

On the use of calcium deconvolution algorithms in practical contexts

Mathew H. Evans^{1,2}, Rasmus S. Petersen² & Mark D. Humphries^{1,2}

¹School of Psychology, University of Nottingham, UK

²Faculty of Biology, Medicine and Health, University of Manchester, UK

November 20, 2019

Abstract

Calcium imaging is a powerful tool for capturing the simultaneous activity of large populations of neurons. Studies using it to address scientific questions of population dynamics and coding often use the raw time-series of changes in calcium fluorescence at the soma. But somatic calcium traces are both contaminated with multiple noise sources and are non-linearly related to spiking. A suite of processing methods are available to recover spike-evoked events from the raw calcium, from simple deconvolution to inferring the spikes themselves. Here we determine the extent to which our inferences of neural population activity, correlations, and coding depend on our choice of processing method. To this end, we processed calcium imaging data obtained from barrel cortex during a pole-detection task using eight commonly used methods. We show that a substantial fraction of the processing methods fail to recover simple features of population activity in barrel cortex already established by electrophysiological recordings. Raw calcium time-series contain an order of magnitude more neurons tuned to features of the pole task; yet there is also qualitative disagreement between deconvolution methods on which neurons are tuned to the task. Finally, we show that raw and processed calcium time-series qualitatively disagree on the structure of correlations within the population and the dimensionality of its joint activity. Collectively, our results show that properties of neural activity, correlations, and coding inferred from calcium imaging are highly sensitive to the choice of method for recovering spike-evoked events. We suggest that quantitative results obtained from population calcium-imaging be verified across multiple processed forms of the calcium time-series.

1 Introduction

Calcium imaging is a wonderful tool for high yield recordings of large neural populations (Harris et al., 2016; Stringer et al., 2019; Ahrens et al., 2013; Portugues et al., 2014). Many pipelines are available for moving from pixel intensity across frames of video to a time-series of calcium fluorescence in the soma of identified neurons (Mukamel et al., 2009; Vogelstein et al., 2010; Kaifosh et al., 2014; Pachitariu et al., 2016; Deneux et al., 2016; Pnevmatikakis et al., 2016; Friedrich et al., 2017; Keemink et al., 2018; Giovannucci et al., 2019). As somatic calcium is proportional to the release of spikes, so we wish to use these fluorescence time-series as a proxy for spiking activity in large, identified populations of neurons. But raw calcium fluorescence is nonlinearly related to spiking, and contains noise from a range of sources.

These issues have inspired a wide range of deconvolution algorithms (Theis et al., 2016; Berens et al., 2018; Stringer and Pachitariu, 2018), which attempt to turn raw somatic calcium into something more closely approximating spikes. We address here the question facing any systems neuroscientist using calcium imaging: do we use the raw calcium, or attempt to clean it up? Thus our aim is to understand if our choice matters: to what extent do our inferences about neural activity, correlations, and coding depend on our choice of raw or deconvolved calcium time-series.

Deconvolution algorithms themselves range in complexity from simple deconvolution with a fixed kernel of the calcium response (Yaksi and Friedrich, 2006), through detecting spike-evoked calcium events (Jewell and Witten, 2018; Pachitariu et al., 2016), to directly inferring spike times (Vogelstein et al., 2010; Lütcke et al., 2013; Deneux et al., 2016). This continuum of options raise the further question of the extent to which we should process the raw calcium signals.

We proceed here in two stages. In order to use deconvolution algorithms, the experimenter needs to choose their parameters. An open question is whether it is worth taking this extra step: how good can these algorithms be in principle, and how sensitive their results are to the choice of parameter values. We thus first evaluate qualitatively different deconvolution algorithms, by optimising their parameters against ground truth data with known spikes. With our understanding of their parameters in hand, we then turn to our main question, by analysing a large-scale population recording from the barrel cortex of a mouse performing a whisker-based decision task. We compare estimates of population coding and correlations obtained using either raw calcium signals, or a range of time-series derived from those calcium signals, covering simple deconvolution, event detection, and spikes.

A substantial fraction of the methods used here fail to recover basic features of population activity in barrel cortex established from electrophysiology. The inferences we draw about coding qualitatively differ between raw and deconvolved calcium signals. In particular, coding analyses based on raw calcium signals detect an order of magnitude more neurons tuned to task features. Yet there is also qualitative disagreement between deconvolution methods on which neurons are tuned. The inferences we draw about correlations between neurons do not distinguish between raw and deconvolved calcium signals, but can qualitatively differ between deconvolution methods. Our results thus suggest care is needed in drawing inferences from population recordings of somatic calcium, and that one solution is to replicate all results in both raw and deconvolved calcium signals.

2 Results

2.1 Performance of deconvolution algorithms on ground-truth data-sets

We select here three deconvolution algorithms that infer discrete spike-like events, each an example of the state of the art in qualitatively different approaches to the problem: Suite2p (Pachitariu et al., 2016), a peeling algorithm that matches a scalable kernel to the calcium signal to detect spike-triggered calcium events; LZero (Jewell and Witten, 2018), a change-point detection algorithm, which finds as events the step-like changes in the calcium signal that imply spikes; and MLspike (Deneux et al., 2016), a forward model, which fits an explicit model of the spike-to-calcium dynamics in order to find spike-evoked changes in the calcium signal, and returns spike times. We emphasize that these methods were chosen as exemplars of their approaches, and are each innovative takes on the problem; we are not here critiquing individual methods, but using an array of methods to illustrate

81 the problems and decisions facing the experimentalist when using calcium imaging data.

82 We first ask if these deconvolution methods work well in principle, by testing if there
83 exists parameter sets for which they each successfully recover known spike times from
84 calcium traces. We fit the parameters of each method to a data-set of 21 ground-truth
85 recordings (Chen et al., 2013), where the spiking activity of a neuron is recorded simulta-
86 neously with 60 Hz calcium imaging and a high-signal-to-noise cell-attached glass pipette
87 (Figure 1a). To fit the parameters for each recording, we sweep each method’s parameter
88 space to find the parameter value(s) with the best match between the true and inferred
89 spike train.

90 The best-fit parameters depend strongly on how we evaluate the match between true
91 and inferred spikes. The Pearson correlation coefficient between the true and inferred
92 spike train is a common choice (Brown et al., 2004; Paiva et al., 2010; Theis et al., 2016;
93 Reynolds et al., 2018; Berens et al., 2018), typically with both trains convolved with a
94 Gaussian kernel to allow for timing errors. However, we find that choosing parameters to
95 maximise the correlation coefficient can create notable errors. The inferred spike trains
96 from MLSpike have too many spikes on average (mean error over recordings: 31.72%),
97 and the accuracy of recovered firing rates widely varies across recordings (Fig 1b, blue
98 symbols). We attribute these errors to the noisy relationship between the correlation
99 coefficient and the number of inferred spikes (Figure 1c): for many recordings, there is
100 no well-defined maximum coefficient, especially for the amplitude parameter A , so that
101 near-maximum correlation between true and inferred trains is consistent with a wide range
102 of spike counts in the inferred trains. We see the same sensitivity for the event rates from
103 recordings optimised using Suite2p (Figure 1f). If we compare their inferred event rates
104 to true firing rates (Fig 1b), we see Suite2p estimates far more events than spikes (mean
105 error 79.47%) and LZero fewer events than spikes (mean error: -21.14%). These further
106 errors are problematic: there cannot be more spike-driven calcium events than spikes, and
107 LZero’s underestimate is considerably larger than the fraction of frames with two or more
108 spikes ($<2e^{-4}$ % frames).

109 To address the weaknesses of the Pearson correlation coefficient, we instead optimise
110 parameters using the Error Rate metric of Deneux et al. (2016). Error Rate returns a
111 normalised score between 0 for a perfect match between two spike trains, and 1 when all
112 the spikes are missed. This comparison between inferred and true spike trains is most
113 straightforward for algorithms like MLSpike that directly return spike times; for the other
114 algorithms, we use here their event times as inferred spikes, a reasonable choice given the
115 low firing rate and well separated spikes in the ground truth data. Choosing parameters
116 to minimise the Error Rate between the true and inferred spike-trains results in excellent
117 recovery of the true number of spikes for all three deconvolution methods (Fig 1b, green
118 symbols), with mean errors of 12% for Suite2P, 7.3% for MLSpike, and 5% for LZero.
119 As we show in Figure 1e for MLSpike and Figure 1f for Suite2p, the Error Rate has a
120 well-defined minima for almost every recording. Consequently, all deconvolution methods
121 can, in principle, accurately recover the true spike-trains given an appropriate choice of
122 parameters.

123 A potential caveat here is that the ground-truth data are single neurons imaged at a
124 frame-rate of 60Hz, an order of magnitude greater than is typically achievable in popula-
125 tion recordings (Peron et al., 2015a). Such a high frame-rate could allow for more accurate
126 recovery of spikes than is possible in population recordings. To test this, we downsample
127 the ground-truth data to a 7Hz frame-rate, and repeat the parameter sweeps for each
128 deconvolution method applied to each recording. As we show in Figure 1c, optimising pa-
129 rameters using the minimum Error Rate still results in excellent recovery of the true spike

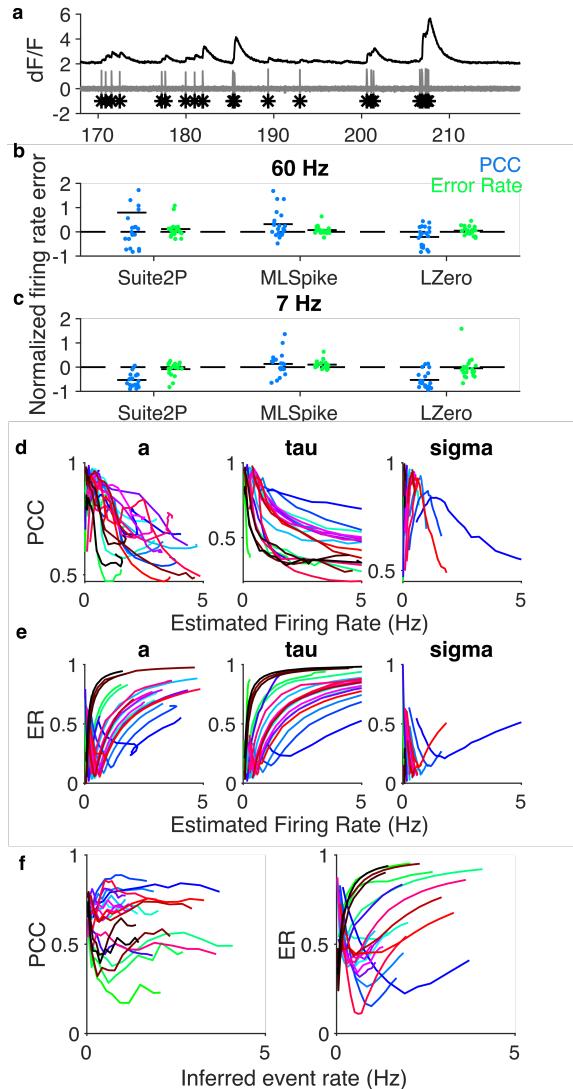


Figure 1: Ground truth data analysis.

