thalamic nuclei convey diverse contextual infor-
888 mation to layer 1 of visual cortex. *Nat Neurosci*. 2016; 19(2):299–307.
- 889 **Saalman YB**, Kastner S. Cognitive and perceptual functions of the visual thalamus. *Neuron*. 2011; 71(2):209–
890 223.
- 891 **Seabold S**, Perktold J. statsmodels: Econometric and statistical modeling with python. In: *9th Python in Science*
892 *Conference*; 2010. .
- 893 **Senzai Y**, Buzsáki G. Physiological Properties and Behavioral Correlates of Hippocampal Granule Cells and
894 Mossy Cells. *Neuron*. 2017; 93(3):691–704.e5. doi: [10.1016/j.neuron.2016.12.011](https://doi.org/10.1016/j.neuron.2016.12.011).
- 895 **Siegle JH**, Jia X, Durand S, Gale S, Bennett C, Graddis N, Heller G, Ramirez TK, Choi H, Luviano JA, Groblewski PA,
896 Ahmed R, Arkhipov A, Bernard A, Billeh YN, Brown D, Buice MA, Cain N, Caldejon S, Casal L, et al. Survey of
897 spiking in the mouse visual system reveals functional hierarchy. *Nature*. 2021; 592(7852):86–92.
- 898 **Silva D**, Feng T, Foster DJ. Trajectory events across hippocampal place cells require previous experience. *Nature*
899 *neuroscience*. 2015; 18(12):1772–1779.
- 900 **Steinmetz NA**, Aydin C, Lebedeva A, Okun M, Pachitariu M, Bauza M, Beau M, Bhagat J, Böhm C, Broux M, Chen
901 S, Colonell J, Gardner RJ, Karsh B, Kloosterman F, Kostadinov D, Mora-Lopez C, O'Callaghan J, Park J, Putzeys
902 J, et al. Neuropixels 2.0: A miniaturized high-density probe for stable, long-term brain recordings. *Science*.
903 2021; 372(6539):eabf4588.
- 904 **Steinmetz NA**, Zatka-Haas P, Carandini M, Harris KD. Distributed coding of choice, action and engagement
905 across the mouse brain. *Nature*. 2019 Dec; 576(7786):266–273.
- 906 **The International Brain Laboratory**, iblvideo; 2021. <https://github.com/int-brain-lab/iblvideo>.
- 907 **The International Brain Laboratory**, pykilosort; 2021. <https://github.com/int-brain-lab/pykilosort>.
- 908 **The International Brain Laboratory**. Appendix 1: IBL electrophysiological (Ephys Neuropixels) rig setup in-
909 structions: Hardware and Software. figshare. 2022; doi: [10.6084/m9.figshare.17307077](https://doi.org/10.6084/m9.figshare.17307077).
- 910 **The International Brain Laboratory**. Appendix 2: IBL protocol for electrophysiology recording using Neu-
911 ropixels probe. figshare. 2022; doi: [10.6084/m9.figshare.19697896](https://doi.org/10.6084/m9.figshare.19697896).

- 912 **The International Brain Laboratory.** Appendix 3: IBL protocol for craniotomy surgery. figshare. 2022; doi:
913 [10.6084/m9.figshare.19697827](https://doi.org/10.6084/m9.figshare.19697827).
- 914 **The International Brain Laboratory.** Appendix 4: Protocol for labeling the tip of Neuropixels probes. figshare.
915 2022; doi: [10.6084/m9.figshare.19698130](https://doi.org/10.6084/m9.figshare.19698130).
- 916 **The International Brain Laboratory.** Appendix 5: IBL protocol for perfusion and shipment of brain sample.
917 figshare. 2022; doi: [10.6084/m9.figshare.19698061](https://doi.org/10.6084/m9.figshare.19698061).
- 918 **The International Brain Laboratory.** Appendix 6: IBL protocol for registering the electrode location using
919 LASAGNA. figshare. 2022; doi: [10.6084/m9.figshare.19698166](https://doi.org/10.6084/m9.figshare.19698166).
- 920 **The International Brain Laboratory.** Appendix 7: IBL protocol for mouse brain reconstruction and registra-
921 tion. figshare. 2022; doi: [10.6084/m9.figshare.19698895](https://doi.org/10.6084/m9.figshare.19698895).
- 922 **The International Brain Laboratory,** Aguillon-Rodriguez V, Angelaki D, Bayer H, Bonacchi N, Carandini M,
923 Cazettes F, Chapuis G, Churchland AK, Dan Y, Dewitt E, Faulkner M, Forrest H, Haetzel L, Hausser M, Hofer
924 SB, Hu F, Khanal A, Krasniak C, Laranjeira I, et al. Standardized and reproducible measurement of decision-
925 making in mice. *eLife*. 2021; 10:e63711. doi: [10.7554/eLife.63711](https://doi.org/10.7554/eLife.63711).
- 926 **The International Brain Laboratory,** Banga K, Boussard J, Chapuis G, Faulkner M, Harris K, Huntenburg J,
927 Hurwitz C, Lee HD, Paninski L, Rossant C, Roth N, Steinmetz N, Windolf C, Winter O. Spike sorting pipeline
928 for the International Brain Laboratory. figshare. 2022; doi: [10.6084/m9.figshare.19705522](https://doi.org/10.6084/m9.figshare.19705522).
- 929 **The International Brain Laboratory,** Birman D, Bonacchi N, Buchanan K, Chapuis G, Huntenburg J, Meijer G,
930 Paninski L, Schartner M, Svoboda K, Whiteway M, Wells M, Winter O. Video hardware and software for the
931 International Brain Laboratory. figshare. 2022; doi: [10.6084/m9.figshare.19694452](https://doi.org/10.6084/m9.figshare.19694452).
- 932 **Tolhurst DJ,** Movshon JA, Dean AF. The statistical reliability of signals in single neurons in cat and monkey
933 visual cortex. *Vision Res*. 1983; 23(8):775–85. [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=6623937)
934 [db=PubMed&dopt=Citation&list_uids=6623937](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=6623937), doi: 0042-6989(83)90200-6 [pii].
- 935 **Tomar S.** Converting video formats with FFmpeg. *Linux Journal*. 2006; 2006(146):10.
- 936 **Tsui J,** Schwartz N, Ruthazer ES. A developmental sensitive period for spike timing-dependent plasticity in the
937 retinotectal projection. *Frontiers in synaptic neuroscience*. 2010; 2:13.
- 938 **Turk-Browne NB.** The hippocampus as a visual area organized by space and time: A spatiotemporal similarity
939 hypothesis. *Vision research*. 2019; 165:123–130.
- 940 **Urai AE,** Doiron B, Leifer AM, Churchland AK. Large-scale neural recordings call for new insights to link brain
941 and behavior. *Nature Neuroscience*. 2022 Jan; 25(1):11–19. <https://doi.org/10.1038/s41593-021-00980-9>, doi:
942 [10.1038/s41593-021-00980-9](https://doi.org/10.1038/s41593-021-00980-9).
- 943 **Voelkl B,** Altman NS, Forsman A, Forstmeier W, Gurevitch J, Jaric I, Karp NA, Kas MJ, Schielzeth H, Van de Castele
944 T, Würbel H. Reproducibility of animal research in light of biological variation. *Nature Reviews Neuroscience*.
945 2020; 21(7):384–393. doi: [10.1038/s41583-020-0313-3](https://doi.org/10.1038/s41583-020-0313-3).
- 946 **de Vries SEJ,** Lecoq JA, Buice MA, Groblewski PA, Ocker GK, Oliver M, Feng D, Cain N, Ledochowitsch P, Millman
947 D, Roll K, Garrett M, Keenan T, Kuan L, Mihalas S, Olsen S, Thompson C, Wakeman W, Waters J, Williams D,
948 et al. A large-scale standardized physiological survey reveals functional organization of the mouse visual
949 cortex. *Nature Neuroscience*. 2020; 23(1):138–151. doi: [10.1038/s41593-019-0550-9](https://doi.org/10.1038/s41593-019-0550-9).
- 950 **Waaga T,** Agmon H, Normand VA, Nagelhus A, Gardner RJ, Moser MB, Moser EI, Burak Y. Grid-cell modules
951 remain coordinated when neural activity is dissociated from external sensory cues. *Neuron*. 2022; .
- 952 **Wang Q,** Ding SL, Li Y, Royall J, Feng D, Lesnar P, Graddis N, Naeemi M, Facer B, Ho A, Dolbeare T, Blanchard
953 B, Dee N, Wakeman W, Hirokawa KE, Szafer A, Sunkin SM, Oh SW, Bernard A, Phillips JW, et al. The Allen
954 Mouse Brain Common Coordinate Framework: A 3D Reference Atlas. *Cell*. 2020; 181(4):936–953.e20. doi:
955 [10.1016/j.cell.2020.04.007](https://doi.org/10.1016/j.cell.2020.04.007).
- 956 **West SJ,** BrainRegister; 2021. <https://github.com/stevenjwest/brainregister>.
- 957 **Zhang LI,** Tao HW, Holt CE, Harris WA, Poo Mm. A critical window for cooperation and competition among
958 developing retinotectal synapses. *Nature*. 1998; 395(6697):37–44.

