

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

June 25, 2020

Abstract

Calcium imaging is a powerful tool for capturing the simultaneous activity of large populations of neurons. Here we determine the extent to which our inferences of neural population activity, correlations, and coding depend on our choice of whether and how we deconvolve the calcium time-series into spike-driven events. To this end, we use a range of deconvolution algorithms to create nine versions of the same calcium imaging data obtained from barrel cortex during a pole-detection task. Seeking suitable values for the deconvolution algorithms' parameters, we optimise them against ground-truth data, and find those parameters both vary by up to two orders of magnitude between neurons and are sensitive to small changes in their values. Applied to the barrel cortex data, 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 sensitive to the choice of if and how spike-evoked events are recovered. 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., 2019a; 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 slow on the time-scale of spikes, 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 cal-

36 cium into something more closely approximating spikes. Deconvolution algorithms them-
37 selves range in complexity from simple deconvolution with a fixed kernel of the calcium
38 response (Yaksi and Friedrich, 2006), through detecting spike-evoked calcium events (Jew-
39 ell and Witten, 2018; Pachitariu et al., 2016), to directly inferring spike times (Vogelstein
40 et al., 2010; Lütcke et al., 2013; Deneux et al., 2016). This continuum of options raises
41 the further question of the extent to which we should process the raw calcium signals. We
42 address here the question facing any systems neuroscientist using calcium imaging: do we
43 use the raw calcium, or attempt to clean it up? Thus our aim is to understand if our
44 choice matters: to what extent do our inferences about neural activity, correlations, and
45 coding depend on our choice of raw or deconvolved calcium time-series.

46 We proceed here in two stages. In order to use deconvolution algorithms, the data
47 analyst needs to choose their parameters. We thus first address how good these algorithms
48 can be in principle with optimised parameters, and how sensitive their results are to the
49 choice of parameter values. To do so, we evaluate qualitatively different deconvolution
50 algorithms by optimising their parameters against ground truth data with known spikes.

51 With our understanding of their parameters in hand, we then turn to our main ques-
52 tion, by analysing a large-scale population recording from the barrel cortex of a mouse
53 performing a whisker-based decision task. We compare estimates of population coding and
54 correlations obtained using either raw calcium signals, or a range of time-series derived
55 from those calcium signals, covering simple deconvolution, event detection, and spikes.

56 We find that a substantial fraction of the deconvolution methods used here fail to
57 recover basic features of population activity in barrel cortex established from electro-
58 physiology. The inferences we draw about coding qualitatively differ between raw and
59 deconvolved calcium signals. In particular, coding analyses based on raw calcium sig-
60 nals detect an order of magnitude more neurons tuned to task features. Yet there is also
61 qualitative disagreement between deconvolution methods on which neurons are tuned.
62 The inferences we draw about correlations between neurons do not distinguish between
63 raw and deconvolved calcium signals, but can qualitatively differ between deconvolution
64 methods. Our results thus suggest care is needed in drawing inferences from population
65 recordings of somatic calcium, and that one solution is to replicate all results in both raw
66 and deconvolved calcium signals.

67 **2 Results**

68 **2.1 Performance of deconvolution algorithms on ground-truth data-sets**

69 We select here three deconvolution algorithms that infer discrete spike-like events, each
70 an example of the state of the art in qualitatively different approaches to the problem:
71 Suite2p (Pachitariu et al., 2016), a peeling algorithm that matches a scalable kernel to the
72 calcium signal to detect spike-triggered calcium events; LZero (Jewell and Witten, 2018), a
73 change-point detection algorithm, which finds as events the step-like changes in the calcium
74 signal that imply spikes; and MLspike (Deneux et al., 2016), a forward model, which fits
75 an explicit model of the spike-to-calcium dynamics in order to find spike-evoked changes
76 in the calcium signal, and returns spike times. We emphasize that these methods were
77 chosen as exemplars of their approaches, and are each innovative takes on the problem; we
78 are not here critiquing individual methods, nor are we seeking a “best” method. Rather,
79 we are using an array of methods to illustrate the problems and decisions facing the data
80 analyst when using calcium imaging data.

81 We first ask if these deconvolution methods work well in principle, by testing if there

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

89 The best-fit parameters depend strongly on how we evaluate the match between true
90 and inferred spike trains. The Pearson correlation coefficient between the true and inferred
91 spike train is a common choice (Brown et al., 2004; Paiva et al., 2010; Theis et al., 2016;
92 Reynolds et al., 2018; Berens et al., 2018), typically with both trains convolved with a
93 Gaussian kernel to allow for timing errors. However, we find that choosing parameters to
94 maximise the correlation coefficient can create notable errors. The inferred spike trains
95 from MLSpike have too many spikes on average (mean error over recordings: 31.72%),
96 and the accuracy of recovered firing rates widely varies across recordings (Fig 1b, blue
97 symbols). We attribute these errors to the noisy relationship between the correlation
98 coefficient and the number of inferred spikes (Figure 1c): for many recordings, there is
99 no well-defined maximum coefficient, especially for the amplitude parameter A , so that
100 near-maximum correlation between true and inferred trains is consistent with a wide range
101 of spike counts in the inferred trains. We see the same sensitivity for the event rates from
102 recordings optimised using Suite2p (Figure 1e) and LZero (Figure 1f, top). If we compare
103 their inferred event rates to true firing rates (Fig 1b), we see Suite2p estimates far more
104 events than spikes (mean error 79.47%) and LZero fewer events than spikes (mean error: -
105 21.14%). These further errors are problematic: there cannot be more spike-driven calcium
106 events than spikes, and LZero’s underestimate is considerably larger than the fraction of
107 frames with two or more spikes ($< 0.002\%$ frames).

108 To address the weaknesses of the Pearson correlation coefficient, we instead optimise
109 parameters using the error rate metric of Deneux et al. (2016). The error rate is derived
110 from the proportions of missed and excess spikes (see Methods), and returns a normalised
111 score between 0 for a perfect match between two spike trains, and 1 when all the spikes are
112 missed. This comparison between inferred and true spike trains is most straightforward
113 for algorithms like MLSpike that directly return spike times; for the other algorithms, we
114 use here their event times as inferred spikes, a reasonable choice given the low firing rate
115 and well separated spikes in the ground truth data. Choosing parameters to minimise the
116 error rate between the true and inferred spike-trains results in excellent recovery of the
117 true number of spikes for all three deconvolution methods (Fig 1b, green symbols), with
118 mean errors in spike counts of 12% excess spikes for Suite2P, 7.3% for MLSpike, and 5%
119 for LZero. As we show in Figure 1d-f, for all three deconvolution methods the error rate
120 has a well-defined minimum for almost every recording. Consequently, all deconvolution
121 methods can, in principle, accurately recover the true spike-trains given an appropriate
122 choice of 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 population
125 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 1g, optimising
129 parameters using the minimum error rate still results in excellent recovery of the true spike
130 rate (and interestingly for some recordings reduces the error when using the correlation

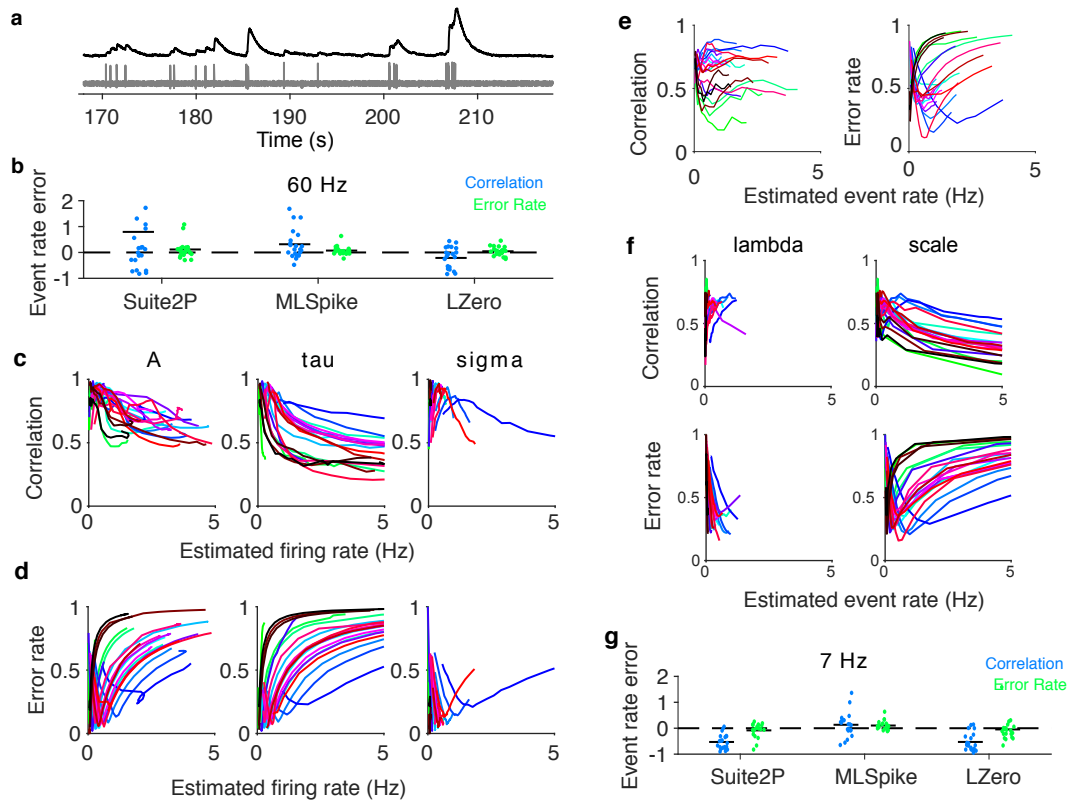


Figure 1: Deconvolution algorithms can accurately recover spiking events in principle

(a) Example simultaneous recording of somatic voltage (grey) and calcium activity (black) imaged at 60Hz.