- (a) Example simultaneous recording of somatic voltage (grey) and calcium activity (black) imaged at 60Hz. Spikes are marked with asterisks.
- (b) Error in estimating the true firing rate when using optimised parameters, across all three methods. One symbol per recording. We separately plot errors for parameters optimised to maximise the correlation coefficient (PCC), the errors for parameters optimised to minimise the error rate (ER). Horizontal black bars are means. Error is computed relative to the true firing rate: $(Rate_{true} - Rate_{estimated})/Rate_{true}$. For LZero and Suite2p, $Rate_{estimated}$ is computed from event times.
- (c) As for (b), but with the somatic calcium down-sampled to 7Hz before optimising parameters for the deconvolution methods.
- (d) Dependence of MLspike's deconvolution performance on the firing rate of the inferred spike train. For each of ML Spike's free parameters, we plot the correlation coefficient between true and inferred spikes as a function of the firing rate estimated from the inferred spikes. One line per recording. Parameters: A : calcium transient amplitude per spike ($\Delta F/F$); τ calcium decay time constant (s); σ : background (photonic) noise level ($\Delta F/F$)
- (e) as in (d), but using Error Rate between the true and inferred spikes.
- (f) Dependence of Suite2p's deconvolution performance on the firing rate of the inferred event train as a detection threshold parameter is varied. Left: correlation coefficient; right: Error Rate.

130 rate (and interestingly for some recordings reduces the error when using the correlation
131 coefficient). Lower frame-rates need not then be an impediment to using deconvolution
132 methods.

133 **2.2 Parameters optimised on ground-truth are widely distributed and 134 sensitive**

135 What might be an impediment to using deconvolution methods on population recordings is
136 if the best parameter values vary widely between neurons. If so, then parameters optimised
137 for one neuron would generalise poorly to the rest of the population.

138 Figure 2a-b plots the best-fit parameter values for each recording across deconvolution
139 methods and sampling rates. Each method has at least one parameter with substantial
140 variability across recordings, varying by an order of magnitude or more. This suggests
141 that the best parameters for one neuron may not apply to another. The in turn could
142 mean that analysis of population recordings created from a single set of deconvolution
143 parameters would potentially include many aberrant time-series.

144 The problem of between-neuron variation in parameter values would be compensated
145 somewhat if the quality of the inferred spike or event trains is robust to changes in those
146 values. However, we find performance is highly sensitive to changes in some parameters.
147 Figure 2b-c shows that for most recordings the quality of the inferred spike train abruptly
148 worsens with small increases or decreases in the best parameter. Thus using deconvolution
149 algorithms on population recordings comes with the potential issues that parameters can
150 be both sensitive and vary considerably across neurons.

151 **2.3 Deconvolution of population imaging in barrel cortex during a de- 152 cision task**

153 The above results imply that the insights we gain from analysing large-scale population
154 calcium imaging data would depend crucially on which deconvolution methods we use. We
155 now test the extent of disagreement using a large-scale population recording from barrel
156 cortex. Applying 8 different deconvolution methods to the same raw calcium time-series,
157 we compare the resulting statistics of neural activity, properties of neural coding, and the
158 extent and structure of between neuron correlations.

159 The data we use are two-photon calcium imaging time-series from a head-fixed mouse
160 performing a whisker-based two-alternative decision task (Fig. 3a-b), from the study of
161 Peron et al. (2015b). We analyse here a single session with 1552 simultaneously recorded
162 pyramidal neurons in L2/3 of a single barrel in somatosensory cortex, imaged at 7 Hz for
163 just over 56 minutes, giving 23559 frames in total across 335 trials of the task.

164 Our primary goal is to understand how the choices of deconvolving these calcium-
165 imaging data alter the scientific inferences we can draw. As our baseline, we use the
166 “raw” $\Delta F/F$ time-series of changes in calcium indicator fluorescence. We use the above
167 three discrete deconvolution methods to extract spike counts (MLSpike), event occurrence
168 (LZero), or event magnitude (Suite2p) per frame. For comparison, we use Peron et al.
169 (2015b)’s own version of denoised calcium time-series, created using a custom version
170 of the peeling algorithm (Lütcke et al., 2013), a greedy template-fitting algorithm with
171 variable decay time constants across events and neurons, with parameters chosen to result
172 in the same proportion of silent neurons as has been shown previously with unbiased
173 electrophysiology. Given the above-demonstrated dependence of these algorithms on their
174 parameters, we also use Yaksi and Friedrich (2006)’s simple deconvolution of the raw
175 calcium with a fixed kernel of the calcium response to a single spike, whose only free

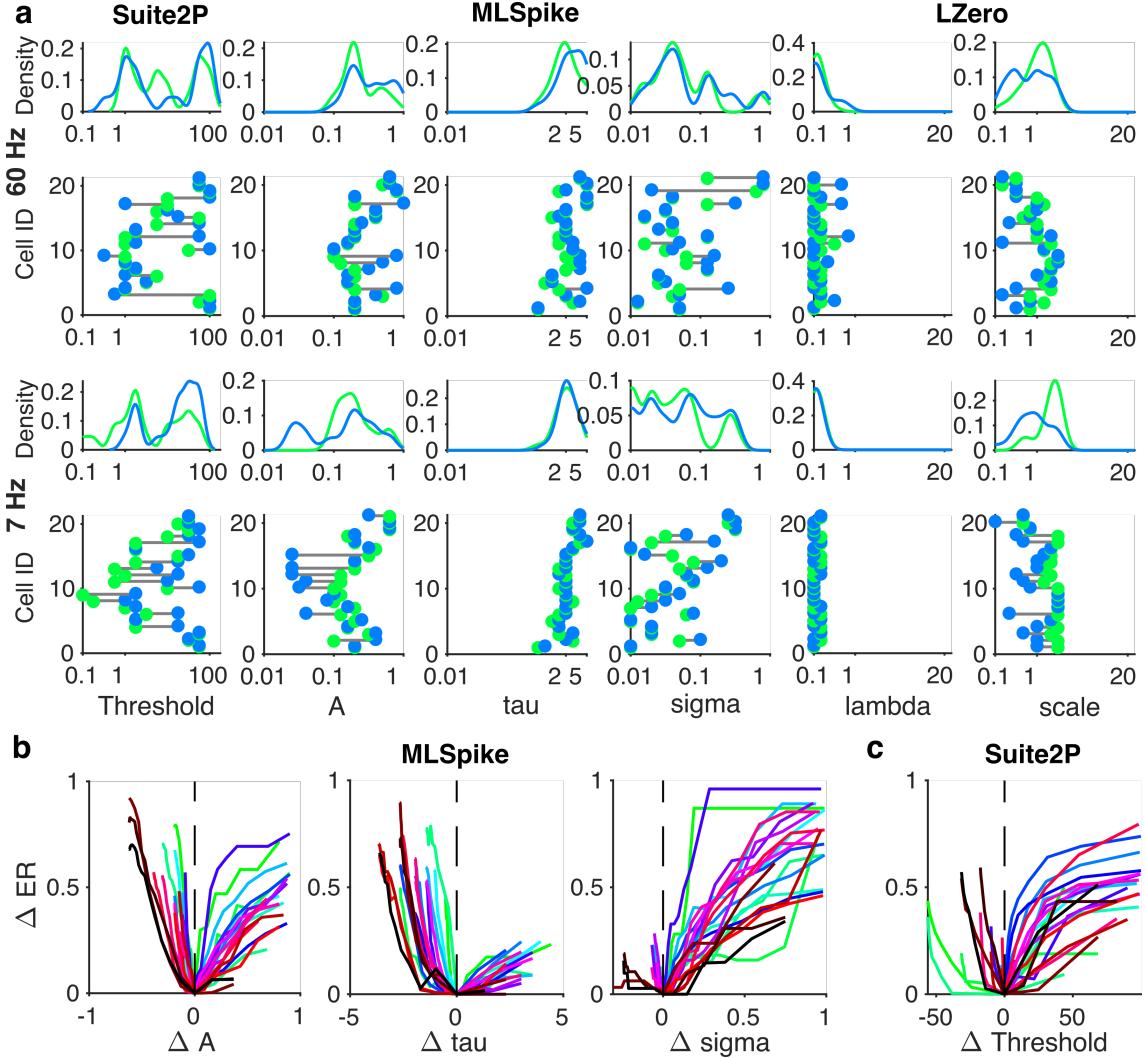


Figure 2: Variation in best-fit spike deconvolution parameters across ground-truth recordings.

(a) Distributions of optimised parameter values across recordings. For each parameter (a column), the bottom panel plots the found parameter values on the x-axis against the recording ID on the y-axis (in an arbitrary but consistent order); the top panel plots the marginal distribution of the parameter value over all neurons. We plot for each recording the optimised parameter value found using correlation coefficient (blue) and Error Rate (green).

(b) As for panel (a), fits to the same ground-truth data down-sampled to 7 Hz.

(c) Change in error rate as a function of the change away from a parameter’s optimum value, for each of ML Spike’s free parameters. One line per recording.

(d) Change in the error rate with change in Suite2p’s threshold value away from its optimum for each recording. One line per recording.

parameters are the response-kernel which are fixed from data. And finally we create smoothed versions of the discrete-deconvolution methods, by convolving their recovered spikes/events with a fixed spike-response kernel. Figure 3c show an example raw calcium time-series for one neuron, and the result of applying each of these 8 processing methods. We thus repeat all analyses on 9 different sets of time-series extracted from the same population recording.