959 **Supplementary figures**

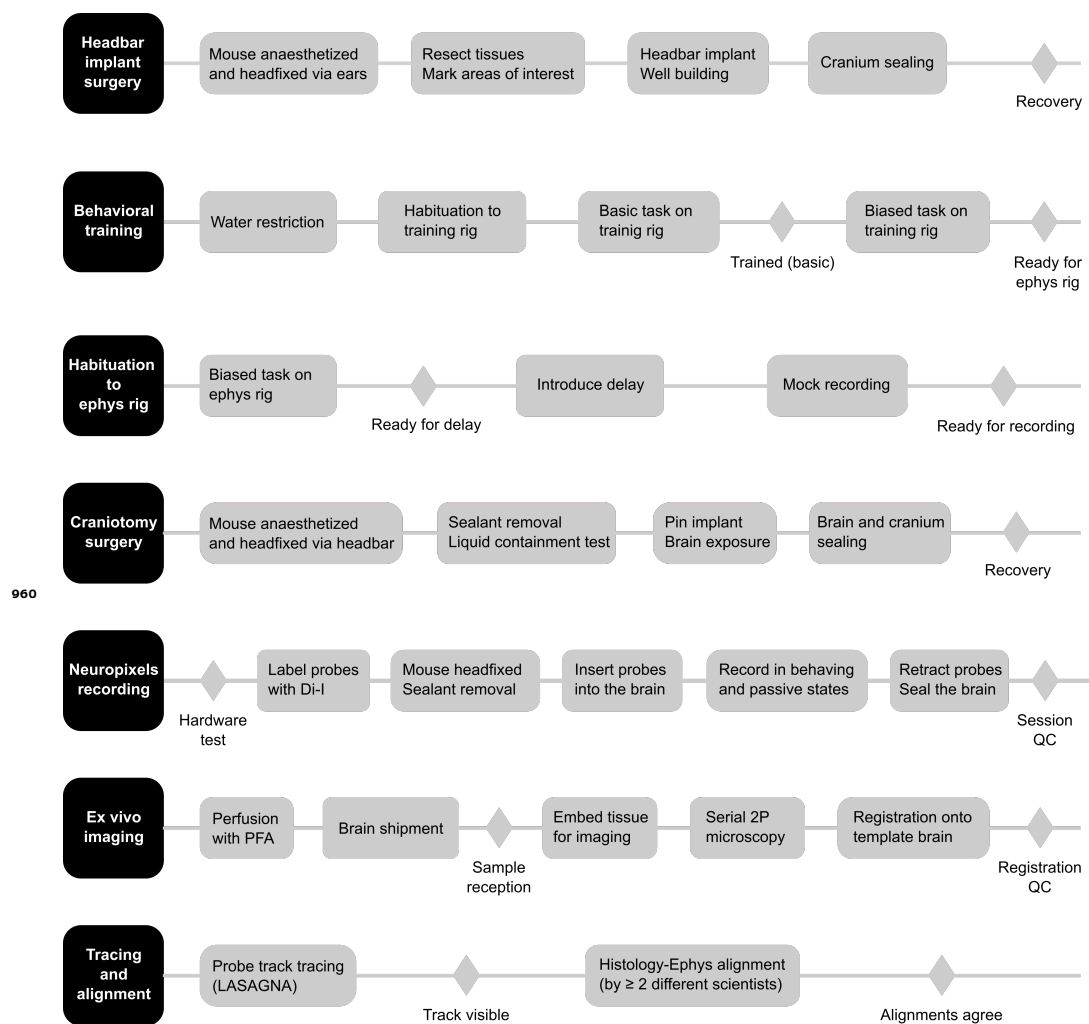


Figure 1-Figure supplement 1. Detailed experimental pipeline for the Neuropixels experiment. The experiment follows the steps indicated in the left-hand black squares in chronological order from top to bottom. Within each, actions are undertaken from left to right; diamond markers indicate points of control.

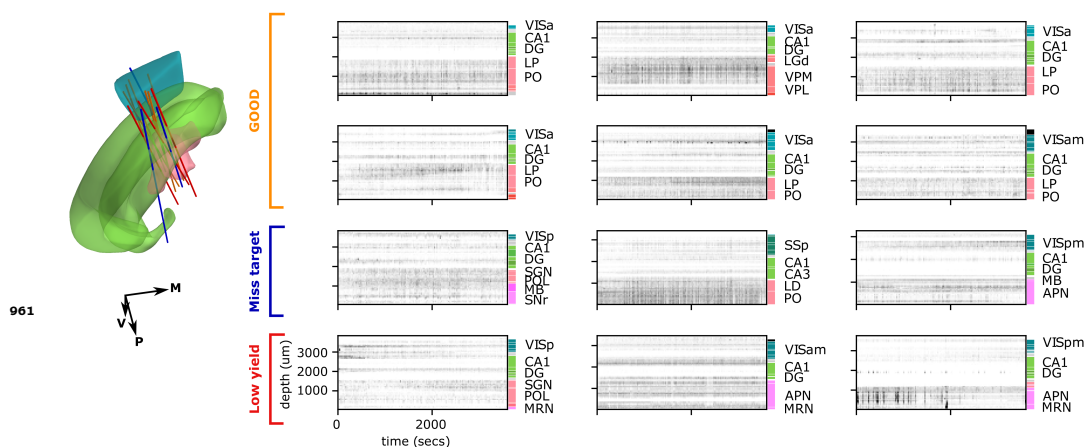


Figure 1-Figure supplement 2. Spiking activity qualitatively appears as heterogeneous across recordings. Example raster plots of neural activity recorded from the repeated site in N=12 mice. The raster plots in the first top two rows originate from sessions marked as being of good quality. The middle and bottom rows are raster plots from recordings that were excluded, based either on the probe misplacement, or the low number of detected units.