(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, and the errors for parameters optimised to minimise the error rate. Horizontal black bars are means. Error is computed relative to the true firing rate: $(Rate_{true} - Rate_{estimated})/Rate_{true}$; and error of 1 thus corresponds to twice as many estimated spikes as there are in the ground-truth data. For LZero and Suite2p, $Rate_{estimated}$ is computed from event times.

(c) Dependence of MLspike’s deconvolution performance on the firing rate of the inferred spike train. For each of MLSpike’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 obtained at each tested parameter value. One line per recording; colours are used solely to help distinguish the lines. Parameters: A , calcium transient amplitude per spike ($\Delta F/F$); τ , calcium decay time constant (s); σ , background (photonic) noise level ($\Delta F/F$)

(d) as in (c), but using error rate between the true and inferred spikes.

(e) Dependence of Suite2p’s deconvolution performance on the firing rate of the inferred event train as its detection threshold parameter is varied. Left: correlation coefficient; right: error rate.

(f) Dependence of LZero’s deconvolution performance on the firing rate of the inferred event train, as its two parameters are varied: λ , sparsity of spike events; scale, the magnitude of a single spike-induced fluorescence change.

(g) As for (b), but with the somatic calcium down-sampled to 7Hz before optimising parameters for the deconvolution methods.

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

145 The problem of between-neuron variation in parameter values would be compensated
146 somewhat if the quality of the inferred spike or event trains is robust to changes in those
147 values. However, we find performance is highly sensitive to changes in some parameters.
148 Figure 2c-e shows that for most recordings the quality of the inferred spike train abruptly
149 worsens with small increases or decreases in the best parameter, regardless of the decon-
150 volution method used. As we show in Figure 2f, the inferred spikes for a single neuron
151 can vary dramatically as we change a parameter value, even when we restrict ourselves to
152 just the range of optimised values across the recordings. That the parameters are sensitive
153 and vary considerably across neurons has the significant implication that, unless ground
154 truth data is available for every neuron being analysed, deconvolution algorithms could
155 be substantially inaccurate.

156 **2.3 Deconvolution of population imaging in barrel cortex during a de-** 157 **cision task**

158 We turn now to the core problem facing any analyst of population calcium imaging data:
159 there are rarely ground truth data, and never for every neuron. In the absence of ground-
160 truth data, there is no way of selecting a “best” deconvolution algorithm or a “best”
161 set of parameters for analysing a population recording. Yet the above results imply that
162 the insights we gain about population activity would indeed depend crucially on which
163 deconvolution method we use. We now test the extent of this dependence by applying
164 8 different deconvolution methods to the same raw calcium time-series, and compare the
165 resulting statistics of neural activity, properties of neural coding, and the extent and
166 structure of correlations between neurons.

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

172 Our primary goal is to understand how the choices of deconvolving these calcium-
173 imaging data alter the scientific inferences we can draw. As our baseline, we use the
174 “raw” $\Delta F/F$ time-series of changes in calcium indicator fluorescence. We use the above
175 three discrete deconvolution methods to extract spike counts (MLSpike), event occurrence

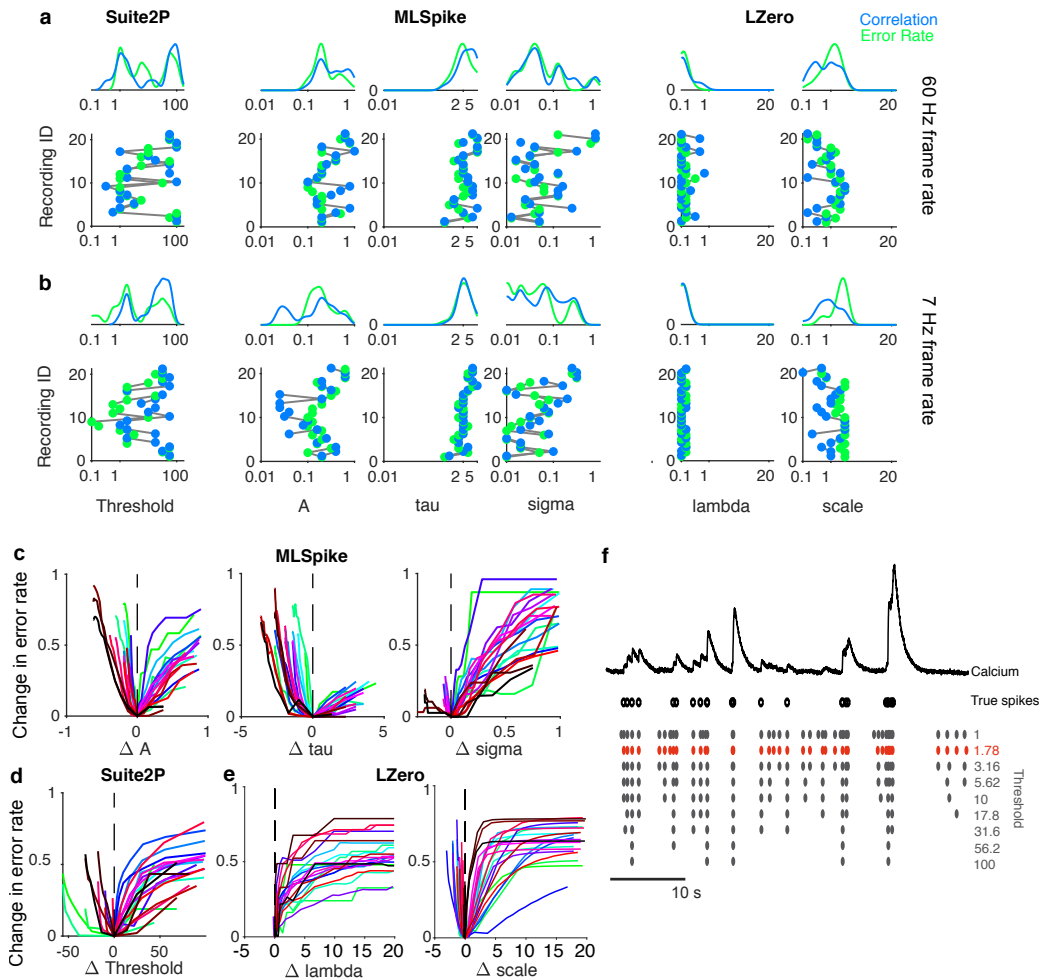


Figure 2: Variation in best-fit spike deconvolution parameters across ground-truth recordings.

(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 recordings. We plot for each recording the optimised parameter value found using correlation coefficient and error rate. Lines join recordings from the same neuron.

(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 MLSpike’s free parameters. One line per recording, at 60 Hz frame rate.

(d) As for panel (c), for changes in Suite2p’s threshold parameter.

(e) As for panel (c), for changes in LZero’s two parameters.

(f) Example of the range of inferred spike event trains possible when applying plausible but wrong parameter values to a recording. For one recording, we plot in red the inferred spike events detected using its optimised threshold parameter for Suite2p. Alongside we plot the inferred trains of spike events that result if we vary the threshold parameter across the range of optimised values found within the set of 21 recordings (values in panel (a), optimised using error rate).

176 (LZero), or event magnitude (Suite2p) per frame. Given the above-demonstrated depen-
 177 dence of these algorithms on their parameters, we use [Yaksi and Friedrich \(2006\)](#)’s simple
 178 deconvolution of the raw calcium with a fixed kernel of the GCaMP6s response to a single
 179 spike, whose only free parameters are fixed from data. For comparison, we use [Peron et al.](#)

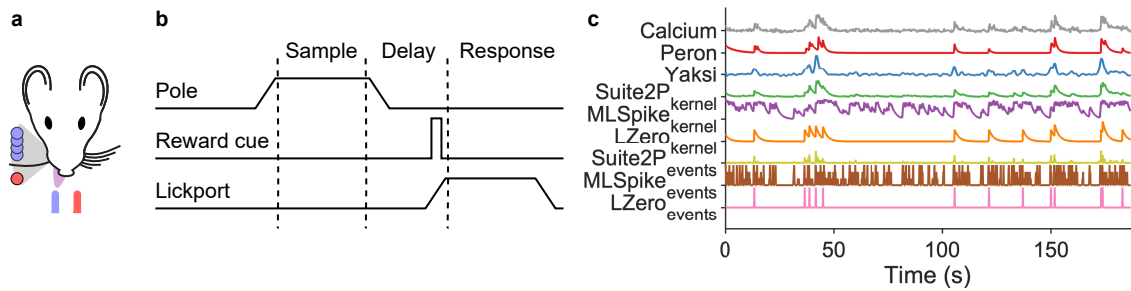


Figure 3: Experimental data from Peron et al. (2015b).

