

On the use of calcium deconvolution algorithms in practical contexts

Mathew H. Evans, Rasmus S. Petersen & Mark D. Humphries [add affiliations]

May 24, 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 methods are available to recover spike-evoked events from the raw calcium, from simple deconvolution to inferring the spikes themselves. Here we explore the extent to which our choice of raw or deconvolved calcium time-series affects the scientific inferences we can draw. Our results show the choice qualitatively changes the potential scientific inferences we draw about neural activity, coding, and correlation structure. 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 cells tuned to task features; yet there is also qualitative disagreement between deconvolution methods on which neurons are tuned. 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. We suggest that quantitative results obtained from population calcium-imaging be verified across multiple forms of the calcium time-series.

1 Introduction

Calcium imaging is a wonderful tool for high yield recordings of large neural populations [Harris; Stringer; Ahrens; Orger; others]. 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 [cite loads; including van Rossum latest Sci Report paper].

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 [cite Theis benchmarking; Stringer Curr Opin], 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: how do our scientific inferences depend on our choice of raw or deconvolved calcium time-series.

Deconvolution algorithm themselves range in complexity from simple deconvolution with a fixed kernel of the calcium response [Yaksi], through detecting spike-evoked calcium events [LZero, Suite2p], to directly inferring spike times [MLspike, Peeling]. This continuum of options raise the further question of the extent to which we should process the raw calcium signals.

37 We proceed in two stages. In order to use deconvolution algorithms, we need to choose
38 their parameters. We'd like to know whether it is worth taking this extra step: how good
39 can these algorithms be in principle, and how sensitive their results are to the choice of
40 parameter values. We thus first evaluate qualitatively different deconvolution algorithms,
41 by optimising their parameters against ground truth data with known spikes. With our
42 understanding of their parameters in hand, we then turn to our main question, by analysing
43 a large-scale population recording from the barrel cortex of a mouse performing a whisker-
44 based decision task. We compare the scientific inferences about population coding and
45 correlations we obtain using either raw calcium signals, or a range of time-series derived
46 from those calcium signals, covering simple deconvolution, event detection, and spikes.

47 We find contrasting answers. A substantial fraction of the methods used here fail
48 to recover basic features of population activity in barrel cortex established from electro-
49 physiology. The inferences we draw about coding qualitatively differ between raw and decon-
50 volved calcium signals. In particular, coding analyses based on raw calcium signals
51 detect an order of magnitude more cells tuned to task features. Yet there is also qualitative
52 disagreement between deconvolution methods on which neurons are tuned. The inferences
53 we draw about correlations between neurons do not distinguish between raw and decon-
54 volved calcium signals, but can qualitatively differ between deconvolution methods. Our
55 results thus suggest care is needed in drawing inferences from population recordings of so-
56 matic calcium, and that one solution is to replicate all results in both raw and deconvolved
57 calcium signals.

58 **2 Results**

59 **2.1 Performance of deconvolution algorithms on ground-truth data-sets**

60 We select here three deconvolution algorithms that infer discrete spike-like events, each an
61 example of the state of the art in qualitatively different approaches to the problem: Suite2p
62 ([Pachitariu et al.](#)), a peeling algorithm that matches a scalable kernel to the calcium signal
63 to detect spike-triggered calcium events; LZero ([Jewell and Witten, 2017](#)), a change-point
64 detection algorithm, which finds as events the step-like changes in the calcium signal the
65 imply spikes [[Check](#)]; and MLspike ([Deneux et al., 2016](#)), a forward model, which fits an
66 explicit model of the spike-to-calcium dynamics in order to find spike-evoked changes in
67 the calcium signal, and returns spike times. We emphasize that these methods were chosen
68 as exemplars of their approaches, and are each innovative takes on the problem; we are
69 not here critiquing individual methods, but using an array of methods to illustrate the
70 problems and decisions facing the experimentalist when using calcium imaging data.

71 We first ask if these deconvolution methods work well in principle. We fit the parame-
72 ters of each method to a data-set of 21 ground-truth recordings ([Chen et al., 2013](#)), where
73 the spiking activity of a cell is recorded simultaneously with 60 Hz calcium imaging using
74 high-signal-to-noise juxtacellular recording techniques (Figure 1a). To fit the parameters
75 for each recording, we sweep each method's parameter space to find the parameter value(s)
76 with the best match between the true and inferred spike train.