We choose the algorithm parameters as follows. Simple deconvolution (Yaksi and

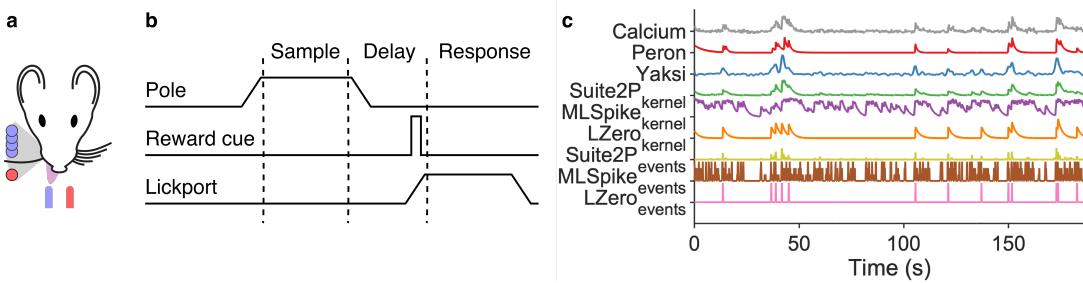


Figure 3: Experimental data from Peron et al. (2015b).

- (a) Schematic of task set-up. A pole was raised within range of the single left-hand whisker; its position, forward (red) or backward (blue) indicated whether reward would be available from the left or right lick-port.
- (b) Schematic of trial events. The pole was raised and lowered during the sample period; a auditory cue indicated the start of the response period.
- (c) All deconvolution methods applied to one raw calcium signal from the same neuron.

183 Friedrich, 2006) involves taking a parameterised kernel of the GCaMP6s response to a sin-
 184 gle spike. For the three discrete deconvolution methods, we choose the modal values of the
 185 best-fit parameters that optimised the Error Rate over the ground-truth recordings. This
 186 seems a reasonably consistent way of comparing methods, by using the most consistently
 187 performing values obtained from comparable data: neurons in the same layer (L2/3) in
 188 the same species (mouse), in another primary sensory area (V1). Most importantly for our
 189 purposes, choosing the modal values means we avoid pathological regions of the parameter
 190 space.

191 2.4 Deconvolution methods disagree on estimates of simple neural statistics

193 We first check how well each approach recovers the basic statistics of neural activity
 194 event rates in L2/3 of barrel cortex. Electrophysiological recordings have shown that the
 195 distribution of firing rates across neurons in a population is consistently long-tailed, and
 196 often log-normal, all across rodent cortex (Wohrer et al., 2013). Cell-attached recordings
 197 of L2/3 neurons in barrel cortex are no different (O'Connor et al., 2010), with median
 198 firing rates less than 1 Hz, and a long right-hand tail of rarer high-firing neurons. We thus
 199 test if the calcium event rates or spike rates from our time-series follow such a distribution.
 200 (Event rates for raw calcium, Peron, Yaksi and the continuous (kernel) versions of the data
 201 was obtained by thresholding the calcium time-series)

202 Figure 4a shows that the raw calcium and two of the discrete deconvolution methods
 203 (Suite2p, LZero) qualitatively match the expected distributions of event rates (median
 204 near zero, long right-hand tails). The Peron time-series also have the correct distribution
 205 of event rates, which is unsurprising as it was tuned to do so. All other methods give wrong
 206 distributions of spike rates (MLSpike) or event rates (all other methods). There is also
 207 little overlap in the distributions of spike rates between the three discrete deconvolution
 208 methods. Applying a kernel to their inferred spikes/events shifts rather than smooths
 209 the firing rate distributions (Suite2P_{kernel}, MLSpike_{kernel}, LZero_{kernel}), suggesting noise
 210 in the deconvolution process is amplified through the additional steps of convolving with
 211 a kernel and thresholding.

212 Cell-attached recordings in barrel cortex have shown that ~26% of L2/3 pyramidal
 213 neurons are silent during a similar pole localisation task, with silence defined as emitting

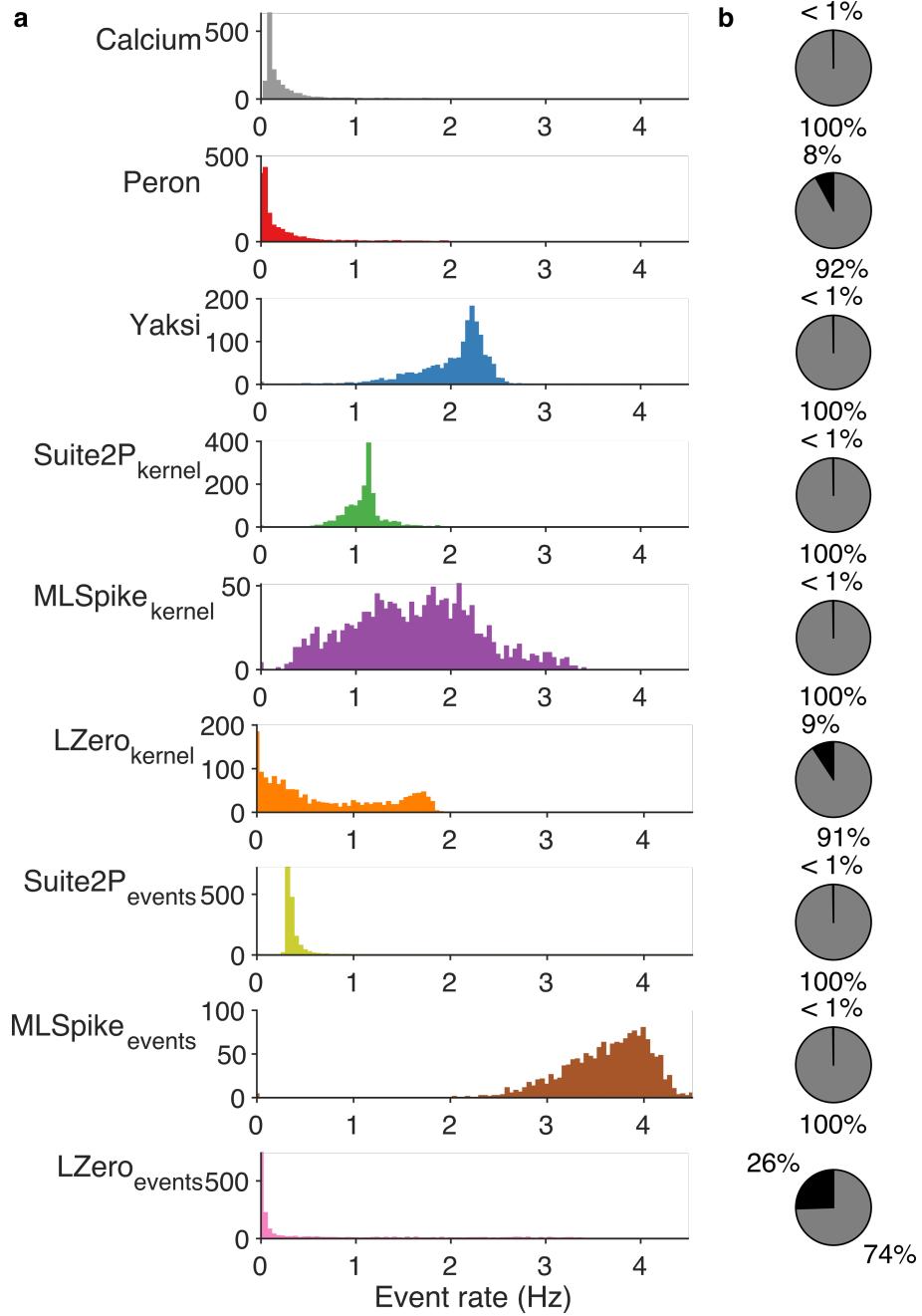


Figure 4: Estimates of population-wide event rates vary qualitatively across deconvolution methods.

(a) The distribution of event rate per neuron across the recorded population, according to each deconvolution method. For raw calcium and the five denoising methods (upper 6 panels), events are detected as fluorescence transients greater in magnitude than three standard deviations of background noise. The discrete deconvolution methods (lower 3 panels) return per frame: a spike count (MLSpike), a binary event detection (LZero), or an event magnitude (Suite2p); these time-series were thus sparse, with most frames empty.

(b) Proportion of active (gray) and silent (black) neurons for each method. Silent neurons are defined following (Peron et al., 2015b) as those with an event rate less than 0.0083Hz.

214 fewer than one spike every two minutes (O'Connor et al., 2010). For the nine approaches
215 we test here, six estimated the proportion of silent neurons to be less than 1%, including
216 two of the discrete deconvolution methods (Figure 4b). For raw calcium and methods
217 returning continuous time-series, raising the threshold for defining events will lead to more
218 silent neurons, but at the cost of further shifting the event rate distributions towards zero.
219 Even for simple firing statistics of neural activity, the choice of deconvolution method gives
220 widely differing, and sometimes wrong, results.

221 **2.5 Inferences of single neuron tuning differ widely between raw calcium**
222 **and deconvolved methods**

223 In any paradigm where one records the responses of neurons as an animal performs some
224 task, a basic question is what fraction of neurons in a target brain region are selective to
225 some aspect of the task. Here we ask how the detection of task-tuned neurons depends on
226 our choice of processing method for the raw calcium time-series.

227 The decision task facing the mouse (Fig. 3a) requires that it moves its whisker back-
228 and-forth to detect the position of the pole, delay for a second after the pole is withdrawn,
229 and then make a choice of the left or right lick-port based on the pole's position (Fig. 3b).
230 As the imaged barrel corresponds to the single spared whisker (on the contralateral side
231 of the face), so the captured population activity during each trial likely contains neurons
232 tuned to different aspects of the task.

233 Following Peron et al. 2015a, we define a task-tuned neuron as one for which the peak
234 in its trial-averaged histogram of activity exceeds the predicted upper limit from shuffled
235 data (see Methods). When applied to the raw calcium time-series, close to half the neurons
236 are tuned (734/1552; Fig.5a). This is more than double the proportion of tuned neurons
237 found for the next nearest method (Yaksi's simple deconvolution), and at least a factor of
238 5 greater than the proportion of tuned neurons resulting from any discrete deconvolution
239 method, which each report less than 10% of the neurons are tuned.

240 As these numbers suggest, we cannot assume that the tuned neurons in the raw cal-
241 cium are a good guide. Of the 734 tuned neurons in the raw calcium time-series, half (364,
242 49.5%) are unique, detected only in those time-series. By contrast, across all 8 deconvolu-
243 tion methods only 6 neurons are found tuned by one method alone. That we lose at least
244 half of the tuned neurons as soon as we apply any attempt to recover spike-events implies
245 that much of the detected tuning in the raw calcium time-series are noise.