(a) Schematic of task set-up. A pole was raised within range of the single right-side whisker; the pole's position, forward (red circle) or backward (blue circles) 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.

180 (2015b)'s own version of denoised calcium time-series, which they created using a custom
181 version of the peeling algorithm (Lütcke et al., 2013), a greedy template-fitting event-
182 detection algorithm with variable rise and decay time constants across events. The Peron
183 time-series are then the detected spike-events convolved with a kernel of the detected rise
184 and decay time. And finally, for comparison with the Peron time-series, we create equiva-
185 lent versions for our three discrete-deconvolution methods, by convolving their recovered
186 spikes/events with a fixed GCaMP6s spike-response kernel. Figure 3c show an example
187 raw calcium time-series for one neuron, and the result of applying each of these 8 process-
188 ing methods. We thus repeat all analyses on 9 different sets of time-series extracted from
189 the same population recording.

190 We choose the algorithm parameters as follows. Simple deconvolution (Yaksi and
191 Friedrich, 2006) uses a parameterised kernel of the GCaMP6s response to a single spike.
192 For the three discrete deconvolution methods, we choose the modal values of the best-fit
193 parameters that optimised the error rate over the ground-truth recordings. This seems
194 a reasonably consistent way obtaining comparable results between methods, by using
195 the most consistently performing values obtained from comparable data: neurons in the
196 same layer (L2/3) in the same species (mouse), in another primary sensory area (V1).
197 Most importantly for our purposes, choosing the modal values means we avoid extreme
198 and potentially pathological regions of the parameter space. Again this recapitulates the
199 problem facing any analyst of population calcium imaging data, of how to choose the
200 parameters for a deconvolution algorithm in the absence of any ground-truth recordings.

201 2.4 Deconvolution methods disagree on estimates of simple neural statis- 202 tics

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

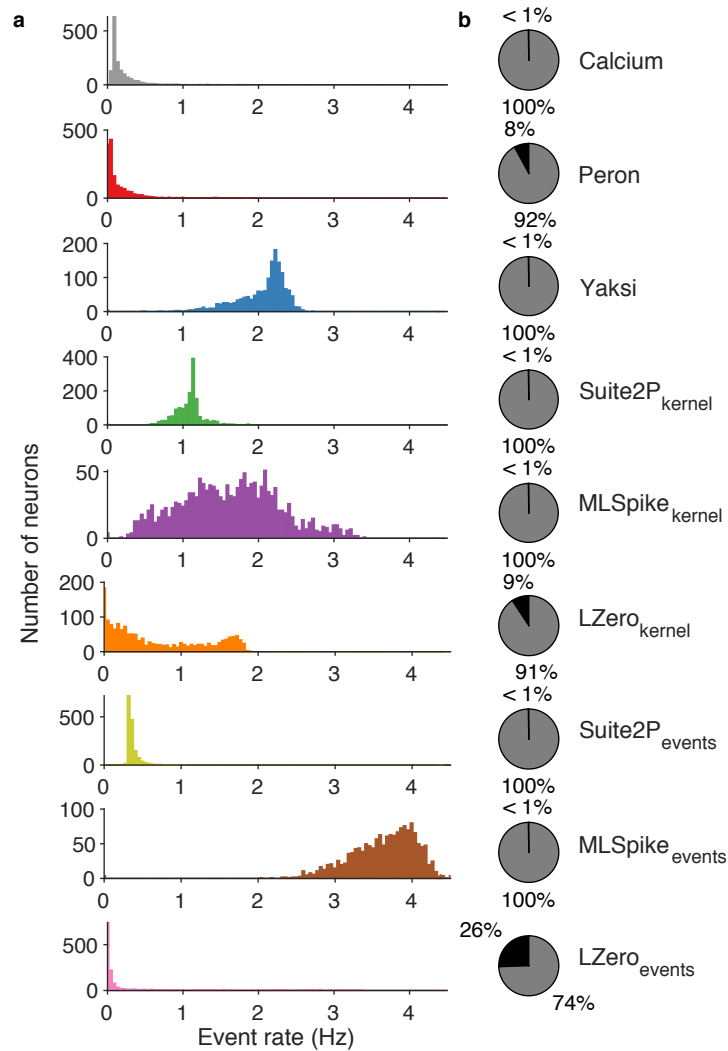


Figure 4: Estimates of population-wide event rates vary qualitatively across deconvolution methods.

(a) The distribution of event rate per neuron across the recorded population, according to each deconvolution method. For raw calcium and the five continuous versions of the time-series (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.

211 was obtained by thresholding the calcium time-series)

212 Figure 4a shows that the raw calcium and two of the discrete deconvolution methods
213 (Suite2p_{events}, LZero_{events}) qualitatively match the expected distributions of event rates
214 (median near zero, long right-hand tails). The Peron time-series also have the correct
215 distribution of event rates. All other methods give wrong distributions, whether of spike
216 rates (MLSpike) or event rates (all other methods). There is also little overlap in the
217 distributions of spike rates between the three discrete deconvolution methods. Applying a
218 kernel to their inferred spikes/events shifts rather than smooths the firing rate distributions
219 (Suite2P_{kernel}, MLSpike_{kernel}, LZero_{kernel}), suggesting noise in the deconvolution process
220 is amplified through the additional steps of convolving with a kernel and thresholding.

221 Cell-attached recordings in barrel cortex have shown that ~26% of L2/3 pyramidal
222 neurons are silent during a similar pole localisation task, with silence defined as emitting
223 fewer than one spike every two minutes (O'Connor et al., 2010). For the nine approaches
224 we test here, six estimated the proportion of silent neurons to be less than 1%, including
225 two of the discrete deconvolution methods (Figure 4b). For raw calcium and methods
226 returning continuous time-series, raising the threshold for defining events will lead to more
227 silent neurons, but at the cost of further shifting the event rate distributions towards zero.
228 Even for simple firing statistics of neural activity, the choice of deconvolution method gives
229 widely differing, and sometimes wrong, results.

230 2.5 Inferences of single neuron tuning differ widely between raw calcium 231 and deconvolved methods

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

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

242 Following Peron et al. 2015a, we define a task-tuned neuron as one for which the peak
243 in its trial-averaged activity exceeds the predicted upper limit from shuffled data (Fig. 5a).
244 When we apply this definition to the raw calcium time-series, close to half the neurons are
245 tuned (734/1552; Fig.5b). This is more than double the proportion of tuned neurons we
246 find for the next nearest method (Yaksi's simple deconvolution), and at least a factor of
247 5 greater than the proportion of tuned neurons resulting from any discrete deconvolution
248 method ("events"), which each report less than 10% of the neurons are tuned.

249 The wide variation in numbers reflects little consistent agreement between the nine
250 sets of time-series about which neurons are tuned. A substantial fraction of the neurons
251 are found tuned in only one time-series of the nine (Fig.5c). And that time-series is
252 overwhelmingly the raw calcium: of the 734 tuned neurons in the raw calcium time-series,
253 half (364, 49.5%) are unique, detected only in those time-series. By contrast, across all 8
254 deconvolution methods only 6 neurons are found tuned by one method alone. Thus either
255 the raw calcium time-series contains many erroneously-detected tuned neurons, or the 8
256 deconvolution methods combined miss many tuned neurons, or both.

257 One likely source of this broad disagreement is that the raw calcium time-series allows a
258 generous definition of "tuned". Spike-evoked changes to the somatic calcium concentration