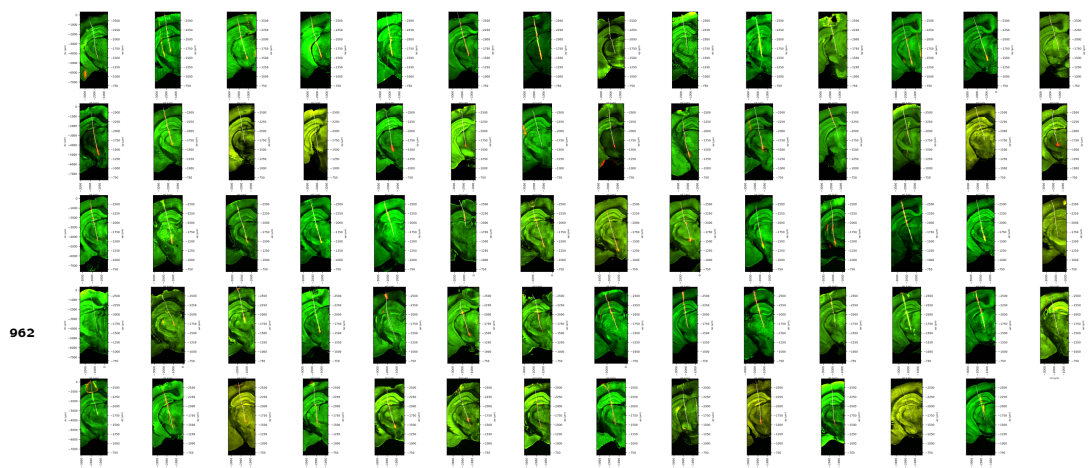


Figure 2-Figure supplement 1. Plots of all subjects with a repeated site insertion that were included in analysis of probe placement. Coronal tilted slices are made along the linearly interpolated best-fit to the histology insertion, shown through the raw histology (green: auto-fluorescence data for image registration; red: CM-Dil fluorescence signal marking probe tracks). Traced probe tracks are highlighted in white.

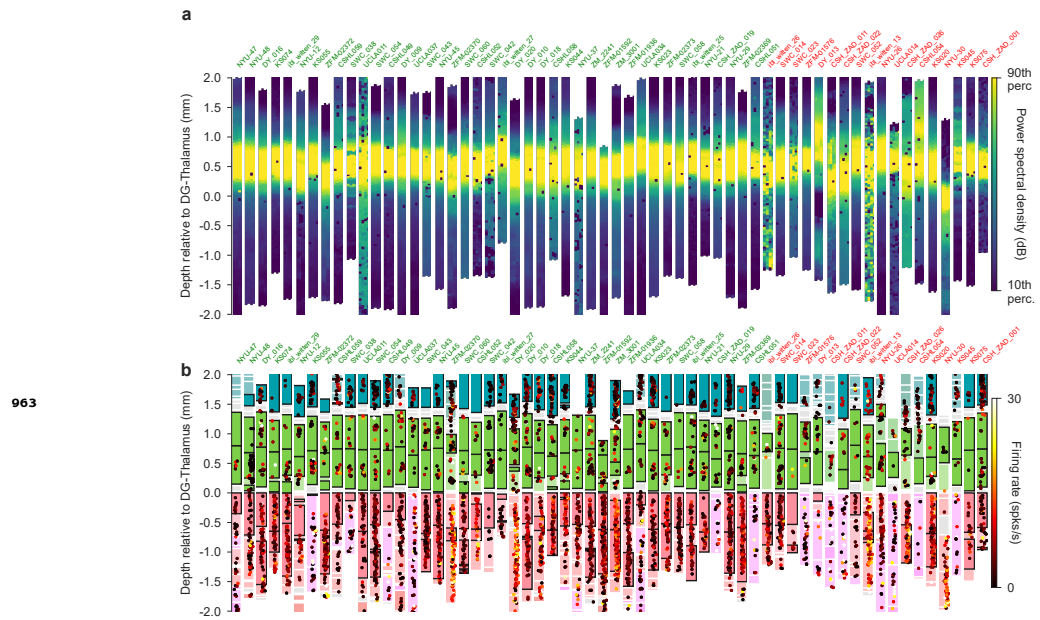


Figure 3-Figure supplement 1. Recordings that failed quality control were often visible outliers. **a**, Power spectral density between 20 and 80 Hz of all insertions, including those that failed to meet quality criteria. Recordings are labelled with the subject name above them; names in green passed quality control whereas names in red did not. **b**, Plots as in **a** but with firing rates of single neurons according to the depth at which they were recorded.

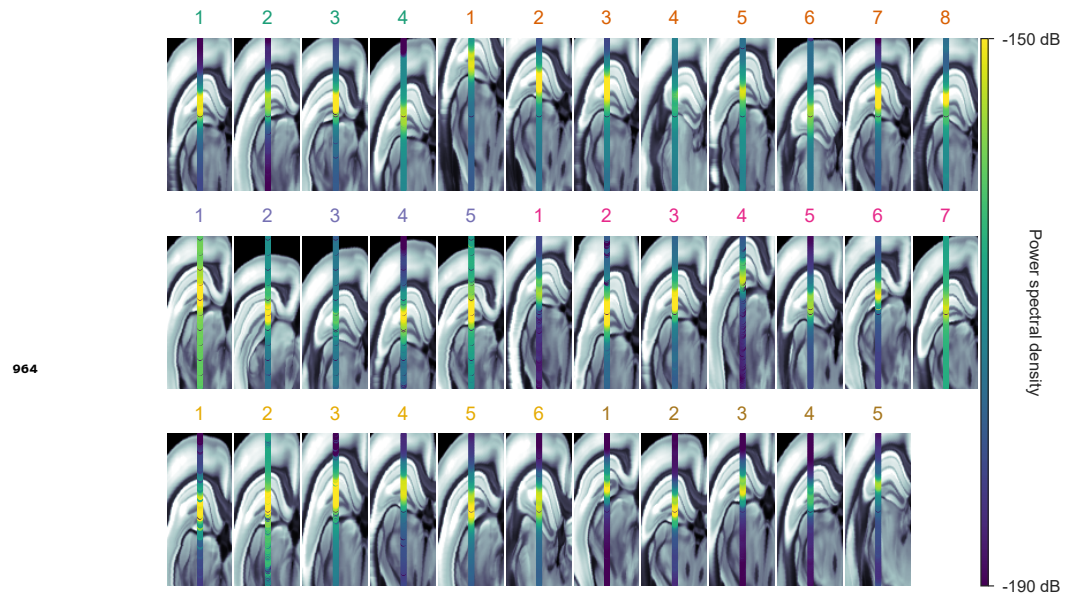


Figure 3-Figure supplement 2. Power spectral density between 20 and 80 Hz recorded along each probe shown in figure 3 overlaid on a coronal slice through brain. Each coronal slice has been rotated such that the probe lies along the vertical axis.