246 There is little consistent agreement between the nine sets of time-series about which
247 neurons are tuned (Fig.5b). Only 104 neurons (6.7%) are labelled as tuned in at least
248 two out of the nine sets of time-series, and just 21 (1.35%) are labelled as tuned in all
249 nine. Even separately considering the continuous and discrete time-series, we find only 38
250 (2.4%) neurons are tuned across all six continuous methods, and 25 (1.6%) neurons for
251 all three discrete deconvolution methods (Fig.5c). Figure 5d illustrates the diversity of
252 detected tuning even amongst the neurons with the greatest agreement between methods.

253 These results suggest that raw calcium alone over-estimates tuning in the population,
254 but also that there can be substantial disagreement between deconvolution methods. One
255 solution for robust detection of tuned neurons is to find those agreed between the raw
256 calcium time-series and more than one deconvolution method. In Figure 5e-h, we show how
257 increasing the number of methods required to agree on a neuron's tuned status creates clear
258 agreement between time-series processed with all methods, even if a particular method did
259 not reach significance for that neuron. Even requiring agreement between the raw calcium
260 and just two other methods is enough to see tuning of many neurons. The identification
261 of unambiguously task-tuned neurons could thus be achieved by triangulating the raw

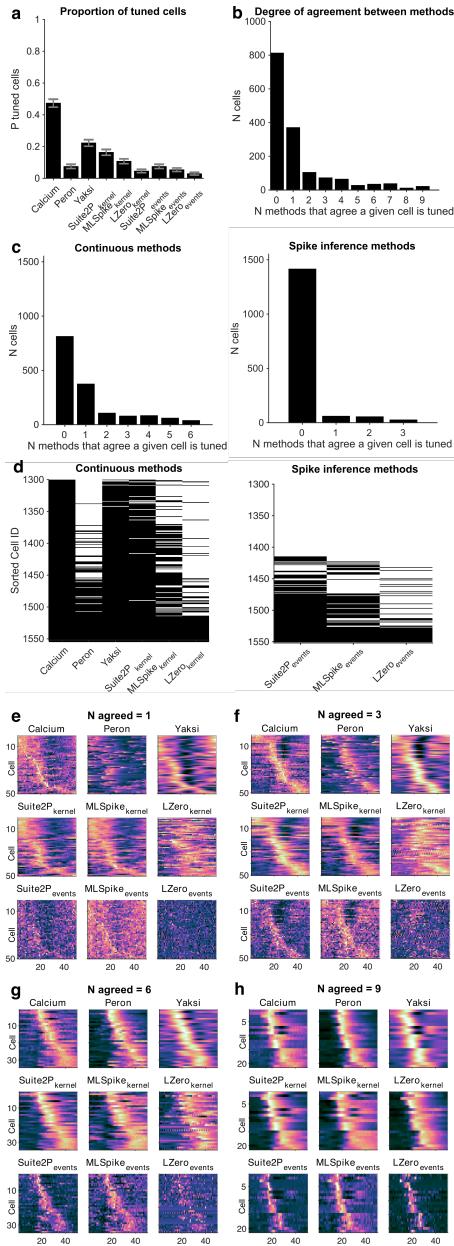


Figure 5: Inferences of single neuron tuning show poor agreement between raw calcium and deconvolution methods, and between methods.

(a0) Examples of a tuned (left) and non-tuned (right) neuron from the raw calcium time-series. X: Data; Y: upper 95% interval from shuffled data.

(a) Number of tuned neurons per deconvolution method. Error bars are 95% binomial confidence intervals.

(b) Agreement between methods. For each neuron, we count the number of methods (including raw calcium) for which it is labelled as tuned. Bars show the number of neurons classified as tuned by exactly N methods.

(c) Similar to (b), but breaking down the neurons into: agreement between methods (raw or denoising) resulting in continuous signals (left panel); and agreement between discrete deconvolution methods (right panel).

(d) Comparison of neuron tuning across methods. Each row shows whether that neuron is tuned (black) or not (white) under that deconvolution method. Cells are ordered from bottom to top by the number of methods that classify that neuron as tuned.

(e-h) Identifying robust neuron tuning. Panel groups (e) to (h) show neurons classed as tuned by increasing numbers of deconvolution methods. Each panel within a group plots one neuron's normalised (z-scored) trial-average histogram per row, ordered by the time of peak activity. The first panel in a group of 9 shows histograms from raw calcium signals; each of the 8 subsequent panel shows trial-average histograms for the same neurons, but following processing by each of the eight deconvolution methods.

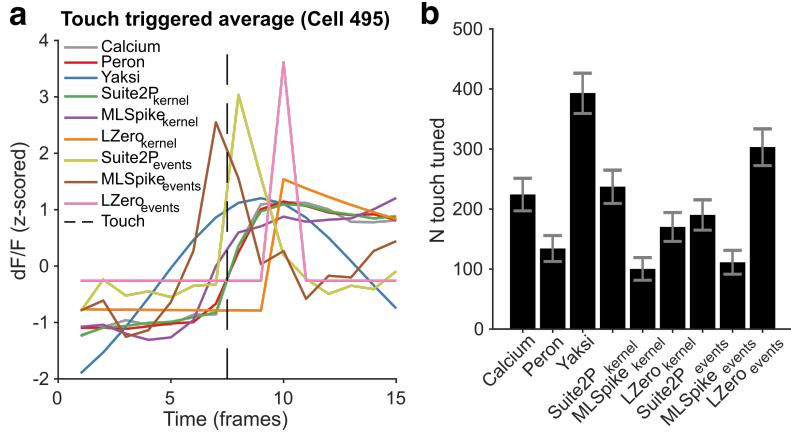


Figure 6: Touch-triggered neuron responses.

(a) Touch-triggered average activity from one neuron, across all deconvolution methods. The dotted line is the imaging frame in which the whisker touched the pole.

(b) Number of touch-tuned neurons across deconvolution methods. A neuron is classed as touch-tuned if its peak touch-triggered activity is significantly greater than shuffled data. Error bars are 95% Jeffreys confidence intervals for binomial data (Brown et al., 2001).

262 calcium with the output of multiple deconvolution methods.

263 In the pole detection task considered here, neurons tuned to pole contact are potentially
 264 crucial to understanding the sensory information used to make a decision. Touch onset is
 265 known to drive a subset of neurons to spike with short latency and low jitter (O'Connor
 266 et al., 2010; Hires et al., 2015). Detecting such rapid, precise responses in the slow kinetics
 267 of calcium imaging is challenging, suggesting discrete-deconvolution methods might be
 268 necessary to detect touch-tuned neurons. To test this, in each of the 9 sets of time-series
 269 we identify touch-tuned neurons by a significant peak in their touch-triggered activity
 270 (Fig 6a). Figure 6b shows that, while all data-sets have touch-tuned neurons, the number
 271 of such neurons differs substantially between them. And rather than being essential to
 272 detecting fast responding touch-tuned neurons, discrete deconvolution methods disagree
 273 strongly on touch-tuning, with LZero (events) finding more touch-tuned neurons than in
 274 the raw calcium, but MLSpike (events) finding less than half that number. Thus our
 275 inferences of the coding of task-wide or specific sensory events crucially depends on our
 276 choice of calcium imaging time-series.

277 2.6 Inconsistent recovery of population correlation structure across de- 278 convolution approaches

279 The high yield of neurons from calcium imaging is ideal for studying the dynamics and
 280 coding of neural populations (Harvey et al., 2012; Huber et al., 2012; Kato et al., 2015).
 281 Many analyses of populations start from pairwise correlations between neurons, whether
 282 as measures of a population's synchrony or joint activity, or as a basis for further analyses
 283 like clustering and dimension reduction (Cunningham and Yu, 2014). We now ask how our
 284 inferences of population correlation structure also depend on the choice of deconvolution
 285 method.

286 Figure 7a shows that the distributions of pairwise correlations qualitatively differ be-
 287 tween the sets of time-series we derived from the same calcium imaging data. The con-
 288 siderably narrower distributions from the discrete deconvolution time-series compared to
 289 the others is expected, as these time-series are sparse. Nonetheless, there are qualitative

290 differences within the sets of discrete and continuous time-series. Some distributions are
291 approximately symmetric, with broad tails; some asymmetric with narrow tails; the corre-
292 lation distribution from the Peron method time-series is the only one with a median below
293 zero. These qualitative differences are not due to noisy estimates of the pairwise correla-
294 tions: for all our sets of time-series the correlations computed on a sub-set of time-points
295 in the session agree well with the correlations computed on the whole session (Figure 7b).
296 Thus pairwise correlation estimates for each method are stable, but their distributions
297 differ between methods.

298 Looking in detail at the full correlation matrix shows that even for methods with similar
299 distributions, their agreement on correlation structure is poor. Some neuron pairs that ap-
300 pear correlated from time-series processed by one deconvolution method are uncorrelated
301 when processed with another method (Figure 7c). Over the whole population, the cor-
302 relation structure obtained from the raw calcium, Yaksi and Suite2p (kernel) time-series
303 all closely agree, but nothing else does (Figure 7d): the correlation structure obtained
304 from LZero agrees with nothing else; and the discrete deconvolution methods all generate
305 dissimilar correlation structures (Figure 7e). Our inferences about the extent and identity
306 of correlations within the population will differ qualitatively depending on our choice of
307 imaging time-series.

308 2.7 Deconvolution methods show the same population activity is both 309 low and high dimensional

310 Dimensionality reduction techniques, like principal components analysis (PCA), allow re-
311 searchers to make sense of large scale neuroscience data (Chapin and Nicolelis, 1999;
312 Briggman et al., 2005; Churchland et al., 2012; Harvey et al., 2012; Cunningham and Yu,
313 2014; Kobak et al., 2016), by reducing the data from N neurons to $d < N$ dimensions.
314 Key to such analyses is the choice of d , a choice guided by how much of the original data
315 we can capture. To assess such inferences of population dimensionality, we apply PCA to
316 our 9 sets of imaging time-series to estimate the dimensionality of the imaging data (which
317 for PCA is the variance explained by each eigenvector of the data's covariance matrix).