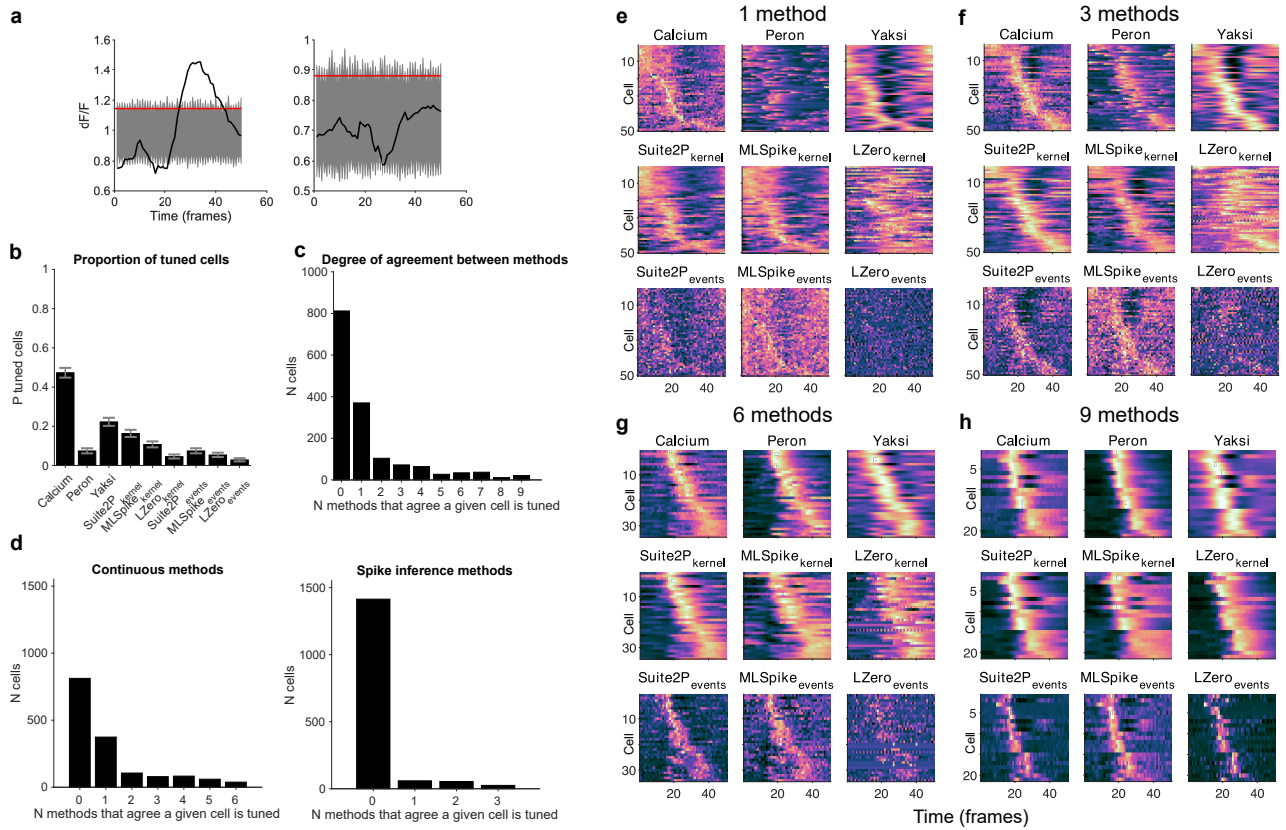


Figure 5: Inferences of single neuron tuning show poor agreement between raw calcium and deconvolution methods, and between methods.

(a) Examples of a tuned (left) and non-tuned (right) neuron from the raw calcium time-series. Black: trial-averaged calcium fluorescence. Grey shading: full range of $\Delta F/F$ from the shuffled data. Red line: 95th percentile of the peak $\Delta F/F$ value across the shuffled data.

(b) Number of tuned neurons per deconvolution method. Error bars are 95% Jeffreys confidence intervals for binomial data (Brown et al., 2001).

(c) 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.

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

(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 panels shows trial-average histograms for the same neurons, but following processing by each of the eight deconvolution methods.

259 are slow on the time-scale of spikes, and the calcium sensors are slower still: the GCaMP6s
260 sensor’s response to a single spike has a rise-to-maximum time of around 0.2 seconds, and
261 a half-width decay time of at least 0.5 seconds (Chen et al., 2013; Dana et al., 2014). This
262 creates a strong low-pass filtering of the underlying spike train, leading to weak correlation
263 between the timing of the spikes and the timing of the calcium changes (see also Sabatini
264 (2019)). Consequently, in the raw calcium time-series a neuron could emit a spike on each
265 trial that are over a second apart between trials, yet each would still contribute to a “peak”
266 in the trial-averaged signal. (Indeed in Fig.5b we see this low-pass filtering effect in the
267 convolved versions of our discrete-deconvolution time-series: we always get more tuned
268 neurons in the “kernel” versions despite them having identical underlying spike-events to
269 the “event” versions.) Finding a tuned neuron in the raw calcium time-series tells us only
270 that the neuron was active during the trial, not that its spikes were specifically tuned to
271 some event in the world.

272 Indeed, the need to correct the low correlation between raw calcium time-series and
273 behaviour or events in the world is a major reason why deconvolution methods have been
274 developed. Simply convolving the raw calcium signal with a fixed-parameter kernel as in
275 the Yaksi method immediately halves the number of apparently tuned neurons, potentially
276 because the immediate rise time and fixed decay time of the kernel reduce the low-pass
277 filtering of the spike train. And when recovering discrete spike-events a neuron can only be
278 “tuned” if those spike-events align in time across trials, leading to far fewer tuned neurons.
279 We can see then our choice of time-series processing creates a continuum of definitions of
280 neuron tuning in this analysis.

281 But the implicit definition of tuning is not the only source of disagreement between
282 the 9 time-series. Of the neurons found tuned in more than one time-series, the agreement
283 is still poor. Just 21 (1.35%) are labelled as tuned in all nine (Fig.5c). Even separately
284 considering the continuous and spike-event time-series, we find only 38 (2.4%) neurons are
285 tuned across all six continuous methods, and 25 (1.6%) neurons for all three spike-event
286 deconvolution methods (Fig.5d). Even between time-series with similar implicit definitions
287 of “tuned”, there is inconsistency about which neurons fit that definition.

288 An approach for the consistent detection of tuned neurons is to find those agreed
289 between the raw calcium time-series and more than one deconvolution method. In Figure
290 5e-h, we show how increasing the number of methods required to agree on a neuron’s tuned
291 status creates clear agreement between time-series processed with all methods, even if a
292 particular method did not reach significance for that neuron. Even requiring agreement
293 between the raw calcium and just two other methods is enough to see tuning of many
294 neurons. More reliable identification of task-tuned neurons could potentially be achieved
295 by triangulating the raw calcium with the output of multiple deconvolution methods.

296 In the pole detection task considered here, neurons tuned to pole contact are potentially
297 crucial to understanding the sensory information used to make a decision. Touch onset is
298 known to drive a subset of neurons in barrel cortex to spike with short latency and low
299 jitter (O’Connor et al., 2010; Hires et al., 2015). Detecting such rapid, precise responses
300 in the slow kinetics of calcium imaging is challenging, suggesting discrete-deconvolution
301 methods might be necessary to detect touch-tuned neurons. To test this, in each of the
302 9 sets of time-series we identify touch-tuned neurons by a significant peak in their touch-
303 triggered activity (Fig 6a). Figure 6b shows that, while all data-sets have touch-tuned
304 neurons, the number of such neurons differs substantially between them. And rather than
305 being essential to detecting fast responding touch-tuned neurons, discrete deconvolution
306 methods disagree strongly on touch-tuning, with LZero (events) finding more touch-tuned
307 neurons than in the raw calcium, but MLSpike (events) finding less than half that number.

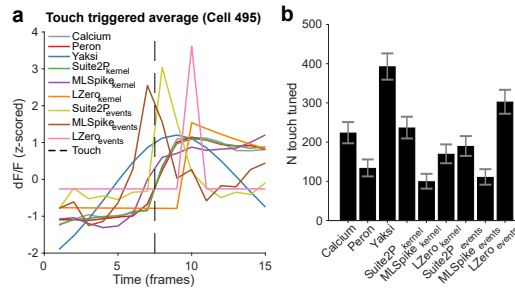


Figure 6: Touch-triggered neuron responses.

(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 randomly resampled data. Error bars are 95% Jeffreys confidence intervals for binomial data.

308 Thus our inferences of the coding of task-wide or specific sensory events crucially depends
 309 on both whether we deconvolve the raw calcium time-series or not, and on which algorithm
 310 we choose to do so.

311 2.6 Inconsistent recovery of population correlation structure across de- 312 convolution approaches

313 The high yield of neurons from calcium imaging is ideal for studying the dynamics and
 314 coding of neural populations (Harvey et al., 2012; Huber et al., 2012; Kato et al., 2015).
 315 Many analyses of a population’s dynamics or coding start from pairwise correlations be-
 316 tween its neurons, whether as measures of a population’s synchrony or joint activity, or as
 317 a basis for further analyses like clustering and dimension reduction (Cunningham and Yu,
 318 2014; Humphries, 2017; Stringer et al., 2019b). Consequently, differences in correlation
 319 estimates will play out as different inferences of population dynamics or population cod-
 320 ing. For example, weak correlations between neurons in primary visual cortex would be
 321 evidence of sparse coding of visual information (Stringer et al., 2019a; Rumyantsev et al.,
 322 2020). We now ask how our inferences of population correlation structure in the barrel
 323 cortex data also depend on the choice of deconvolution method.

324 Figure 7a shows that the distributions of pairwise correlations qualitatively differ be-
 325 tween the sets of time-series we derived from the same calcium imaging data. The con-
 326 siderably narrower distributions from the discrete deconvolution time-series compared to
 327 the others is expected, as these time-series are sparse. Nonetheless, there are qualitative
 328 differences within the sets of discrete and continuous time-series. Some distributions are
 329 approximately symmetric, with broad tails; some asymmetric with narrow tails; the corre-
 330 lation distribution from the Peron method time-series is the only one with a median below
 331 zero. These qualitative differences are not due to noisy estimates of the pairwise correla-
 332 tions: for all our sets of time-series the correlations computed on a sub-set of time-points
 333 in the session agree well with the correlations computed on the whole session (Figure 7b).
 334 (Although we note that, as expected, the three spike-event time-series require far more
 335 time-points to obtain stable correlation estimates, because of their sparse events). Thus
 336 pairwise correlation estimates for each method are stable, but their distributions differ
 337 between methods.