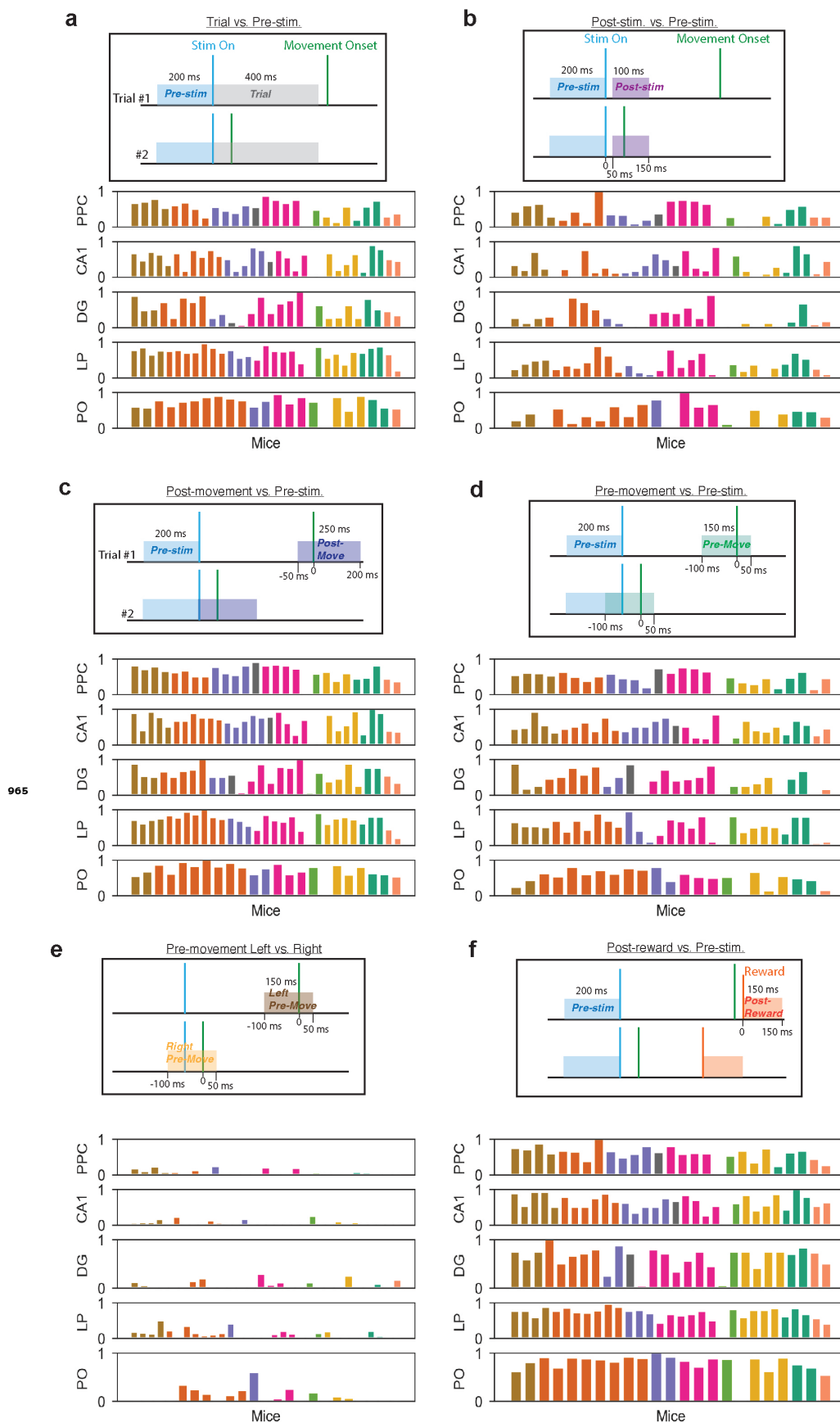


Figure 5-Figure supplement 1. (a)-(g) Schematics of six different tests performed (in addition to the test in Figure 5b) for finding task-modulated neurons. The two example trials show potential caveats of using each method; for instance, in **a**, the trial period may or may not include movement, depending on the reaction time in each trial. Below each schematic, the proportion of task-modulated neurons for the test is shown, across mice and brain regions, colored by lab ID.

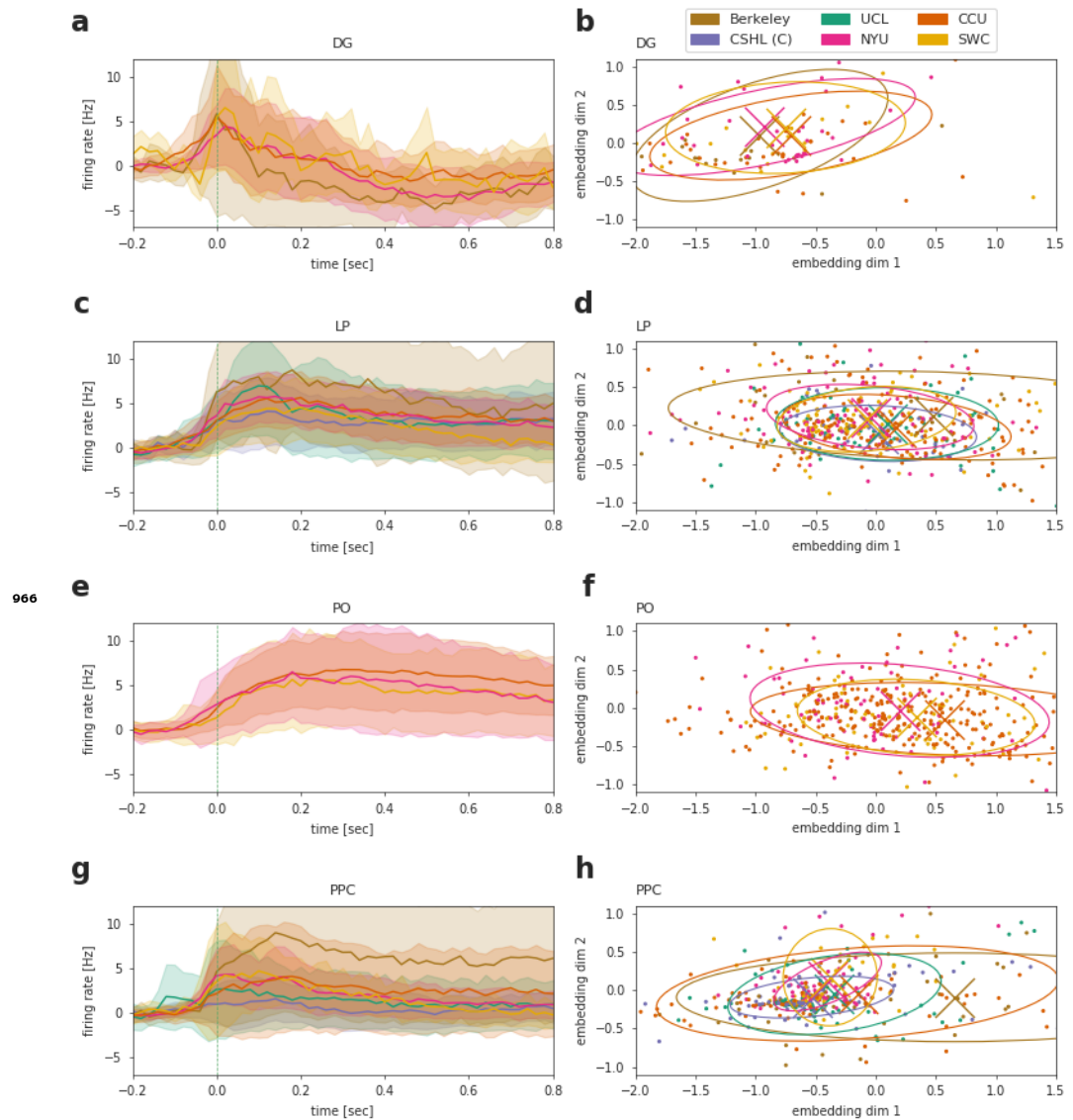


Figure 6-Figure supplement 1. Same as 6(e,f), for the remaining regions. Note that only Berkeley lab in region PPC differs significantly from the mean of all labs.