77 The best-fit parameters depend strongly on how we evaluate the match between true
78 and inferred spikes. The Pearson correlation coefficient between the true and inferred
79 spike train is a common choice ([Brown et al., 2004; Paiva et al., 2010; Theis et al., 2016;](#)
80 [Reynolds et al., 2017; Berens et al., 2018](#)), typically with both trains convolved with a
81 Gaussian kernel to allow for timing errors. However, we find that choosing parameters to
82 maximise the correlation coefficient can create notable errors. The inferred spike trains

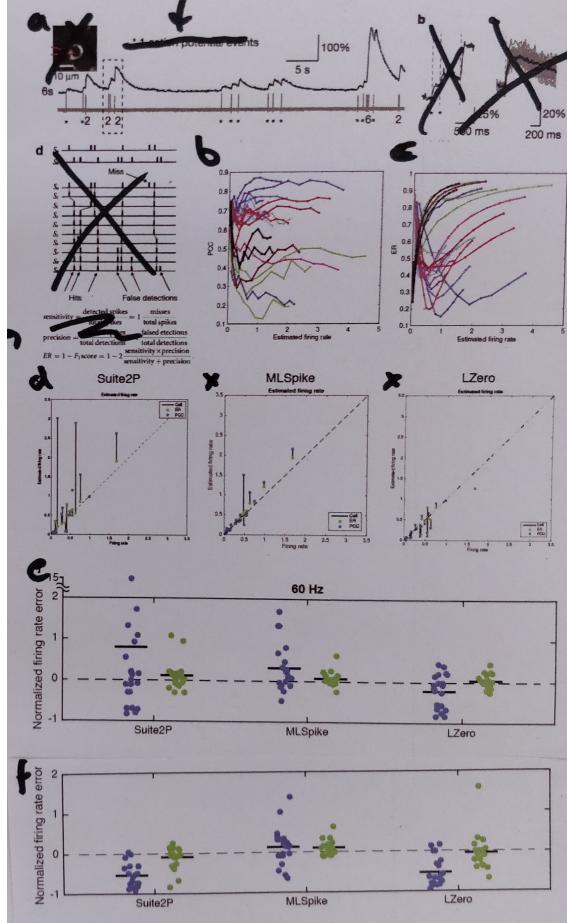


Figure 1: Ground truth data analysis.

- (a) Example simultaneous recording of somatic voltage and calcium activity imaged at 60Hz from ([Chen et al., 2013](#)). 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 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. One line per recording. Parameters: A : [explain]; τ [explain]; σ : explain [Full range on y axis: (0,1)] [Here and throughout: lose the box axes (upper and right)]
- (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. Left: correlation coefficient; right: Error Rate.

83 from ML Spike have too many spikes on average (mean error: 31.72%), and the accuracy of
84 recovered firing rates widely varies across recordings (Fig 1b, blue symbols). We attribute
85 these errors to the noisy relationship between the correlation coefficient and the number
86 of inferred spikes (Figure 1c): for many recordings, there is no well-defined maximum
87 coefficient, especially for the amplitude parameter A [Correct description?], so that near-
88 maximum correlation between true and inferred trains is consistent with a wide range of
89 spike counts in the inferred trains. We see the same sensitivity for the event rates from
90 recordings optimised using Suite2p and LZero (Figure 1f). If we compare their inferred
91 event rates to true firing rates (Fig 1b), we see Suite2p estimates far more events than
92 spikes (mean error 79.47%) and LZero fewer events than spikes (mean error: -21.14%).
93 These further errors are problematic: there cannot be more spike-driven calcium events
94 than spikes, and LZero's underestimate is considerably larger than the fraction of frames
95 with two or more spikes [I would guess: WHICH IS?].

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

110 A potential caveat here is that the ground-truth data are single neurons imaged at a
111 frame-rate of 60Hz, an order of magnitude greater than is typically achievable in popula-
112 tion recordings (Peron et al., 2015a). Such a high frame-rate could allow for more accurate
113 recovery of spikes than is possible in population recordings. To test this, we downsample
114 the ground-truth data to a 7Hz frame-rate, and repeat the parameter sweeps for each
115 deconvolution method applied to each recording. As we show in Figure 1c, optimising pa-
116 rameters using the minimum Error Rate still results in excellent recovery of the true spike
117 rate (and interestingly for some recordings reduces the error when using the correlation
118 coefficient). Lower frame-rates need not then be an impediment to using deconvolution
119 methods.

120 2.2 Parameters optimised on ground-truth are widely distributed and 121 sensitive

122 What might be an impediment to using deconvolution methods on population recordings
123 is that the best parameter values vary widely between cells. Figure 2a-b plots the best-fit
124 parameter values for each recording across deconvolution methods and sampling rates.
125 Each method has at least one parameter with substantial variability across recordings,
126 [SOMETHING HERE ON RANGE: are these ranges a high proportion of the total range
127 of that parameter?]. This suggests that the best parameters for one cell may perform
128 poorly for another cell.