338 Looking in detail at the full correlation matrix shows that even for methods with similar
 339 distributions, their agreement on correlation structure is poor. Some neuron pairs that ap-

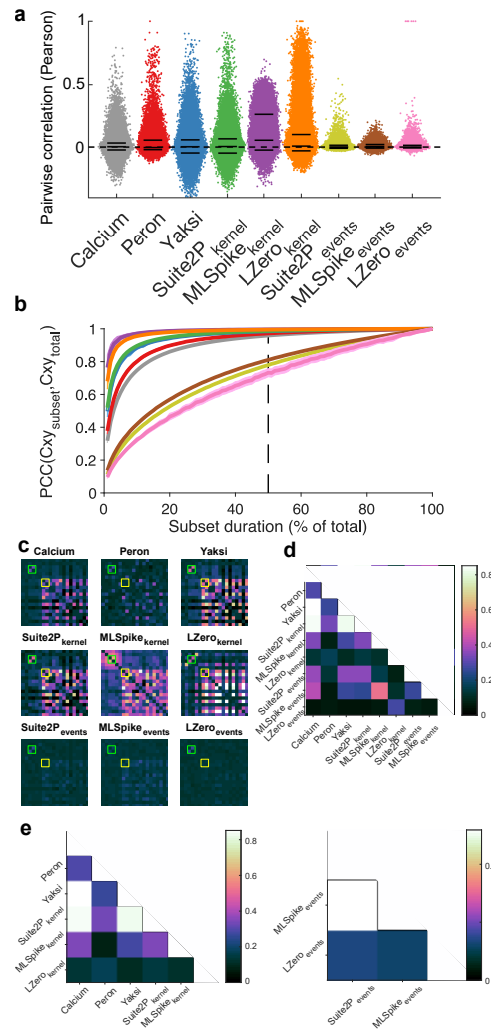


Figure 7: Effects of deconvolution on pairwise correlations between neurons.

(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 (Cxy_{total}) to the same matrix constructed on a subset of the data (Cxy_{subset}). Each data-point is the mean correlation between Cxy_{total} and Cxy_{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 boxes); the other not (yellow 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).

340 pear correlated from time-series processed by one deconvolution method are uncorrelated
341 when processed with another method (Figure 7c). Over the whole population, the cor-
342 relation structure obtained from the raw calcium, Yaksi and Suite2p (kernel) time-series
343 all closely agree, but nothing else does (Figure 7d): the correlation structure obtained
344 from LZero agrees with nothing else; and the discrete deconvolution methods all generate
345 dissimilar correlation structures (Figure 7e). Our inferences about the extent and identity
346 of correlations within the population will differ qualitatively depending on our choice of
347 deconvolution method.

348 2.7 Deconvolution methods show the same population activity is both 349 low and high dimensional

350 Dimensionality reduction techniques, like principal components analysis (PCA), allow re-
351 searchers to make sense of large scale neuroscience data (Chapin and Nicolelis, 1999;
352 Briggman et al., 2005; Churchland et al., 2012; Harvey et al., 2012; Cunningham and Yu,
353 2014; Kobak et al., 2016), by reducing the data from N neurons to $d < N$ dimensions. Key
354 to such analyses is the choice of d dimensions, a choice guided by how much of the origi-
355 nal data we can capture. Differences in dimensionality imply different computations: for
356 example, low-dimensional activity implies a sensory population uses a redundancy code,
357 while high-dimensional activity implies the population uses a sparse code (Wohrer et al.,
358 2013). To assess such inferences of population dimensionality, we apply PCA to our 9 sets
359 of imaging time-series to estimate the dimensionality of the imaging data (which for PCA
360 is the variance explained by each eigenvector of the data’s covariance matrix).

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

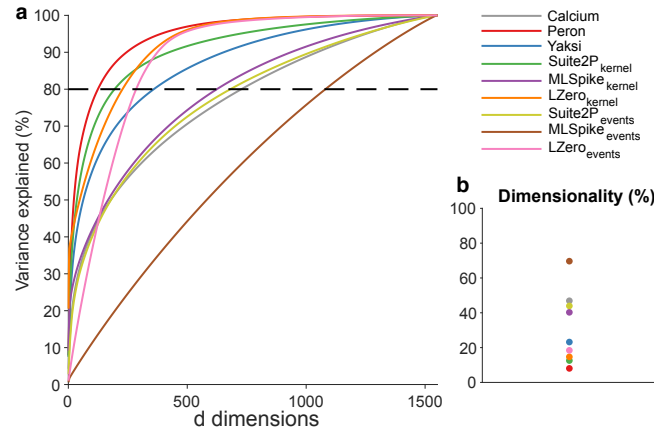


Figure 8: Dimensionality of population activity.

(a) Cumulative variance explained by each dimension of the data’s covariance matrix, one line per deconvolution method. Dashed line is the 80% threshold used in panel (b).

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

373 3 Discussion

374 Imaging of somatic calcium is a remarkable tool for capturing the simultaneous activity of
375 hundreds to thousands of neurons. But the time-series of each neuron’s calcium fluo-
376 rescence is inherently noisy and non-linearly related to its spiking. We sought here to address
377 how our choice of corrections to these time-series – to use them raw, deconvolve them into
378 continuous time-series, or deconvolve them into discrete events – affect the quality and
379 reliability of the scientific inferences drawn. Our approach was to replicate the process of
380 a typical population calcium-imaging study: choose an algorithm, choose its parameters
381 using some reasonable heuristics, and analyse the resulting time-series.

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

394 3.1 Accurate discrete deconvolution is possible, but sensitive

395 We find much that is encouraging. In fitting discrete deconvolution methods to ground-
396 truth data, we found they can in principle accurately recover known spike-times from raw
397 calcium time-series. A caveat here is that the choice of metric for evaluation and fitting
398 of parameters is of critical importance. The widely-used Pearson correlation coefficient is
399 a poor choice of metric as it returns inconsistent results with small changes in algorithm
400 parameters, and leads to poor estimates of simple measures such as firing rate when used
401 across methods and sampling rates. By contrast, the Error Rate metric (Deneux et al.,

402 2016; Victor and Purpura, 1996) resulted in excellent recovery of ground-truth spike trains.
403 Other recently developed methods for comparing spike-trains based on information theory
404 (Theis et al., 2016) or fuzzy set theory (Reynolds et al., 2018), may also be appropriate.

405 However, while good estimates of ground-truth spike times can be achieved with mod-
406 ern discrete deconvolution methods (Berens et al., 2018; Pachitariu et al., 2018), the best
407 parameters vary substantially between cells, and small changes in parameter values result
408 in poor performance. This variation and sensitivity of parameters played out as widely-
409 differing results between the three discrete deconvolution methods in analyses of neural
410 activity, coding, and correlation structure.

411 3.2 Choosing parameters for deconvolution methods

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

422 3.3 Ways forward

423 How then to solve the problem of the wide disagreements we report here, both between the
424 raw calcium and the deconvolved time-series, and between the outputs of the deconvolution
425 methods?

426 The simplest approach is to side-step the issue, and just use the raw calcium time-
427 series. Many studies use the raw calcium signal as the basis for all their analyses (Harvey
428 et al., 2012; Huber et al., 2012; Chu et al., 2016), perhaps assuming this is the least biased
429 approach. Our results suggest caution: the discrepancy between the raw and deconvolved
430 calcium on single neuron coding suggests an extraordinary range of possible results, from
431 about half of all neurons tuned to the task down to less 5 percent. The qualitative
432 conclusion – there is coding – is not satisfactory. Moreover, as noted by (Sabatini, 2019),
433 the raw calcium fluorescence signal is a low-pass filtered version of the underlying spike
434 train, which places strong limits on the maximum correlation between the raw signal
435 and underlying spikes, and hence on any correlations between the raw signal and the
436 behavioural variables related to those spikes. Indeed, the desire for better recovery of the
437 spikes and their correlations with behaviour is one of the principle reasons for developing
438 deconvolution methods.

439 A natural step then is to improve deconvolution methods with better forward models,
440 like MLSpike, for the link from spiking to calcium fluorescence (Greenberg et al., 2018).
441 Indeed, as sensors with faster kinetics (though fundamentally limited by kinetics of calcium
442 release itself) and higher signal-to-noise ratios are developed (Badura et al., 2014; Dana
443 et al., 2016, 2019), so the accuracy and robustness of de-noising and deconvolution should
444 improve; and as the neuron yield continues to increase (Ahrens et al., 2013; Stringer et al.,
445 2019a), so the potential for insights from inferred spikes or spike-driven events grows.
446 Developing further advanced deconvolution algorithms will harness these advances, but
447 are potentially always limited by the lack of ground-truth to fit their parameters (Wei

448 [et al., 2019](#)). Worse, no matter how good the forward model for a single neuron, our
449 results suggest the wide variation in the model parameters needed for each neuron would
450 make population analyses challenging to interpret.

451 A simple alternative approach to the inconsistencies between different forms of decon-
452 volved time-series is to triangulate them, and take the consensus across their results. For
453 example, our finding of a set of tuned neurons across multiple methods is strong evidence
454 that neurons in L2/3 of barrel cortex are responsive across the stages of the decision task.
455 Further examples of such triangulation in the literature are rare; Klaus and colleagues
456 ([Klaus et al., 2017](#)) used two different pipelines to derive raw $\Delta F/F$ of individual neurons
457 from one-photon fibre-optic recordings in the striatum, and replicated all analyses using
458 the output of both pipelines. Our results encourage the further use of triangulation to
459 create robust inference: obtaining the same result in the face of wide variation increases
460 our belief in its reliability ([Munafò and Davey Smith, 2018](#)).