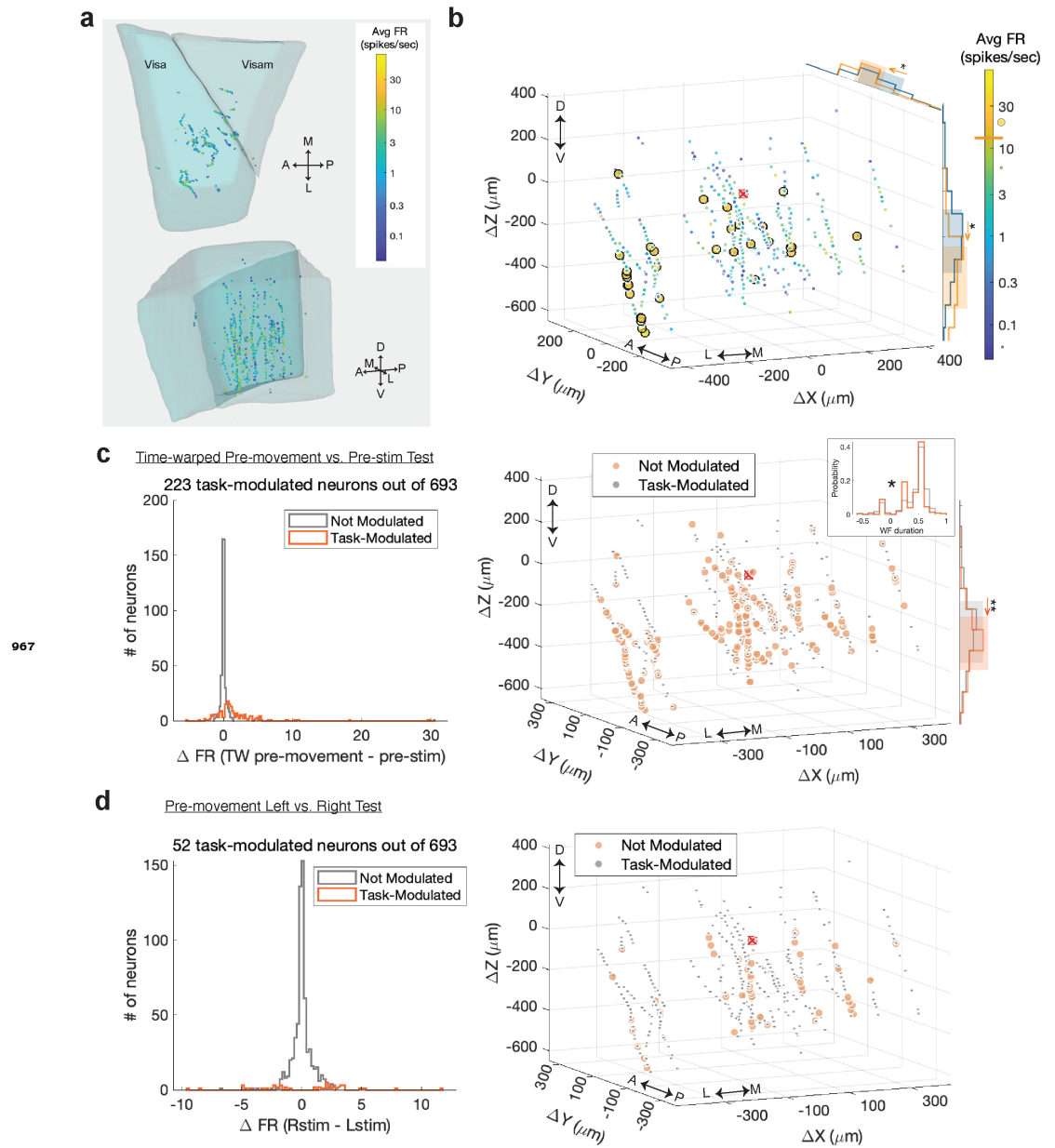


Figure 7-Figure supplement 1. High-firing and task-modulated PPC neurons are located in deeper layers than other PPC neurons. **(a-d)** Similar to Figure 7 but for PPC.

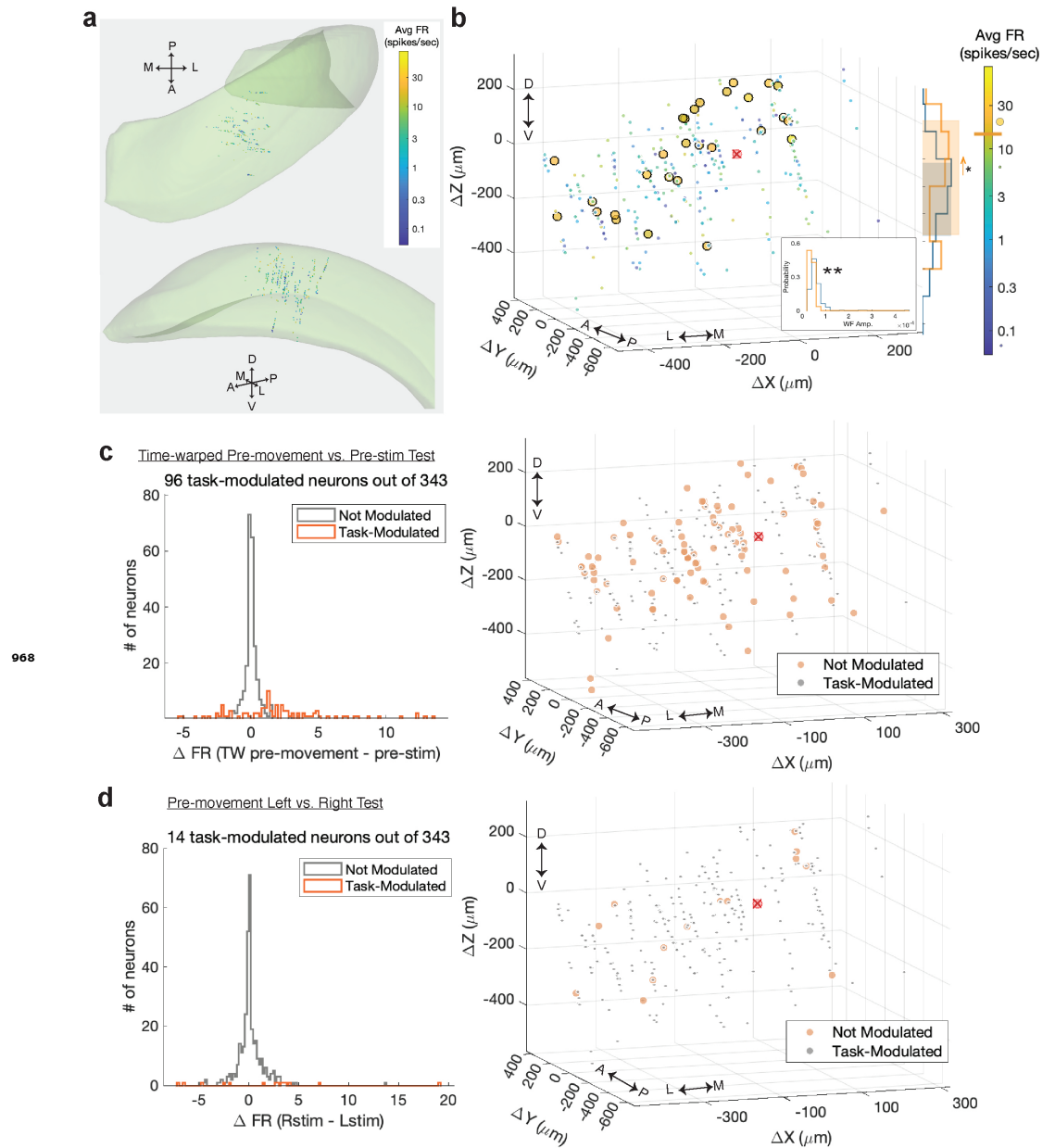


Figure 7-Figure supplement 2. High-firing, but not task-modulated, CA1 neurons are positioned more dorsally and have lower spike amplitudes than other CA1 neurons. **(a-d)** Similar to Figure 7 but for CA1.