318 Figure 8a plots for each deconvolution method the cumulative variance explained when
319 increasing the number of retained dimensions. Most deconvolution methods qualitatively
320 disagree with the raw calcium data-set on the relationship between dimensions and vari-
321 ance. This relationship is also inconsistent across deconvolution methods; indeed the
322 discrete deconvolution methods result in the shallowest (MLSpike_{events}) and amongst the
323 steepest (LZero_{events}) relationships between increasing dimensions and variance explained.
324 The number of dimensions required to explain 80% of the variance in the data ranges
325 from $d = 125$ (Peron) to $d = 1081$ (MLSpike_{events}), a jump from 8% to 70% of all pos-
326 sible dimensions (Fig 8b). Thus we could equally infer that the same L2/3 population
327 activity is low dimensional (<10% dimensions required to explain 80% of the variance)
328 or high-dimensional (>50% of dimensions required) depending on our choice of imaging
329 time-series.

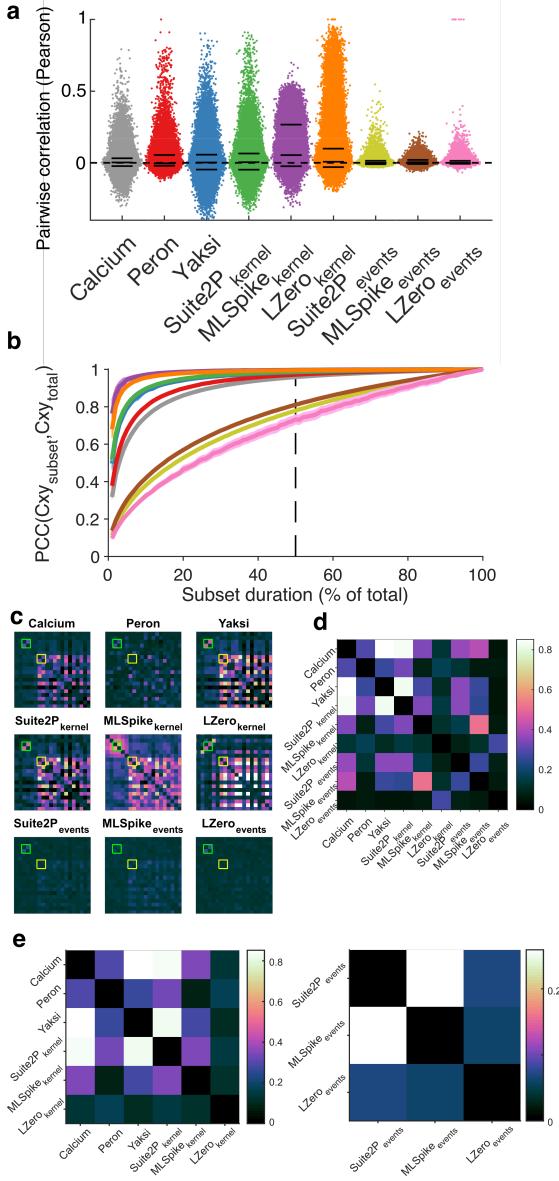


Figure 7: Effects of deconvolution on pairwise correlations between neurons.

- (a) Distributions of pairwise correlations between all neurons, for each deconvolution method (one dot per neuron pair, x-axis jitter added for clarity). Solid black lines are 5th, 50th and 95th percentiles.
- (b) Stability of correlation structure in the population. We quantify here the stability of the pairwise correlation estimates, by comparing the correlation matrix constructed on the full data ($C_{xy, total}$) to the same matrix constructed on a subset of the data ($C_{xy, subset}$). Each data-point is the mean correlation between $C_{xy, total}$ and $C_{xy, subset}$; one line per deconvolution method. Shaded error bars are one standard deviation of the mean across 100 random subsets.
- (c) Examples of qualitatively differing correlation structure across methods. Each panel plots the pairwise correlations for the same 50 neurons on the same colour scale. As examples, we highlight two pairs of neurons: one consistently correlated across different methods (green arrow and boxes); the other not (yellow arrow and boxes).
- (d) Comparison of pairwise correlation matrices between deconvolution methods. Each square is the Spearman's rank correlation between the full-data correlation matrix for that pair of methods. We use rank correlation to compare the ordering of pairwise correlations, not their absolute values.
- (e) as in (d), but split to show continuous methods (left) or discrete deconvolution methods (right).



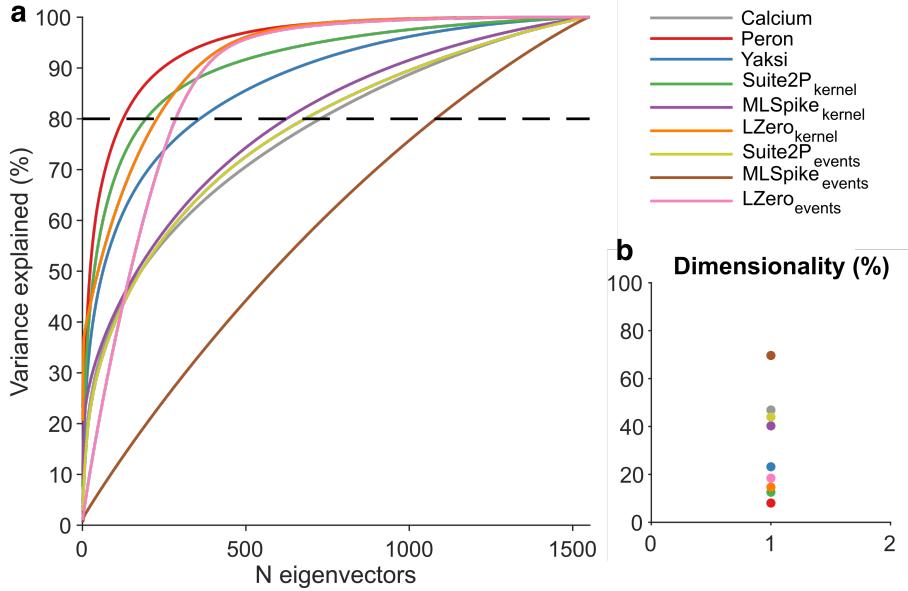


Figure 8: Dimensionality of population activity.

(a) Cumulative variance explained by each dimension of the `datas` covariance matrix, one line per deconvolution method. Dimensions are obtained from principal components analysis, and are ordered by decreasing contribution to the total variance explained. Dashed line is the 80% threshold used in panel (b).

(b) Proportion of dimensions required to explain 80% of the variance in the data.

3 Discussion

330 Imaging of somatic calcium is a remarkable tool for capturing the simultaneous activity of
 331 hundreds to thousands of neurons. But the time-series of each neuron's calcium fluorescence
 332 is inherently noisy and non-linearly related to its spiking. We sought here to address
 333 how our choice of corrections to these time-series – to use them raw, deconvolve them into
 334 continuous time-series, or deconvolve them into discrete events – affect the quality and
 335 reliability of the scientific inferences drawn. Our approach was to replicate the process of
 336 a typical population calcium-imaging study: choose an algorithm, choose its parameters
 337 using some reasonable heuristics, and analyse the resulting time-series.

338 Our results show the choice of processing qualitatively changes the potential scientific
 339 inferences we draw about the activity, coding, and correlation structure of a neural popula-
 340 tion in barrel cortex. Only the raw calcium and two of the processed time-series correctly
 341 capture the expected long-tailed distribution of spiking activity across the population.
 342 Neurons identified as being tuned to any feature of a pole-detection task differ widely be-
 343 tween processing methods. Few methods agree on the pairwise correlation structure of the
 344 population. Moreover, the apparent dimensionality of the population activity can differ by
 345 an order of magnitude across the processing methods. Across all analyses, we consistently
 346 observe that the results differ sharply between the raw calcium and most, if not all, of the
 347 processed time-series. However, the deconvolved time-series also consistently disagreed
 348 with each other, even between methods of the same broad class (continuous or discrete
 349 time-series).

351 3.1 Accurate discrete deconvolution is possible, but sensitive

352 We find much that is encouraging. In fitting discrete deconvolution methods to ground-
353 truth data, we found they can in principle accurately recover known spike-times from raw
354 calcium time-series. A caveat here is that the choice of metric for evaluation and fitting
355 of parameters is of critical importance. The widely-used Pearson correlation coefficient is
356 a poor choice of metric as it returns inconsistent results with small changes in algorithm
357 parameters, and leads to poor estimates of simple measures such as firing rate when used
358 across methods and sampling rates. By contrast, the Error Rate metric (Deneux et al.,
359 2016; Victor and Purpura, 1996) resulted in excellent recovery of ground-truth spike trains.
360 Other recently developed methods for comparing spike-trains based on information theory
361 (Theis et al., 2016) or fuzzy set theory (Reynolds et al., 2018), may also be appropriate.

362 However, while good estimates of ground-truth spike times can be achieved with mod-
363 ern discrete deconvolution methods (Berens et al., 2018; Pachitariu et al., 2018), the best
364 parameters vary substantially between cells, and small changes in analysis parameters
365 result in poor performance. This variation and sensitivity of parameters played out as
366 widely-differing results between the three discrete deconvolution methods in analyses of
367 neural activity, coding, and correlation structure.

368 3.2 Choosing parameters for deconvolution methods

369 A potential limitation of our study is that we use a single set of parameter values for each
370 discrete deconvolution method applied to the population imaging data from barrel cortex.
371 But then our situation is the same as that facing any experimentalist: in the absence of
372 ground-truth, how do we set the parameters? Our solution here was to use the modal
373 parameter values from ground-truth fitting. We also felt these were a reasonable choice
374 for the population imaging data from barrel cortex, given that the ground-truth recordings
375 came from the same species (mouse) in the same layer (2/3) of a different bit of primary
376 sensory cortex (V1). It would be instructive in future work to quantify the dependence
377 of analyses of neural activity, coding, and correlation on varying the parameters of each
378 deconvolution methods.

379 Rather than use the most general parameters values, another solution would be to
380 tune the parameters to obtain known gross statistics of the neural activity. This was
381 the approach used by Peron and colleagues (Peron et al., 2015b) to obtain the denoised
382 Peron time-series we included here. But as we've seen, this approach can lead to its own
383 problems: for example, in the Peron time-series, it created a distributions of correlations
384 that differed from any other set of time-series. Indeed, finding good parameter values may
385 be an intractable problem, as our results suggest each neuron requires individual fitting,
386 to reflect the combination of its expression of fluorescent protein, and its particular non-
387 linearity between voltage and calcium.