461 There are caveats to triangulating by using a full consensus across three or more
462 versions of the time series. For single neuron analyses, such a full consensus inevitably
463 comes at the price of reducing the yield of neurons to which we can confidently assign roles.
464 There is also an assumption that all contributions to the consensus contain useful data:
465 if one deconvolution method returns time-series with no relation to the underlying spike
466 events, then including its outputs in the consensus would inevitably worsen the results.
467 An alternative version of triangulation partially circumventing these problems would be
468 to separately take the consensus between the raw calcium time-series and each of two
469 or more spike-event deconvolution methods, and then combine the results. Future work
470 on triangulation approaches would also need to look at how to combine more complex
471 analyses than single neuron properties, such as pairwise correlations.

472 Another approach, little explored to date, would be to use data constraints to tune
473 the deconvolution algorithm parameters. One option would be to use known properties
474 of neural activity in a recorded population as constraints. We showed, for example, that
475 some deconvolution methods did not recover the expected population-wide distribution of
476 activity in layer 2/3 of barrel cortex; so constraining all algorithms to reproduce the long-
477 tailed activity distribution may improve agreement between them in measures of coding
478 and correlation. Another option would be to tune deconvolution parameters to maximise
479 consistency within the deconvolved data. For example, [Pachitariu et al. \(2018\)](#) recently
480 proposed maximising the correlation between deconvolved traces from the same neuron
481 obtained between trials of the same visual stimulus. Such an approach needs a suitable
482 task design to ensure consistent conditions within which to compare responses of the same
483 neuron (such as identical duration repeats of identical visual stimuli) – and which therefore
484 could not be applied to the pole-detection task considered here. It would also require that
485 the known variations in a neuron’s response between repeats of the same task condition
486 or stimulus is not large enough to prevent meaningful correlations between repeats.

487 Our results provide impetus for different directions of research, not just to improving
488 our modelling of the relationship between spikes and the somatic calcium signal, but
489 also focussing on how we can verify results across the output of different deconvolution
490 algorithms, and thus provide robust scientific inferences about neural populations.

491 4 Methods

492 Ground truth data

493 Ground truth data were accessed from crcns.org (Svoboda, 2015), and the experiments
494 have been described previously (Chen et al., 2013). Briefly, mouse visual cortical neu-
495 rons expressing the fluorescent calcium reporter protein GCaMP6s were imaged with two-
496 photon microscopy at 60Hz. Loose-seal cell-attached recordings were performed simul-
497 taneously at 10kHz. Recordings were made in awake mice during 5 trials (4s blank, 4s
498 stimulus) of the optimal moving grating stimulus (1 of 8 directions) for the cell-attached
499 neuron. The data-set contains twenty one recordings from nine neurons.

500 Neuropil subtraction was performed as described in Chen et al. (2013), based on ex-
501 ample code provided alongside the data at crcns.org. Neuropil signals – defined as the
502 average fluorescence from all pixels within a 20 μm radius from each cell centre excluding
503 the region of interest (ROI) – were subtracted from cell fluorescence in a weighted fashion,
504 $F_{corrected} = F_{cell} - 0.7F_{neuropil}$.

505 Population imaging data description

506 Population imaging data was accessed from crcns.org and have been described previ-
507 ously (Peron et al., 2015b). Briefly, volumetric two photon calcium imaging of primary
508 somatosensory cortex (S1) was performed in awake head-fixed mice performing a whisker-
509 based object localisation task. In the task a metal pole was presented in one of two loca-
510 tions and mice were motivated with fluid reward to lick at one of two lick ports depending
511 on the location of the pole following a brief delay. Two photon imaging of GCaMP6s
512 expressing neurons in superficial S1 was performed at 7Hz. Images were motion corrected
513 and aligned, before regions of interest were manually set and neuropil-subtracted. A single
514 recording from this dataset was used for population analysis. The example session had
515 1552 neurons recorded for a total of 23559 frames (56 minutes).

516 List of deconvolution methods

517 MLSpike

518 MLSpike (Deneux et al., 2016) was accessed from <https://github.com/mlspike>. MLSpike
519 uses a model-based probabilistic approach to recover spike trains in calcium imaging data
520 by taking baseline fluctuations and cellular properties into account. Briefly, MLSpike
521 implements a model of measured calcium fluorescence as a combination of spike-induced
522 transients, background (photonic) noise and drifting baseline fluctuations. A maximum
523 likelihood approach determines the probability of the observed calcium at each time step
524 given an inferred spike train generated through a particular set of model parameters.
525 MLSpike returns a maximum a posteriori spike train (as used here), or a spike probability
526 per time step.

527 MLSpike has a number of free parameters, of which we optimise three: A , the mag-
528 nitude of fluorescence transients caused by a single spike; τ , calcium fluorescence decay
529 time; σ , background (photonic) noise level. MLSpike also has parameters for different
530 calcium sensor kinetics (for OBG, GCaMP3, GCaMP6 and so on) which we fix to default
531 values for GCaMP6.

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

535 Suite2P

536 Suite2P (Pachitariu et al., 2016, 2018) was accessed from [https://github.com/cortex-](https://github.com/cortex-lab/Suite2P)
537 [lab/Suite2P](https://github.com/cortex-lab/Suite2P). Suite2P was developed as a complete end-to-end processing pipeline for
538 large scale 2-photon imaging analysis - from image registration to spike extraction and
539 visualization - of which we only use the spike extraction step. The spike deconvolution
540 of Suite2P uses a sparse non-negative deconvolution algorithm, greedily identifying and
541 removing calcium transients to minimise the cost function

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

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

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

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

555 LZero

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

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

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

572 Yaksi

573 Yaksi is an implementation of the deconvolution approach of Yaksi and Friedrich (2006).
574 The fluorescence time series is low-pass filtered (4th order butterworth filter, 0.7Hz cutoff)
575 to remove noise before having a calcium kernel (exponential decay of 2 seconds, as used
576 in Suite2P and LZero above) linearly deconvolved out of the signal using Matlab's `deconv`

577 function. The output of Yaksi is a continuous signal approximating spike density per unit
578 time.

579 **Peron events**

580 **Peron events** refer to the de-noised calcium event traces detailed in the original [Peron
581 et al. \(2015b\)](#) paper. Here a version of the ‘peeling’ algorithm ([Lütcke et al., 2013](#)) was
582 developed, a template-fitting algorithm with variable decay time constants across events
583 and neurons. The output for analysis is a continuous signal approximating de-noised
584 calcium concentration per unit time.

585 *Events and kernel versions of spike inference methods*

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

592 **Ground truth spike train metrics**

593 Pearson correlation coefficient was computed between the ground truth and inferred spikes
594 (MLSpike) or events (Suite2P, LZero) following convolution of both with a gaussian kernel
595 (61 samples wide, 1.02 seconds).

596 Error Rate was computed between the ground truth and inferred spikes/events using
597 the [Deneux et al. \(2016\)](#) implementation of normalised error rate, derived from [Victor and
598 Purpura \(1996\)](#) – code available <https://github.com/MLspike>. Briefly, the error rate is 1
599 - F1-score, where the F1-score is the harmonic mean of sensitivity and precision ([Davis
600 and Goadrich, 2006](#)),

$$sensitivity = 1 - \frac{misses}{total\ spikes},$$

$$precision = 1 - \frac{false\ detections}{total\ detections},$$

$$ErrorRate = 1 - 2 \frac{sensitivity \times precision}{sensitivity + precision}.$$

601 Hits, misses and false detections were counted with a temporal precision of 0.5 seconds.

602 For normalised estimation of errors in firing/event rate we compute,

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

603 where spike/event rates are measured in Hz.

604 **Parameter fitting**

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

609 21 values), sigma (0.01:1, 21 values); Suite2P Threshold (0.1:100, 13 values); LZero lambda
610 (0.1:20, 23 values), scale (0.1:20, 23 values).

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

614 Downsampling

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

618 4.1 Event rate estimation

619 Spike inference methods (Suite2P_{events}, MLSpikes_{events}, LZero_{events}) return estimated spike
620 times (MLSpikes), or event times (Suite2P/LZero) which were converted into mean event
621 rates (Hz) per neuron.

622 The event rate for continuous methods (Calcium, Peron, Yaksi, Suite2P_{kernel}, MLSpikes_{kernel},
623 LZero_{kernel}) for each neuron was determined by counting activity/fluorescence transients
624 greater than three standard deviations of the background noise. Background noise was
625 calculated by subtracting a four-frame moving average of the fluorescence from the raw
626 data to result in a ‘noise only’ trace. This operation was done separately for each neuron
627 and each method. Event rate was then computed in Hz.

628 Silent neurons were defined as neurons with event rates below 0.0083Hz (or fewer than
629 one spike per two minutes of recording) as in [O’Connor et al. \(2010\)](#).

630 4.2 Task-tuned neurons

631 Task-tuning was determined for each neuron using the model-free approach of [Peron et al.](#)
632 [\(2015b\)](#). Neurons were classed as task-tuned if their peak trial-average activity exceeded
633 the 95th percentile of a distribution of trial-average peaks from shuffled data (10000 shuffles
634 of time-series order). The shuffle test was done separately for correct lick-left and lick-right
635 trials and neurons satisfying the tuning criteria in either case were counted as task-tuned.