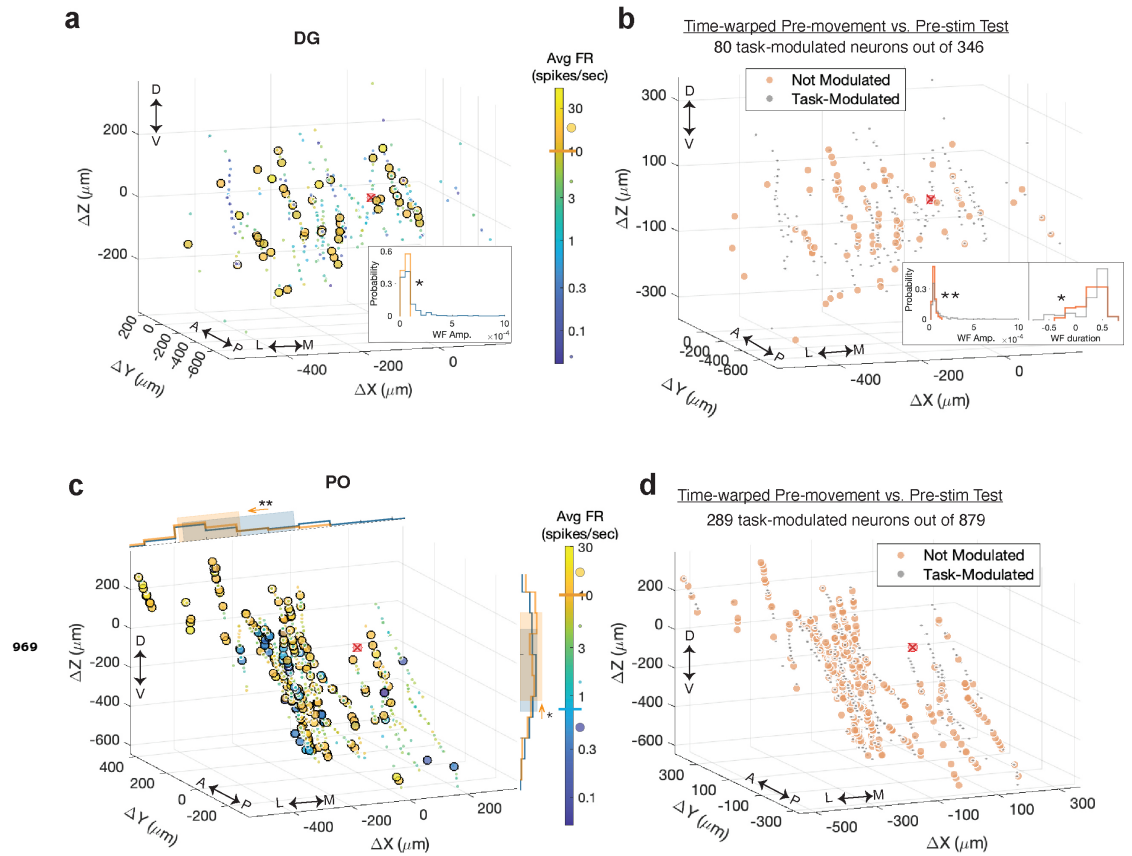


Figure 7-Figure supplement 3. Spatial positions and spike characteristics of outlier and task-modulated neurons in DG and PO are different from other neurons. **(a)** Spatial positions of DG neurons plotted as distance from the planned target center of mass, indicated with the red x. From comparisons of spatial position and waveform features, histogram of only those that were significantly different between the outliers (yellow) and regular neurons (blue) are shown: here, high-firing neurons have smaller waveform amplitudes. **(b)** Spatial positions of task-modulated and non-modulated DG neurons (using the time-warped pre-movement test) with the histogram of significant features shown (here, waveform amplitude and duration). For some other task-modulation tests (not shown), spatial positions of DG neurons were also significantly different. **(c-d)** Same as **a-b** but for PO neurons. In **c**, outliers included high and low firing neurons, making up 209 out of 879 neurons (44 of which are low firing). Shaded areas indicate the 20th and 80th percentiles of the neuron's spatial positions.

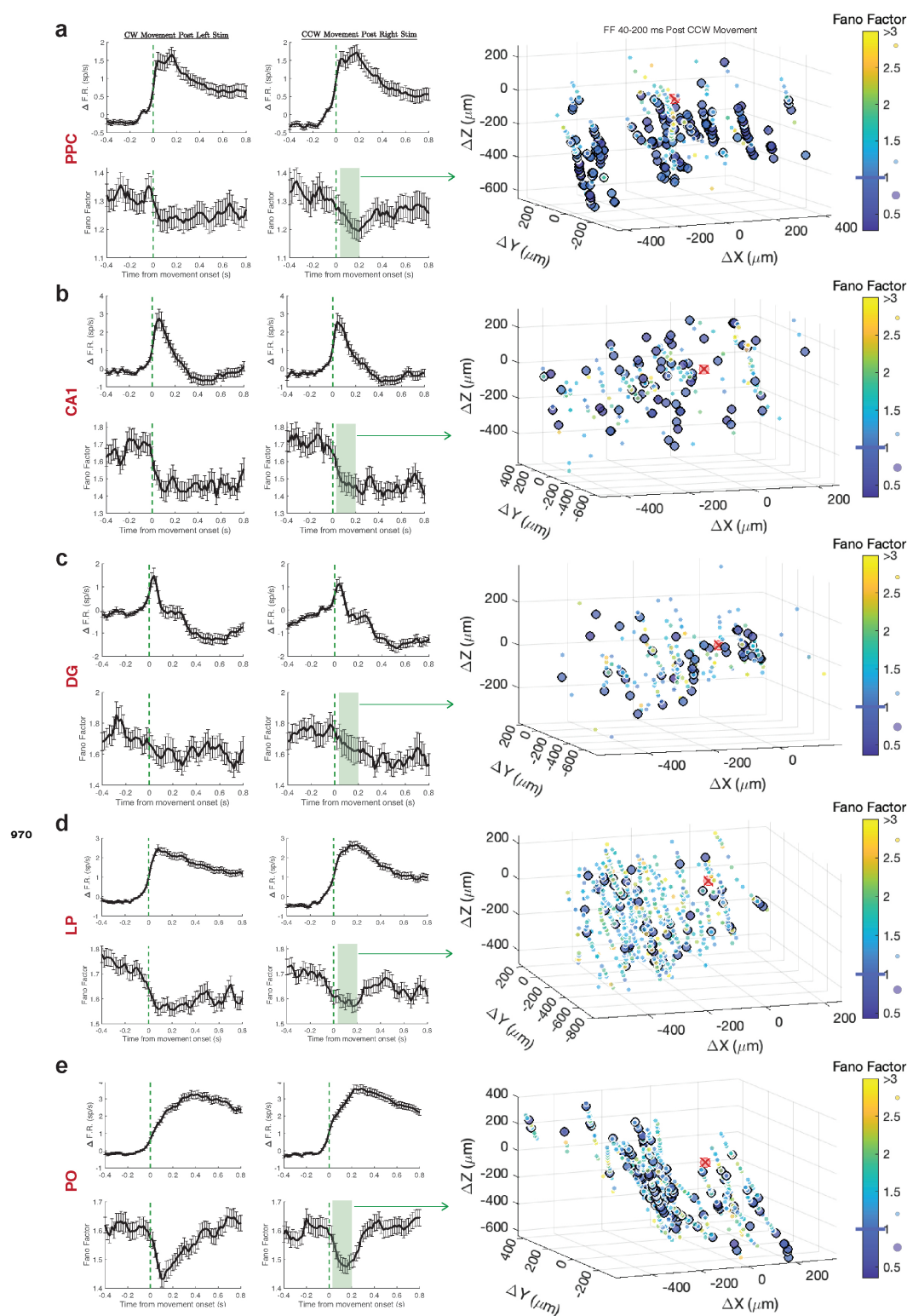


Figure 7-Figure supplement 4. Time-course and spatial position of neuronal Fano Factors. **(a)** *Left column:* Change in firing rate (top) and Fano Factor (bottom) averaged over all PPC neurons when aligned to movement onset after presentation of left or right full-contrast stimuli (correct trials only; Fano Factor calculation limited to neurons with a session-averaged firing rate >1 sp/sec). Error bars: standard error means between neurons. *Right column:* Neuronal Fano Factors (averaged over 40-200 ms post movement onset after right-side full-contrast stimuli) and their spatial positions. Larger circles indicate neurons with Fano Factor <1. **(b-e)** Same as **a** for CA1, DG, LP, and PO. Spatial position between high vs. low Fano Factor neurons was only significantly different in PPC (deeper neurons have lower Fano Factors) possibly due to higher drift in the activity of neurons closer to the surface over long recordings, from drying of the craniotomy. In the thalamus, spike characteristics between high and low Fano Factor neurons were significantly different (not shown).