388 3.3 Ways forward

389 One solution to disagreement between deconvolution methods is to just use the raw calcium
390 time-series. Many studies use the raw calcium signal as the basis for all their analyses
391 (Harvey et al., 2012; Huber et al., 2012; Chu et al., 2016), perhaps assuming this is
392 the least biased approach. Our results show this is not so: the discrepancy between
393 raw and deconvolved calcium on single neuron coding suggests an extraordinary range of
394 possible results, from about half of all neurons tuned to the task down to less 5 percent.
395 The qualitative conclusion – there is coding – is not satisfactory. Moreover, as noted by

396 (Sabatini, 2019), the raw calcium fluorescence signal is a low-pass filtered version of the
397 underlying spike train, which places strong limits on the maximum correlation between the
398 raw signal and underlying spikes, and hence on any correlations between the raw signal
399 and the behavioural variables related to those spikes. Thus our results should not be
400 interpreted as a call to abandon deconvolution methods; rather they serve to delimit how
401 we can interpret their outputs.

402 A simple solution to the inconsistencies between different forms of deconvolved time-
403 series is to triangulate them, and take the consensus across their results. For example,
404 our finding of a set of tuned neurons across multiple methods is strong evidence that
405 neurons in L2/3 of barrel cortex are responsive across the stages of the decision task.
406 Further examples of such triangulation in the literature are rare; Klaus and colleagues
407 (Klaus et al., 2017) used two different pipelines to derive raw $\Delta F/F$ of individual neurons
408 from one-photon fibre-optic recordings in the striatum, and replicated all analyses using
409 the output of both pipelines. Our results encourage the further use of triangulation to
410 create robust inference: obtaining the same result in the face of wide variation increases
411 our belief in its reliability (Munafò and Davey Smith, 2018).

412 There are caveats to using triangulation. For single neuron analyses, triangulation
413 inevitably comes at the price of reducing the yield of neurons to which we can confidently
414 assign roles. A further problem for triangulation is how to combine more complex analyses,
415 such as pairwise correlations; the alternative is to rely on qualitative comparisons. There
416 is also an assumption that all contributions to the consensus contain useful data: if one
417 deconvolution method returns time-series with no relation to the underlying spike events,
418 then including its outputs in the consensus would inevitably worsen the results.

419 Another proposed solution is better forward models, like ML Spike, for the link from
420 spiking to calcium fluorescence Greenberg et al. (2018). Indeed, as sensors with faster
421 kinetics (though fundamentally limited by kinetics of calcium release itself) and higher
422 signal-to-noise ratios are developed (Badura et al., 2014; Dana et al., 2016, 2019), so
423 the accuracy and robustness of de-noising and deconvolution should improve; and as the
424 neuron yield continues to increase (Ahrens et al., 2013; Stringer et al., 2019), so the
425 potential for insights from inferred spikes or spike-driven events grows. Developing further
426 advanced deconvolution algorithms will harness these advances, but are potentially always
427 limited by the lack of ground-truth to fit their parameters (Wei et al., 2019). Worse,
428 no matter how good the forward model for a single neuron, our results suggest the wide
429 variation in the model parameters needed for each neuron would make population analyses
430 challenging to interpret.

431 Our results provide impetus for different directions of research, not just to improving
432 our modelling of the relationship between spikes and the somatic calcium signal, but
433 also focussing on how we can get consensus among the output of different deconvolution
434 algorithms, and thus provide robust scientific inferences about neural populations.

435 **4 Methods**

436 **Ground truth data**

437 Ground truth data was accessed from crcns.org (Svoboda, 2015), and the experiments have
438 been described previously (Chen et al., 2013). Briefly, mouse visual cortical neurons ex-
439 pressing the fluorescent calcium reporter protein GCaMP6s were imaged with two-photon
440 microscopy at 60Hz. Loose-seal cell-attached recordings were performed simultaneously
441 at 10kHz. Recordings were made in awake mice during 5 trials (4s blank, 4s stimulus) of
442 the optimal moving grating stimulus (1 of 8 directions) for the cell-attached neuron. The
443 data-set contains twenty one recordings from nine neurons.

444 **Population imaging data description**

445 Population imaging data was accessed from crcns.org and have been described previ-
446 ously (Peron et al., 2015b). Briefly, volumetric two photon calcium imaging of primary
447 somatosensory cortex (S1) was performed in awake head-fixed mice performing a whisker-
448 based object localisation task. In the task a metal pole was presented **on** one of two loca-
449 tions and mice were motivated with fluid reward to lick at one of two lick ports depending
450 on the location of the pole following a brief delay. Two photon imaging of GCaMP6s
451 expressing neurons in superficial S1 was performed at 7Hz. Images were motion corrected
452 and aligned, before regions of interest were manually set and neuropil-subtracted. A single
453 recording from this dataset was used for population analysis. The example session had
454 1552 neurons recorded for a total of 23559 frames (56 minutes).

455 **List of deconvolution methods**

456 **MLSpike**

457 MLSpike (Deneux et al., 2016) was accessed from <https://github.com/mlspike>. MLSpike
458 uses a model-based probabilistic approach to recover spike trains in calcium imaging data
459 by taking baseline fluctuations and cellular properties into account. Briefly, MLSpike
460 implements a model of measured calcium fluorescence as a combination of spike-induced
461 transients, background (photonic) noise and drifting baseline fluctuations. A maximum
462 likelihood approach determines the probability of the observed calcium at each time step
463 given an inferred spike train generated through a particular set of model parameters.
464 MLSpike returns a maximum a posteriori spike train (as used here), or a spike probability
465 per time step. MLSpike also returns an estimate of the drifting background fluorescence
466 which is ignored in this work.

467 MLSpike has a number of free parameters, of which we optimise three: *A*, the mag-
468 nitude of fluorescence transients caused by a single spike; *tau*, calcium fluorescence decay
469 time; *sigma*, background (photonic) noise level. MLSpike also has parameters for different
470 calcium sensor kinetics (for OBG, GCaMP3, GCaMP6 and so on) which we fix to default
471 values for GCaMP6.

472 For our analysis of event rate MLSpike's spike train was counted (mean event count
473 per second), and for subsequent analyses was converted to a dense array of spike counts
474 per imaging frame.

475 **Suite2P**

476 Suite2P (Pachitariu et al., 2016, 2018) was accessed from [https://github.com/cortex-lab/Suite2P](https://github.com/cortex-
477 lab/Suite2P). Suite2P was developed as a complete end-to-end processing pipeline for

478 large scale 2-photon imaging analysis - from image registration to spike extraction and
479 visualization - of which we only use the spike extraction step. The spike deconvolution
480 of Suite2P uses a sparse non-negative deconvolution algorithm, greedily identifying and
481 removing calcium transients to minimise the cost function

$$C = \|F - s * k\|^2,$$

482 where the cost C is the squared norm of fluorescence F minus a reconstruction of that
483 signal comprising a sparse array of spiking events s multiplied by a parameterised calcium
484 kernel k . The kernel was parameterised following defaults for GCaMP6s (exponential
485 decay of 2 seconds, though it has been shown the precise value of this parameter does not
486 affect performance for this method (Pachitariu et al., 2018)).

487 Suite2P has a further free parameter which sets the minimum spike event size, the
488 *Threshold*, which determines the stopping criteria for the algorithm.

489 Elements of s are of varying amplitude corresponding to the amplitude of the calcium
490 transients at that time. For ground truth firing rate analysis we are interested in each
491 algorithm's ability to recover spike trains, therefore we treat each event as a 'spike' and
492 optimise the algorithm appropriately. For our analysis of event rate Suite2P's event train
493 was counted (mean event count per second), and for subsequent analyses was converted
494 to a dense array of varying amplitude events (i.e. s) per imaging frame.

495 LZero

496 The method we refer to as LZero was written in Matlab based on an implementation
497 in *R* accessed at <https://github.com/jewellsean/LZeroSpikeInference>. A full description is
498 available in the paper of Jewell and Witten (2018). Briefly, in LZero spike detection is cast
499 as a change-point detection problem, which could be solved with an l_0 optimization algo-
500 rithm. Working backwards from the last time point the algorithm finds time points where
501 the calcium dynamics abruptly change from a smooth exponential rise. These change
502 points correspond to spike event times. Spike inference accuracy is assessed similarly to
503 Suite2P by measuring the fit between observed fluorescence and a reconstruction based
504 on inferred spike times and a fixed calcium kernel.

505 LZero has two free parameters - *lambda*, a tuning parameter that controls the trade-off
506 between the sparsity of the estimated spike event train and the fit of the estimated calcium
507 to the observed fluorescence; and *scale*, the magnitude of a single spike induced change in
508 fluorescence.

509 For our analysis of event rate LZero's spike train was counted (mean event count per
510 second), and for subsequent analyses was converted to a dense array of spikes per imaging
511 frame (maximum one spike per imaging frame due to limitations of the algorithm).

512 Yaksi

513 Yaksi is an implementation of the deconvolution approach of Yaksi and Friedrich (2006).
514 The fluorescence time series is low-pass filtered (4th order butterworth filter, 0.7Hz cutoff)
515 to remove noise before having a calcium kernel (exponential decay of 2 seconds, as used
516 in Suite2P and LZero above) linearly deconvolved out of the signal using Matlab's `deconv`
517 function. The output of Yaksi is a continuous signal approximating spike density per unit
518 time.

519 **Peron events**

520 Peron events refer to the de-noised calcium event traces detailed in the original Peron
521 et al. (2015b) paper. Here a version of the ‘peeling’ algorithm (Lütcke et al., 2013) was
522 developed, a template-fitting algorithm with variable decay time constants across events
523 and neurons. This algorithm was tuned by the authors to generate a low number of false
524 positive detections (a rate of 0.01Hz) on ground truth data, matching firing rate statistics
525 from cell-attached electrophysiology and leading to a hit rate of 54%. The output for
526 analysis is a continuous signal approximating de-noised calcium concentration per unit
527 time.

528 **Events and kernel versions of spike inference methods**

529 Where a spike inference method returns spike counts per time point, these are plotted
530 as Method_{events}. To compare to other methods that return a de-noised dF/F or firing
531 rate estimates, these event traces are convolved with a calcium kernel and plotted as
532 Method_{kernel}. The kernel used is consistent with that used as a default for GCaMP6s
533 in ML Spike, Suite2P and LZero, namely an exponential decay of two seconds duration
534 normalised to have an integral of 1.

535 **Ground truth spike train metrics**

536 Pearson correlation coefficient was computed between the ground truth and inferred spikes
537 (ML Spike) or events (Suite2P, LZero) following convolution of both with a gaussian kernel
538 (61 samples wide, 1.02 seconds).