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

639 4.3 Touch-tuned responses

640 Touch-tuned neurons were determined by first computing touch-triggered average activity
641 for each neuron, then calculating whether the data distribution of peak touch-induced
642 activity exceeds the expected activity of resampled data. In more detail, the time of first
643 touch between the mouse’s whisker and the metal pole on each trial was recorded. For
644 each neuron, one second of activity (seven data samples) was extracted before and after
645 the frame closest to the first touch of each trial (15 frames total per trial); taking the
646 mean touch-triggered activity over trials gave the average touch response for the neuron.
647 To determine whether the neuron was touch tuned or not, we compared the neuron’s
648 peak mean response r_{data} to a null distribution by taking a randomly sampled 15 frame
649 segment of a trial, finding the peak mean response across trials r_{null} , and repeating this
650 calculation for 10000 random samples. A p-value for the data peak response was calculated

651 as $p = \#\{r_{null} < r_{data}\}/10000$. Over all neurons, a neuron was considered touch-tuned if
652 $p < 0.05$ after Benjamini-Hochberg correction.

653 4.4 Pairwise correlations

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

656 Stability of correlation estimates (Fig. 7b) at the recording durations used was assessed
657 by computed the similarity between correlation distributions for the the intact dataset to
658 those from subsets of the dataset. For each deconvolution method, we computed the
659 pairwise correlation matrix using the entire session’s data, as above. We also sampled a
660 subset of time-points (1%-100%) of the full dataset at random without replacement and
661 computed a matrix of pairwise correlations for this subset. We then compute the similarity
662 between the total and subset matrices using Pearsons correlation coefficient. This process
663 was repeated 100 times and the mean (line) and standard deviation (shading) of the 100
664 repeats were plotted.

665 4.5 Correlations between correlation matrices

666 Correlations between correlation matrices (Fig. 7c-e) were computed using Spearman’s
667 rank correlation between the unique pairwise correlations from each method (i.e. the
668 upper triangular entries of the correlation matrix).

669 4.6 Dimensionality

670 To determine the dimensionality of each dataset we performed eigendecomposition of the
671 covariance matrix of each dataset. The resultant eigenvalues were sorted into descending
672 order $\lambda_1 \geq \lambda_2 \geq \dots \lambda_N$, and the cumulative variance explained by d dimensions computed
673 as $\sum_{i=1}^d \lambda_i / \sum_{i=1}^N \lambda_i$.

674 References

675 Misha B Ahrens, Michael B Orger, Drew N Robson, Jennifer M Li, and Philipp J Keller.
676 Whole-brain functional imaging at cellular resolution using light-sheet microscopy. *Nat.*
677 *Methods*, 10(5):413–420, May 2013.

678 Aleksandra Badura, Xiaonan Richard Sun, Andrea Giovannucci, Laura A Lynch, and
679 Samuel S-H Wang. Fast calcium sensor proteins for monitoring neural activity. *Neu-*
680 *rophotonics*, 1(2):025008, October 2014.

681 Philipp Berens, Jeremy Freeman, Thomas Deneux, Nicolay Chenkov, Thomas McColgan,
682 Artur Speiser, Jakob H Macke, Srinivas C Turaga, Patrick Mineault, Peter Rupprecht,
683 Stephan Gerhard, Rainer W Friedrich, Johannes Friedrich, Liam Paninski, Marius Pa-
684 chitariu, Kenneth D Harris, Ben Bolte, Timothy A Machado, Dario Ringach, Jasmine
685 Stone, Luke E Rogerson, Nicolas J Sofroniew, Jacob Reimer, Emmanouil Froudarakis,
686 Thomas Euler, Miroslav Román Rosón, Lucas Theis, Andreas S Tolia, and Matthias
687 Bethge. Community-based benchmarking improves spike rate inference from two-photon
688 calcium imaging data. *PLoS Comput. Biol.*, 14(5):e1006157, May 2018.

689 K L Briggman, H D I Abarbanel, and W B Kristan, Jr. Optical imaging of neuronal
690 populations during decision-making. *Science*, 307(5711):896–901, February 2005.

- 691 Emery N Brown, Robert E Kass, and Partha P Mitra. Multiple neural spike train data
692 analysis: state-of-the-art and future challenges. *Nat. Neurosci.*, 7(5):456–461, May 2004.
- 693 Lawrence D. Brown, T. Tony Cai, and Anirban DasGupta. Interval estimation for a
694 binomial proportion. *Statist Sci*, 16:101–133, 2001.
- 695 J K Chapin and M A Nicolelis. Principal component analysis of neuronal ensemble activity
696 reveals multidimensional somatosensory representations. *J. Neurosci. Methods*, 94(1):
697 121–140, December 1999.
- 698 Tsai-Wen Chen, Trevor J Wardill, Yi Sun, Stefan R Pulver, Sabine L Renninger, Amy
699 Baohan, Eric R Schreiter, Rex A Kerr, Michael B Orger, Vivek Jayaraman, Loren L
700 Looger, Karel Svoboda, and Douglas S Kim. Ultrasensitive fluorescent proteins for
701 imaging neuronal activity. *Nature*, 499(7458):295–300, July 2013.
- 702 Monica W Chu, Wankun L Li, and Takaki Komiyama. Balancing the robustness and
703 efficiency of odor representations during learning. *Neuron*, 92:174–186, Oct 2016. ISSN
704 1097-4199. doi: 10.1016/j.neuron.2016.09.004.
- 705 Mark M Churchland, John P Cunningham, Matthew T Kaufman, Justin D Foster, Paul
706 Nuyujukian, Stephen I Ryu, and Krishna V Shenoy. Neural population dynamics during
707 reaching. *Nature*, 487(7405):51–56, July 2012.
- 708 John P Cunningham and Byron M Yu. Dimensionality reduction for large-scale neural
709 recordings. *Nat. Neurosci.*, 17(11):1500–1509, November 2014.
- 710 Hod Dana, Tsai-Wen Chen, Amy Hu, Brenda C. Shields, Caiying Guo, Loren L. Looger,
711 Douglas S. Kim, and Karel Svoboda. Thy1-gcamp6 transgenic mice for neuronal popula-
712 tion imaging in vivo. *PLoS One*, 9(9):e108697, 2014. doi: 10.1371/journal.pone.0108697.
713 URL <http://dx.doi.org/10.1371/journal.pone.0108697>.
- 714 Hod Dana, Boaz Mohar, Yi Sun, Sujatha Narayan, Andrew Gordus, Jeremy P Hasseman,
715 Getahun Tsegaye, Graham T Holt, Amy Hu, Deepika Walpita, Ronak Patel, John J
716 Macklin, Cornelia I Bargmann, Misha B Ahrens, Eric R Schreiter, Vivek Jayaraman,
717 Loren L Looger, Karel Svoboda, and Douglas S Kim. Sensitive red protein calcium
718 indicators for imaging neural activity. *Elife*, 5, March 2016.
- 719 Hod Dana, Yi Sun, Boaz Mohar, Brad K Hulse, Aaron M Kerlin, Jeremy P Hasseman,
720 Getahun Tsegaye, Arthur Tsang, Allan Wong, Ronak Patel, John J Macklin, Yang
721 Chen, Arthur Konnerth, Vivek Jayaraman, Loren L Looger, Eric R Schreiter, Karel
722 Svoboda, and Douglas S Kim. High-performance calcium sensors for imaging activity
723 in neuronal populations and microcompartments. *Nat. Methods*, 16(7):649–657, July
724 2019.
- 725 Jesse Davis and Mark Goadrich. The relationship between Precision-Recall and ROC
726 curves. In *Proceedings of the 23rd International Conference on Machine Learning, ICML*
727 '06, pages 233–240, New York, NY, USA, 2006. ACM.
- 728 Thomas Deneux, Attila Kaszas, Gergely Szalay, Gergely Katona, Tamás Lakner, Amiram
729 Grinvald, Balázs Rózsa, and Ivo Vanzetta. Accurate spike estimation from noisy calcium
730 signals for ultrafast three-dimensional imaging of large neuronal populations in vivo.
731 *Nat. Commun.*, 7:12190, July 2016.