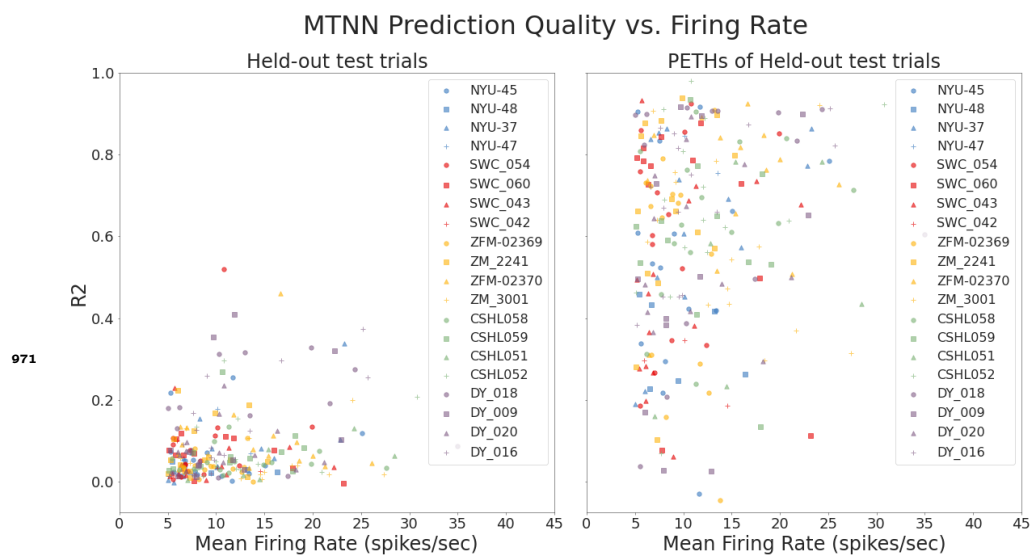


Figure 9–Figure supplement 1. For each unit in each session, we plot the MTNN prediction quality on held-out test trials against the firing rate of the unit averaged over the test trials. Each lab/session is colored/shaped differently. R^2 values on concatenations of the held-out test trials are shown on the left, and those on PETHs of the held-out test trials on the right.

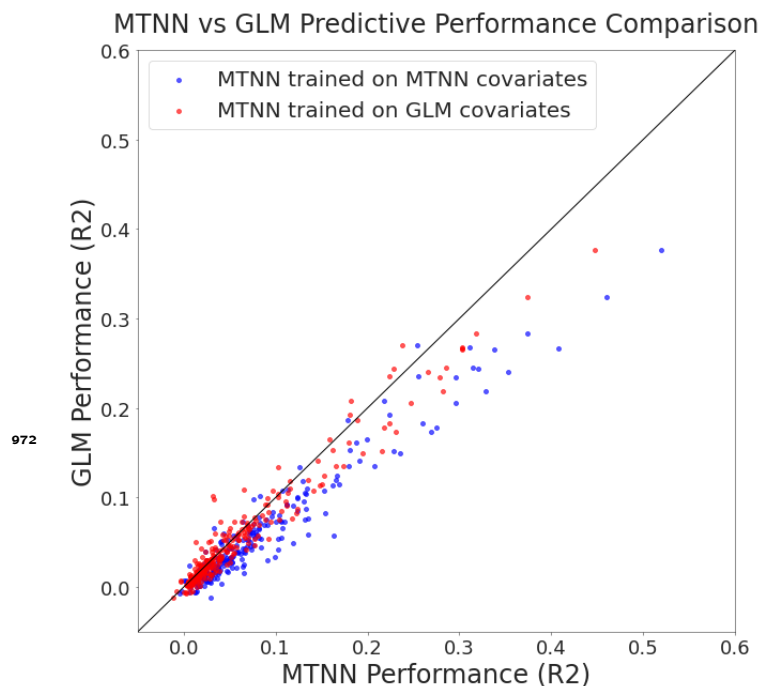


Figure 9–Figure supplement 2. MTNN and GLMs performs similarly on predicting the firing rates of held-out trials when trained on a reduced set of covariates, which includes stimulus onset time, stimulus side and contrast, feedback time and type, first movement onset time, wheel velocity, and mouse’s prior. MTNN trained on the full set of covariates in Table 2 outperforms the MTNN/GLMs trained on the reduced covariate set.

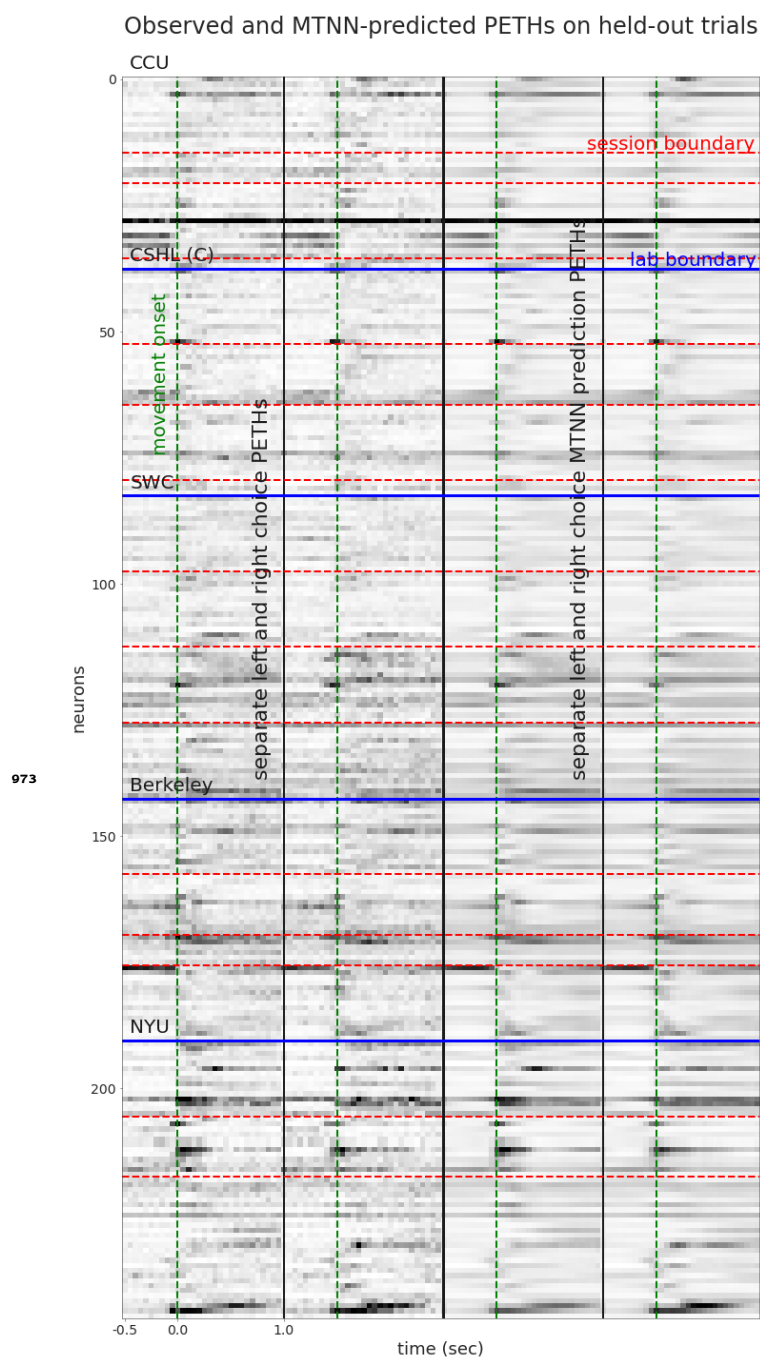


Figure 9-Figure supplement 3. The left half shows for each neuron the trial averaged activity for left choice trials and next to it right choice trials. The vertical green lines show the first movement onset. The horizontal red lines separate recording sessions while the blue lines separate labs. The right half of each of these images shows the MTNN prediction of the left half. The trial-averaged MTNN predictions for held-out test trials captures visible modulations in the PETHs.