539 Error Rate was computed between the ground truth and inferred spikes/events using
540 the Deneux et al. (2016) implementation of normalised error rate, derived from Victor
541 and Purpura (1996) Error Rate (code available <https://github.com/MLspike>). Briefly,
542 the error rate is 1 - F1-score, where the F1-score is the harmonic mean of sensitivity and
543 precision (Davis and Goadrich, 2006),

$$\begin{aligned} \text{sensitivity} &= 1 - \frac{\text{misses}}{\text{total spikes}}, \\ \text{precision} &= 1 - \frac{\text{false detections}}{\text{total detections}}, \\ \text{ErrorRate} &= 1 - 2 \frac{\text{sensitivity} \times \text{precision}}{\text{sensitivity} + \text{precision}}. \end{aligned}$$

544 Hits, misses and false detections were counted with a temporal precision of 0.5 seconds.
545 For normalised estimation of errors in firing/event rate we compute,

$$\frac{\text{estimated rate} - \text{true rate}}{\text{true rate}},$$

546 where spike/event rates are measured in Hz.

547 **Parameter fitting**

548 For each method the best parameters for each neuron were determined by brute force
549 search over an appropriate range (i.e. at least two orders of magnitude encompassing
550 full parameter ranges used in the original publications for each method). The parameter
551 ranges were explored on a log scale as follows: ML Spike A (0.01:1, 21 values), tau (0.01:5,

552 21 values), sigma (0.01:1, 21 values); Suite2P Threshold (0.1:100, 13 values); LZero lambda
553 (0.1:20, 23 values), scale (0.1:20, 23 values).

554 The modal best parameters, as determined using Error Rate on downsampled data,
555 were then fixed for the population imaging data analysis. These were: MLSpike A: 0.1995,
556 tau: 1.9686, sigma: 0.0398; Suite2P Threshold: 1.7783; LZero sigma: 0.1; lambda: 3.1623.

557 **Downsampling**

558 Ground truth calcium data was downsampled from 60Hz to 7Hz in Matlab by up-sampling
559 by 7 (interpolating the signal) and then downsampling the resultant 420Hz time-series of
560 frames to 7 Hz by sampling every 60th frame.

561 **4.1 Event rate estimation**

562 Spike inference methods (Suite2P_{events}, MLSpike_{events}, LZero_{events}) return estimated spike
563 times (MLSpike), or event times (Suite2P/LZero) which were converted into mean event
564 rates (Hz) per neuron.

565 The event rate for continuous methods (Calcium, Peron, Yaksi, Suite2P_{kernel}, MLSpike_{kernel},
566 LZero_{kernel}) for each neuron was determined by counting activity/fluorescence transients
567 greater than three standard deviations of the background noise. Background noise was
568 calculated by subtracting a four-frame moving average of the fluorescence from the raw
569 data to result in a 'noise only' trace. This operation was done separately for each neuron
570 and each method. Event rate was then computed in Hz.

571 Silent neurons were defined as neurons with event rates below 0.0083Hz (or fewer than
572 one spike per two minutes of recording) as in O'Connor et al. (2010).

573 **4.2 Task-tuned neurons**

574 Task-tuning was determined for each neuron using the model-free approach of Peron et al.
575 (2015b). Neurons were classed as task-tuned if their peak trial-average activity exceeded
576 the 95th percentile of a distribution of trial-average peaks from shuffled data (10000 shuf-
577 fles). The shuffle test was done separately for correct lick-left and lick-right trials and
578 neurons satisfying the tuning criteria in either case were counted as task-tuned.

579 Tuned neuron agreement was calculated as the number of methods that agreed to the
580 tuning status of a given neuron, for all methods and separately for continuous and spike
581 inference methods.

582 **4.3 Touch-tuned responses**

583 Touch-tuned neurons were determined by first computing touch-triggered average activity
584 for each neuron, then calculating whether the data distribution of peak touch-induced
585 activity exceeds the expected activity of resampled data. In more detail, the time of first
586 touch between the mouse's whisker and the metal pole on each trial was recorded. For
587 each neuron, one second of activity (seven data samples) was extracted before and after
588 the frame closest to the first touch of each trial (15 frames total per trial); taking the
589 mean touch-triggered activity over trials gave the average touch response for the neuron.
590 To determine whether the neuron was touch tuned or not, we compared the neuron's
591 peak mean response r_{data} to a null distribution by taking a randomly sampled 15 frame
592 segment of a trial, finding the peak mean response across trials r_{null} , and repeating this
593 calculation for 10000 random samples. A p-value for the data peak response was calculated

594 as $p = (\#r_{null} < r_{data})/10000$. Over all neurons, a neuron was considered touch-tuned if
595 $p < 0.05$ after Benjamini-Hochberg correction.

596 4.4 Pairwise correlations

597 Pairwise correlations (Pearson correlation coefficients, Fig. 7a) were calculated between
598 all pairs of neurons at the data sampling rate (7Hz).

599 Stability of correlation estimates (Fig. 7b) at the recording durations used was assessed
600 by computed the similarity between correlation distributions for the the intact dataset to
601 those from subsets of the dataset. For each deconvolution method, we computed the
602 pairwise correlation matrix using the entire sessions data, as above. We also sampled a
603 subset of time-points (1%-100%) of the full dataset at random without replacement and
604 computed a matrix of pairwise correlations for this subset. We then compute the similarity
605 between the total and subset matrices using Pearsons correlation coefficient. This process
606 was repeated 100 times and the mean (line) and standard deviation (shading) of the 100
607 repeats were plotted.

608 4.5 Correlations between correlation matrices

609 Correlations between correlation matrices (Fig. 7c-e) were computed using Spearman's
610 rank correlation between the unique pairwise correlations from each method (i.e. the
611 upper triangular entries of the correlation matrix).

612 4.6 Dimensionality

613 To determine the dimensionality of each dataset we performed eigendecomposition of the
614 covariance matrix of each dataset. The resultant eigenvalues were sorted into descending
615 order, and the cumulative variance explained plotted, the number of eigenvectors required
616 to explain 80% of the variance recorded.

617 References

- 618 Misha B Ahrens, Michael B Orger, Drew N Robson, Jennifer M Li, and Philipp J Keller.
619 Whole-brain functional imaging at cellular resolution using light-sheet microscopy. *Nat.*
620 *Methods*, 10(5):413–420, May 2013.
- 621 Aleksandra Badura, Xiaonan Richard Sun, Andrea Giovannucci, Laura A Lynch, and
622 Samuel S-H Wang. Fast calcium sensor proteins for monitoring neural activity. *Neurophotonics*,
623 1(2):025008, October 2014.
- 624 Philipp Berens, Jeremy Freeman, Thomas Deneux, Nicolay Chenkov, Thomas McColgan,
625 Artur Speiser, Jakob H Macke, Srinivas C Turaga, Patrick Mineault, Peter Rupprecht,
626 Stephan Gerhard, Rainer W Friedrich, Johannes Friedrich, Liam Paninski, Marius Pa-
627 chitariu, Kenneth D Harris, Ben Bolte, Timothy A Machado, Dario Ringach, Jasmine
628 Stone, Luke E Rogerson, Nicolas J Sofroniew, Jacob Reimer, Emmanouil Froudarakis,
629 Thomas Euler, Miroslav Román Rosón, Lucas Theis, Andreas S Tolias, and Matthias
630 Bethge. Community-based benchmarking improves spike rate inference from two-photon
631 calcium imaging data. *PLoS Comput. Biol.*, 14(5):e1006157, May 2018.
- 632 K L Briggman, H D I Abarbanel, and W B Kristan, Jr. Optical imaging of neuronal
633 populations during decision-making. *Science*, 307(5711):896–901, February 2005.

- 634 Emery N Brown, Robert E Kass, and Partha P Mitra. Multiple neural spike train data
635 analysis: state-of-the-art and future challenges. *Nat. Neurosci.*, 7(5):456–461, May 2004.
- 636 Lawrence D. Brown, T. Tony Cai, and Anirban DasGupta. Interval estimation for a
637 binomial proportion. *Statist Sci*, 16:101–133, 2001.
- 638 J K Chapin and M A Nicolelis. Principal component analysis of neuronal ensemble activity
639 reveals multidimensional somatosensory representations. *J. Neurosci. Methods*, 94(1):
640 121–140, December 1999.
- 641 Tsai-Wen Chen, Trevor J Wardill, Yi Sun, Stefan R Pulver, Sabine L Renninger, Amy
642 Baohan, Eric R Schreiter, Rex A Kerr, Michael B Orger, Vivek Jayaraman, Loren L
643 Looger, Karel Svoboda, and Douglas S Kim. Ultrasensitive fluorescent proteins for
644 imaging neuronal activity. *Nature*, 499(7458):295–300, July 2013.
- 645 Monica W Chu, Wankun L Li, and Takaki Komiyama. Balancing the robustness and
646 efficiency of odor representations during learning. *Neuron*, 92:174–186, Oct 2016. ISSN
647 1097-4199. doi: 10.1016/j.neuron.2016.09.004.
- 648 Mark M Churchland, John P Cunningham, Matthew T Kaufman, Justin D Foster, Paul
649 Nuyujukian, Stephen I Ryu, and Krishna V Shenoy. Neural population dynamics during
650 reaching. *Nature*, 487(7405):51–56, July 2012.
- 651 John P Cunningham and Byron M Yu. Dimensionality reduction for large-scale neural
652 recordings. *Nat. Neurosci.*, 17(11):1500–1509, November 2014.
- 653 Hod Dana, Boaz Mohar, Yi Sun, Sujatha Narayan, Andrew Gordus, Jeremy P Hasseman,
654 Getahun Tsegaye, Graham T Holt, Amy Hu, Deepika Walpita, Ronak Patel, John J
655 Macklin, Cornelia I Bargmann, Misha B Ahrens, Eric R Schreiter, Vivek Jayaraman,
656 Loren L Looger, Karel Svoboda, and Douglas S Kim. Sensitive red protein calcium
657 indicators for imaging neural activity. *Elife*, 5, March 2016.
- 658 Hod Dana, Yi Sun, Boaz Mohar, Brad K Hulse, Aaron M Kerlin, Jeremy P Hasseman,
659 Getahun Tsegaye, Arthur Tsang, Allan Wong, Ronak Patel, John J Macklin, Yang
660 Chen, Arthur Konnerth, Vivek Jayaraman, Loren L Looger, Eric R Schreiter, Karel
661 Svoboda, and Douglas S Kim. High-performance calcium sensors for imaging activity
662 in neuronal populations and microcompartments. *Nat. Methods*, 16(7):649–657, July
663 2019.
- 664 Jesse Davis and Mark Goadrich. The relationship between Precision-Recall and ROC
665 curves. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML
666 '06, pages 233–240, New York, NY, USA, 2006. ACM.
- 667 Thomas Deneux, Attila Kaszas, Gergely Szalay, Gergely Katona, Tamás Lakner, Amiram
668 Grinvald, Balázs Rózsa, and Ivo Vanzetta. Accurate spike estimation from noisy calcium
669 signals for ultrafast three-dimensional imaging of large neuronal populations in vivo.
670 *Nat. Commun.*, 7:12190, July 2016.
- 671 Johannes Friedrich, Pengcheng Zhou, and Liam Paninski. Fast online deconvolution of
672 calcium imaging data. *PLoS Comput. Biol.*, 13(3):e1005423, March 2017.
- 673 Andrea Giovannucci, Johannes Friedrich, Pat Gunn, Jérémie Kalfon, Brandon L Brown,
674 Sue Ann Koay, Jiannis Taxidis, Farzaneh Najafi, Jeffrey L Gauthier, Pengcheng Zhou,
675 Baljit S Khakh, David W Tank, Dmitri B Chklovskii, and Eftychios A Pnevmatikakis.