- 732 Johannes Friedrich, Pengcheng Zhou, and Liam Paninski. Fast online deconvolution of
733 calcium imaging data. *PLoS Comput. Biol.*, 13(3):e1005423, March 2017.
- 734 Andrea Giovannucci, Johannes Friedrich, Pat Gunn, Jérémie Kalfon, Brandon L Brown,
735 Sue Ann Koay, Jiannis Taxidis, Farzaneh Najafi, Jeffrey L Gauthier, Pengcheng Zhou,
736 Baljit S Khakh, David W Tank, Dmitri B Chklovskii, and Eftychios A Pnevmatikakis.
737 CaImAn an open source tool for scalable calcium imaging data analysis. *Elife*, 8, January
738 2019.
- 739 David S Greenberg, Damian J Wallace, Kay-Michael Voit, Silvia Wuertenberger, Uwe
740 Czubayko, Arne Monsees, Takashi Handa, Joshua T Vogelstein, Reinhard Seifert,
741 Yvonne Groemping, and Jason ND Kerr. Accurate action potential inference from a cal-
742 cium sensor protein through biophysical modeling. *bioRxiv*, 2018. doi: 10.1101/479055.
743 URL <https://www.biorxiv.org/content/early/2018/11/29/479055>.
- 744 Kenneth D Harris, Rodrigo Quian Quiroga, Jeremy Freeman, and Spencer L Smith. Im-
745 proving data quality in neuronal population recordings. *Nat. Neurosci.*, 19(9):1165–
746 1174, August 2016.
- 747 Christopher D Harvey, Philip Coen, and David W Tank. Choice-specific sequences in
748 parietal cortex during a virtual-navigation decision task. *Nature*, 484(7392):62–68, April
749 2012.
- 750 Samuel Andrew Hires, Diego A Gutnisky, Jianing Yu, Daniel H O’Connor, and Karel
751 Svoboda. Low-noise encoding of active touch by layer 4 in the somatosensory cortex.
752 *Elife*, 4, August 2015.
- 753 Daniel Huber, D A Gutnisky, S Peron, D H O’Connor, J S Wiegert, L Tian, T G Oertner,
754 L L Looger, and K Svoboda. Multiple dynamic representations in the motor cortex
755 during sensorimotor learning. *Nature*, 484(7395):473–478, April 2012.
- 756 Mark D. Humphries. Dynamical networks: Finding, measuring, and tracking neural pop-
757 ulation activity using network science. *Network Neuroscience*, 1:324–338, 2017. doi:
758 10.1162/NETN_a.00020.
- 759 Sean Jewell and Daniela Witten. Exact spike train inference via ℓ_0 optimization. *Ann.*
760 *Appl. Stat.*, 12(4):2457–2482, December 2018.
- 761 Patrick Kaifosh, Jeffrey D Zaremba, Nathan B Danielson, and Attila Losonczy. SIMA:
762 Python software for analysis of dynamic fluorescence imaging data. *Front. Neuroinform.*,
763 8:80, September 2014.
- 764 Saul Kato, Harris S Kaplan, Tina Schrödel, Susanne Skora, Theodore H Lindsay, Eviatar
765 Yemini, Shawn Lockery, and Manuel Zimmer. Global brain dynamics embed the motor
766 command sequence of *Caenorhabditis elegans*. *Cell*, 163(3):656–669, October 2015.
- 767 Sander W Keemink, Scott C Lowe, Janelle M P Pakan, Evelyn Dylida, Mark C W van
768 Rossum, and Nathalie L Rochefort. FISSA: A neuropil decontamination toolbox for
769 calcium imaging signals. *Sci. Rep.*, 8(1):3493, February 2018.
- 770 Andreas Klaus, Gabriela J Martins, Vitor B Paixao, Pengcheng Zhou, Liam Paninski, and
771 Rui M Costa. The spatiotemporal organization of the striatum encodes action space.
772 *Neuron*, 95(5):1171–1180.e7, August 2017.

- 773 Dmitry Kobak, Wieland Brendel, Christos Constantinidis, Claudia E Feierstein, Adam
774 Kepecs, Zachary F Mainen, Xue-Lian Qi, Ranulfo Romo, Naoshige Uchida, and Chris-
775 tian K Machens. Demixed principal component analysis of neural population data. *Elife*,
776 5, April 2016.
- 777 Henry Lütcke, Felipe Gerhard, Friedemann Zenke, Wulfram Gerstner, and Fritjof Helm-
778 chen. Inference of neuronal network spike dynamics and topology from calcium imaging
779 data. *Front. Neural Circuits*, 7:201, December 2013.
- 780 Eran A Mukamel, Axel Nimmerjahn, and Mark J Schnitzer. Automated analysis of cellular
781 signals from large-scale calcium imaging data. *Neuron*, 63(6):747–760, September 2009.
- 782 Marcus R Munafò and George Davey Smith. Robust research needs many lines of evidence.
783 *Nature*, 553(7689):399–401, January 2018.
- 784 Daniel H O’Connor, Simon P Peron, Daniel Huber, and Karel Svoboda. Neural activity
785 in barrel cortex underlying vibrissa-based object localization in mice. *Neuron*, 67(6):
786 1048–1061, September 2010.
- 787 Marius Pachitariu, Carsen Stringer, Sylvia Schröder, Mario Dipoppa, L Federico Rossi,
788 Matteo Carandini, and Kenneth D Harris. Suite2p: beyond 10, 000 neurons with stan-
789 dard two-photon microscopy. *BioRxiv*, Preprint at <http://dx.doi.org/10.1101/061507>,
790 2016.
- 791 Marius Pachitariu, Carsen Stringer, and Kenneth D Harris. Robustness of spike deconvo-
792 lution for neuronal calcium imaging. *J. Neurosci.*, August 2018.
- 793 António R C Paiva, Il Park, and José C Príncipe. A comparison of binless spike train
794 measures. *Neural Comput. Appl.*, 19(3):405–419, April 2010.
- 795 Simon Peron, Tsai-Wen Chen, and Karel Svoboda. Comprehensive imaging of cortical
796 networks. *Curr. Opin. Neurobiol.*, 32:115–123, June 2015a.
- 797 Simon P Peron, Jeremy Freeman, Vijay Iyer, Caiying Guo, and Karel Svoboda. A cellular
798 resolution map of barrel cortex activity during tactile behavior. *Neuron*, 86(3):783–799,
799 May 2015b.
- 800 Eftychios A Pnevmatikakis, Daniel Soudry, Yuanjun Gao, Timothy A Machado, Josh
801 Merel, David Pfau, Thomas Reardon, Yu Mu, Clay Lacefield, Weijian Yang, Misha
802 Ahrens, Randy Bruno, Thomas M Jessell, Darcy S Peterka, Rafael Yuste, and Liam
803 Paninski. Simultaneous denoising, deconvolution, and demixing of calcium imaging
804 data. *Neuron*, January 2016.
- 805 Ruben Portugues, Claudia E Feierstein, Florian Engert, and Michael B Orger. Whole-
806 brain activity maps reveal stereotyped, distributed networks for visuomotor behavior.
807 *Neuron*, 81(6):1328–1343, March 2014.
- 808 Stephanie Reynolds, Therese Abrahamsson, Per Jesper Sjöström, Simon R Schultz, and
809 Pier Luigi Dragotti. CosMIC: A consistent metric for spike inference from calcium
810 imaging. *Neural Comput.*, 30(10):2726–2756, October 2018.
- 811 Oleg I Rumyantsev, Jrme A Lecoq, Oscar Hernandez, Yanping Zhang, Joan Savall, Ra-
812 dos?aw Chrapkiewicz, Jane Li, Hongkui Zeng, Surya Ganguli, and Mark J Schnitzer.
813 Fundamental bounds on the fidelity of sensory cortical coding. *Nature*, 580:100–105,
814 April 2020. ISSN 1476-4687. doi: 10.1038/s41586-020-2130-2.

- 815 Bernardo Sabatini. The impact of reporter kinetics on the interpretation of data gathered
816 with fluorescent reporters. *bioRxiv*, page 834895, 2019. doi: 10.1101/834895. URL
817 <https://www.biorxiv.org/content/early/2019/11/07/834895>.
- 818 Carsen Stringer and Marius Pachitariu. Computational processing of neural recordings
819 from calcium imaging data. *Curr. Opin. Neurobiol.*, 55:22–31, December 2018.
- 820 Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Matteo Carandini, and Ken-
821 neth D Harris. High-dimensional geometry of population responses in visual cortex.
822 *Nature*, 571:361–365, July 2019a. ISSN 1476-4687. doi: 10.1038/s41586-019-1346-5.
- 823 Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Charu Bai Reddy, Matteo Caran-
824 dini, and Kenneth D Harris. Spontaneous behaviors drive multidimensional, brainwide
825 activity. *Science*, 364:255, April 2019b. ISSN 1095-9203. doi: 10.1126/science.aav7893.
- 826 K Svoboda. Simultaneous imaging and loose-seal cell-attached electrical recordings from
827 neurons expressing a variety of genetically encoded calcium indicators. *GENIE project*,
828 *Janelia Farm Campus, HHMI; CRCNS.org*, 2015.
- 829 Lucas Theis, Philipp Berens, Emmanouil Froudarakis, Jacob Reimer, Miroslav
830 Román Rosón, Tom Baden, Thomas Euler, Andreas S Tolias, and Matthias Bethge.
831 Benchmarking spike rate inference in population calcium imaging. *Neuron*, 90(3):471–
832 482, May 2016.
- 833 J D Victor and K P Purpura. Nature and precision of temporal coding in visual cortex:
834 a metric-space analysis. *J. Neurophysiol.*, 76(2):1310–1326, August 1996.
- 835 Joshua T Vogelstein, Adam M Packer, Timothy A Machado, Tanya Sippy, Baktash Babadi,
836 Rafael Yuste, and Liam Paninski. Fast nonnegative deconvolution for spike train infer-
837 ence from population calcium imaging. *J. Neurophysiol.*, 104(6):3691–3704, December
838 2010.
- 839 Ziqiang Wei, Bei-Jung Lin, Tsai-Wen Chen, Kayvon Daie, Karel Svoboda, and Shaul
840 Druckmann. A comparison of neuronal population dynamics measured with calcium
841 imaging and electrophysiology. *bioRxiv*, 2019. doi: 10.1101/840686. URL <https://www.biorxiv.org/content/early/2019/11/15/840686>.
- 843 Adrien Wohrer, Mark D Humphries, and Christian K Machens. Population-wide distri-
844 butions of neural activity during perceptual decision-making. *Prog. Neurobiol.*, 103:
845 156–193, April 2013.
- 846 Emre Yaksi and Rainer W Friedrich. Reconstruction of firing rate changes across neuronal
847 populations by temporally deconvolved ca²⁺ imaging. *Nat. Methods*, 3(5):377–383, May
848 2006.