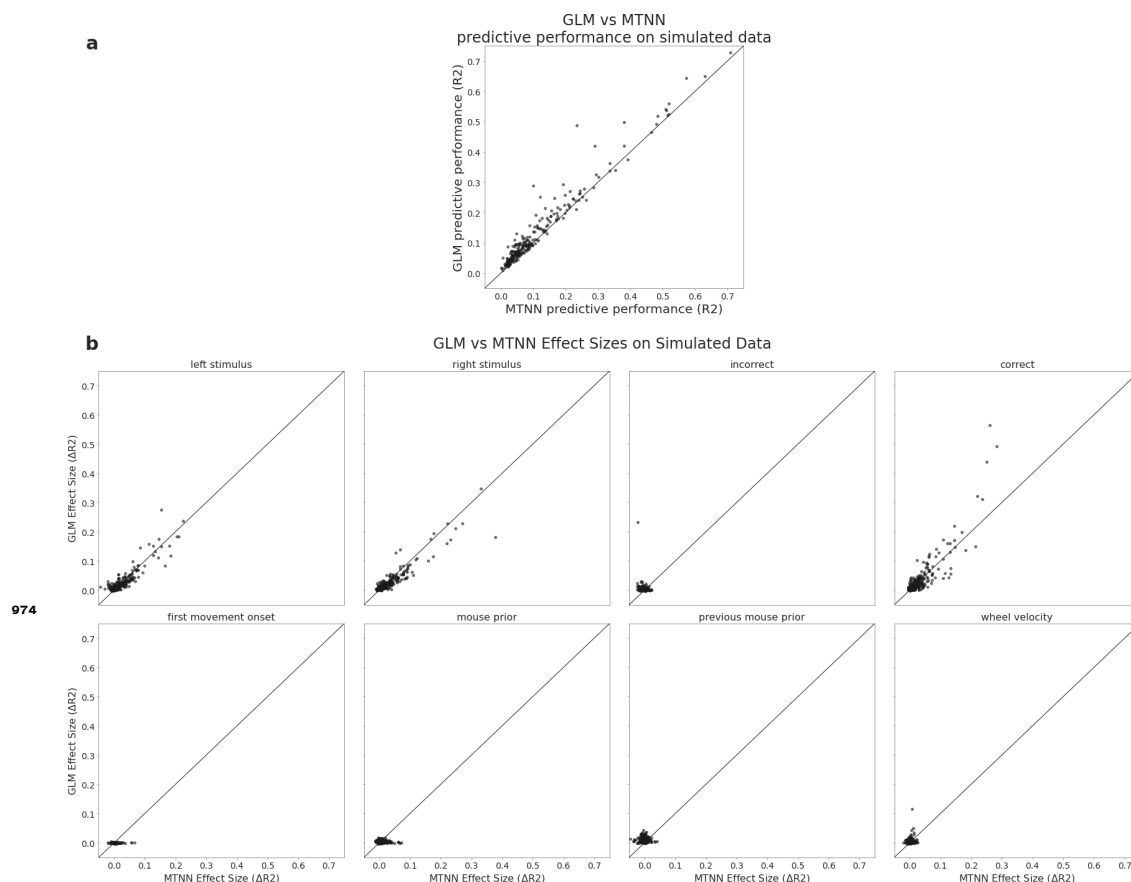


Figure 10-Figure supplement 1. To verify that the MTNN leave-one-out analysis is sensitive enough to capture effect sizes, we simulate data from GLMs and compare the effect sizes estimated by the MTNN and GLM leave-one-out analyses. We first fit GLMs to the same set of sessions that are used for the MTNN effect size analysis and then use the inferred GLM kernels to simulate data. **(a)** We show the scatterplot of the GLM and MTNN predictive performance on held-out test data, where each dot represents the predictive performance for one neural unit. The MTNN prediction quality is comparable to that of GLMs. **(b)** We run GLM and MTNN leave-one-out analyses and compare the estimated effect sizes for 6 covariates. The effect sizes estimated by the MTNN and GLM leave-one-out analyses are comparable.

Pairwise scatterplots of MTNN single-covariate effect sizes

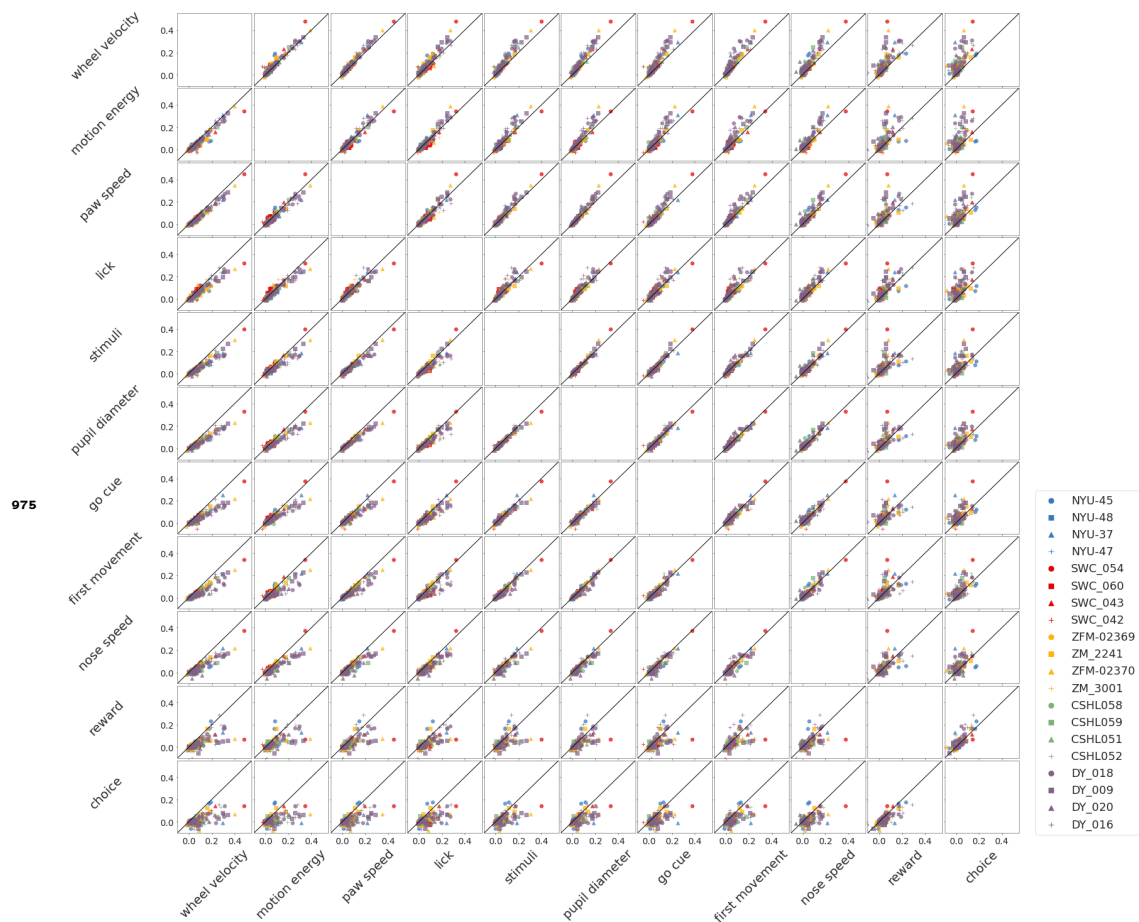


Figure 10-Figure supplement 2. We plot pairwise scatterplots of MTNN single-covariate effect sizes. Each dot represents the effect sizes of one neural unit and is colored by lab. There is no outlier lab. The effect sizes are highly correlated.