- 676 CaImAn an open source tool for scalable calcium imaging data analysis. *Elife*, 8, January
677 2019.
- 678 David S Greenberg, Damian J Wallace, Kay-Michael Voit, Silvia Wuertenberger, Uwe
679 Czubayko, Arne Monsees, Takashi Handa, Joshua T Vogelstein, Reinhard Seifert,
680 Yvonne Groemping, and Jason ND Kerr. Accurate action potential inference from a cal-
681 cium sensor protein through biophysical modeling. *bioRxiv*, 2018. doi: 10.1101/479055.
682 URL <https://www.biorxiv.org/content/early/2018/11/29/479055>.
- 683 Kenneth D Harris, Rodrigo Quian Quiroga, Jeremy Freeman, and Spencer L Smith. Im-
684 proving data quality in neuronal population recordings. *Nat. Neurosci.*, 19(9):1165–
685 1174, August 2016.
- 686 Christopher D Harvey, Philip Coen, and David W Tank. Choice-specific sequences in
687 parietal cortex during a virtual-navigation decision task. *Nature*, 484(7392):62–68, April
688 2012.
- 689 Samuel Andrew Hires, Diego A Gutnisky, Jianing Yu, Daniel H O'Connor, and Karel
690 Svoboda. Low-noise encoding of active touch by layer 4 in the somatosensory cortex.
691 *Elife*, 4, August 2015.
- 692 Daniel Huber, D A Gutnisky, S Peron, D H O'Connor, J S Wiegert, L Tian, T G Oertner,
693 L L Looger, and K Svoboda. Multiple dynamic representations in the motor cortex
694 during sensorimotor learning. *Nature*, 484(7395):473–478, April 2012.
- 695 Sean Jewell and Daniela Witten. Exact spike train inference via ℓ_0 optimization. *Ann.*
696 *Appl. Stat.*, 12(4):2457–2482, December 2018.
- 697 Patrick Kaifosh, Jeffrey D Zaremba, Nathan B Danielson, and Attila Losonczy. SIMA:
698 Python software for analysis of dynamic fluorescence imaging data. *Front. Neuroinform.*,
699 8:80, September 2014.
- 700 Saul Kato, Harris S Kaplan, Tina Schrödel, Susanne Skora, Theodore H Lindsay, Eviatar
701 Yemini, Shawn Lockery, and Manuel Zimmer. Global brain dynamics embed the motor
702 command sequence of *caenorhabditis elegans*. *Cell*, 163(3):656–669, October 2015.
- 703 Sander W Keemink, Scott C Lowe, Janelle M P Pakan, Evelyn Dylda, Mark C W van
704 Rossum, and Nathalie L Rochefort. FISSA: A neuropil decontamination toolbox for
705 calcium imaging signals. *Sci. Rep.*, 8(1):3493, February 2018.
- 706 Andreas Klaus, Gabriela J Martins, Vitor B Paixao, Pengcheng Zhou, Liam Paninski, and
707 Rui M Costa. The spatiotemporal organization of the striatum encodes action space.
708 *Neuron*, 95(5):1171–1180.e7, August 2017.
- 709 Dmitry Kobak, Wieland Brendel, Christos Constantinidis, Claudia E Feierstein, Adam
710 Kepecs, Zachary F Mainen, Xue-Lian Qi, Ranulfo Romo, Naoshige Uchida, and Chris-
711 tian K Machens. Demixed principal component analysis of neural population data. *Elife*,
712 5, April 2016.
- 713 Henry Lütcke, Felipe Gerhard, Friedemann Zenke, Wulfram Gerstner, and Fritjof Helm-
714 chen. Inference of neuronal network spike dynamics and topology from calcium imaging
715 data. *Front. Neural Circuits*, 7:201, December 2013.
- 716 Eran A Mukamel, Axel Nimmerjahn, and Mark J Schnitzer. Automated analysis of cellular
717 signals from large-scale calcium imaging data. *Neuron*, 63(6):747–760, September 2009.

- 718 Marcus R Munafó and George Davey Smith. Robust research needs many lines of evidence.
719 *Nature*, 553(7689):399–401, January 2018.
- 720 Daniel H O’Connor, Simon P Peron, Daniel Huber, and Karel Svoboda. Neural activity
721 in barrel cortex underlying vibrissa-based object localization in mice. *Neuron*, 67(6):
722 1048–1061, September 2010.
- 723 Marius Pachitariu, Carsen Stringer, Sylvia Schröder, Mario Dipoppa, L Federico Rossi,
724 Matteo Carandini, and Kenneth D Harris. Suite2p: beyond 10, 000 neurons with stan-
725 dard two-photon microscopy. *BioRxiv*, Preprint at <http://dx.doi.org/10.1101/061507>,
726 2016.
- 727 Marius Pachitariu, Carsen Stringer, and Kenneth D Harris. Robustness of spike deconvolu-
728 tion for neuronal calcium imaging. *J. Neurosci.*, August 2018.
- 729 António R C Paiva, Il Park, and José C Príncipe. A comparison of binless spike train
730 measures. *Neural Comput. Appl.*, 19(3):405–419, April 2010.
- 731 Simon Peron, Tsai-Wen Chen, and Karel Svoboda. Comprehensive imaging of cortical
732 networks. *Curr. Opin. Neurobiol.*, 32:115–123, June 2015a.
- 733 Simon P Peron, Jeremy Freeman, Vijay Iyer, Caiying Guo, and Karel Svoboda. A cellular
734 resolution map of barrel cortex activity during tactile behavior. *Neuron*, 86(3):783–799,
735 May 2015b.
- 736 Eftychios A Pnevmatikakis, Daniel Soudry, Yuanjun Gao, Timothy A Machado, Josh
737 Merel, David Pfau, Thomas Reardon, Yu Mu, Clay Lacefield, Weijian Yang, Misha
738 Ahrens, Randy Bruno, Thomas M Jessell, Darcy S Peterka, Rafael Yuste, and Liam
739 Paninski. Simultaneous denoising, deconvolution, and demixing of calcium imaging
740 data. *Neuron*, January 2016.
- 741 Ruben Portugues, Claudia E Feierstein, Florian Engert, and Michael B Orger. Whole-
742 brain activity maps reveal stereotyped, distributed networks for visuomotor behavior.
743 *Neuron*, 81(6):1328–1343, March 2014.
- 744 Stephanie Reynolds, Therese Abrahamsson, Per Jesper Sjöström, Simon R Schultz, and
745 Pier Luigi Dragotti. CosMIC: A consistent metric for spike inference from calcium
746 imaging. *Neural Comput.*, 30(10):2726–2756, October 2018.
- 747 Bernardo Sabatini. The impact of reporter kinetics on the interpretation of data gathered
748 with fluorescent reporters. *bioRxiv*, page 834895, 2019. doi: 10.1101/834895. URL
749 <https://www.biorxiv.org/content/early/2019/11/07/834895>.
- 750 Carsen Stringer and Marius Pachitariu. Computational processing of neural recordings
751 from calcium imaging data. *Curr. Opin. Neurobiol.*, 55:22–31, December 2018.
- 752 Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Matteo Carandini, and Ken-
753 neth D Harris. High-dimensional geometry of population responses in visual cortex.
754 *Nature*, June 2019.
- 755 K Svoboda. Simultaneous imaging and loose-seal cell-attached electrical recordings from
756 neurons expressing a variety of genetically encoded calcium indicators. *GENIE project,*
757 *Janelia Farm Campus, HHMI; CRCNS.org*, 2015.

- 758 Lucas Theis, Philipp Berens, Emmanouil Froudarakis, Jacob Reimer, Miroslav
759 Román Rosón, Tom Baden, Thomas Euler, Andreas S Tolias, and Matthias Bethge.
760 Benchmarking spike rate inference in population calcium imaging. *Neuron*, 90(3):471–
761 482, May 2016.
- 762 J D Victor and K P Purpura. Nature and precision of temporal coding in visual cortex:
763 a metric-space analysis. *J. Neurophysiol.*, 76(2):1310–1326, August 1996.
- 764 Joshua T Vogelstein, Adam M Packer, Timothy A Machado, Tanya Sippy, Baktash Babadi,
765 Rafael Yuste, and Liam Paninski. Fast nonnegative deconvolution for spike train infer-
766 ence from population calcium imaging. *J. Neurophysiol.*, 104(6):3691–3704, December
767 2010.
- 768 Ziqiang Wei, Bei-Jung Lin, Tsai-Wen Chen, Kayvon Daie, Karel Svoboda, and Shaul
769 Druckmann. A comparison of neuronal population dynamics measured with calcium
770 imaging and electrophysiology. *bioRxiv*, 2019. doi: 10.1101/840686. URL <https://www.biorxiv.org/content/early/2019/11/15/840686>.
- 772 Adrien Wohrer, Mark D Humphries, and Christian K Machens. Population-wide distri-
773 butions of neural activity during perceptual decision-making. *Prog. Neurobiol.*, 103:
774 156–193, April 2013.
- 775 Emre Yaksci and Rainer W Friedrich. Reconstruction of firing rate changes across neuronal
776 populations by temporally deconvolved ca2+ imaging. *Nat. Methods*, 3(5):377–383, May
777 2006.