129 The problem of between-cell variation in parameter values would be compensated
130 somewhat if the quality of the inferred spike or event trains is robust to changes in those

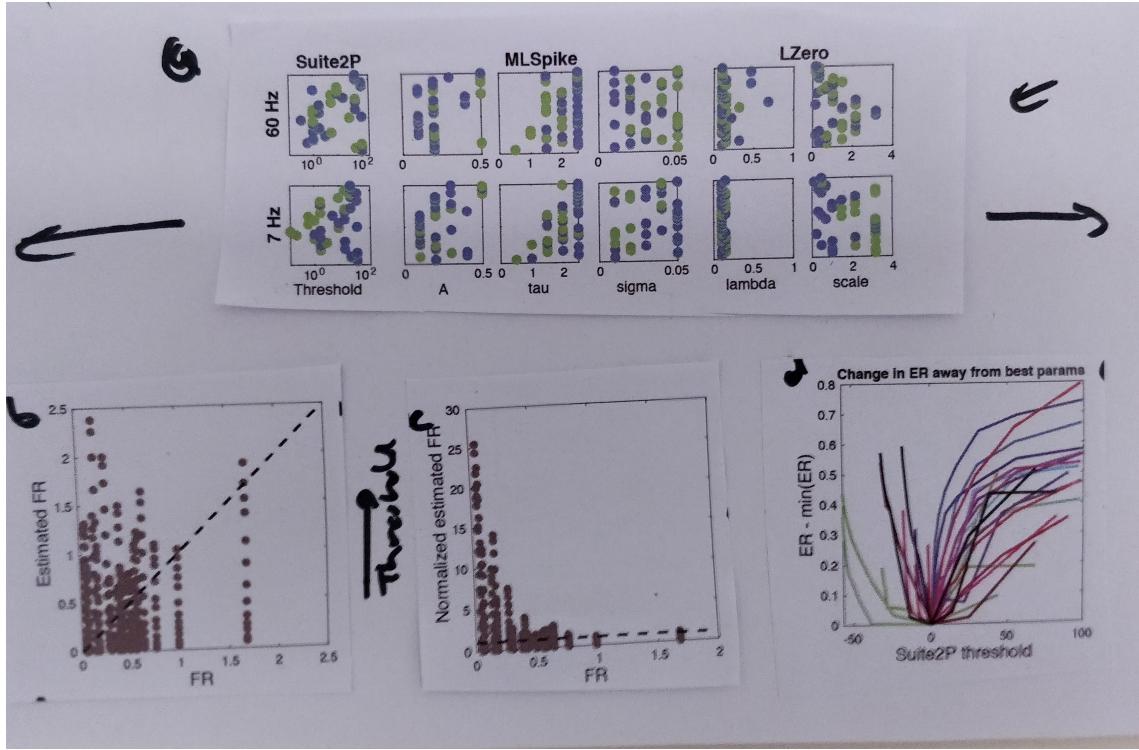


Figure 2: Variation in best-fit spike deconvolution parameters across ground-truth recordings.

(a) Distributions of optimised parameter values across recordings. In each panel, we plot parameter values on the x-axis against the recording ID on the y-axis (in an arbitrary but consistent order). Parameter values are plotted for those optimised using the Error Rate (green). Top row: fits to the original 60 Hz frame-rate data; bottom row: fits to data down-sampled to 7 Hz. [1: link recordings of the same cell] [2: add density histograms to top of each plot, to show distribution of parameters; one for PCC, one for ER]

(b) 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. [3 panels here]

(c) Change in the error rate with change in Suite2p's threshold value away from its optimum for each recording [label x-axis correctly "Change in threshold"]. One line per recording.

131 values. However, we find performance is highly sensitive to changes in some parameters.
 132 Figure ??b-c shows that for most recordings the quality of the inferred spike train abruptly
 133 worsens with small increases or decreases in the best parameter. Thus using deconvolution
 134 algorithms on population recordings comes with the potential issues that parameters can
 135 be both sensitive and vary considerably across cells.

136 2.3 Deconvolution of population imaging in barrel cortex during a de- 137 cision task

138 We turn now to seeing if and how these issues play out when analysing a large-scale
 139 population recording with no ground-truth. The data we use are two-photon calcium
 140 imaging time-series from a head-fixed mouse performing a whisker-based two-alternative
 141 decision task (Fig. 3a-b), from the study of (Peron et al., 2015a). We analyse here a single
 142 session with 1552 simultaneously recorded pyramidal neurons in L2/3 of a single barrel
 143 in somatosensory cortex, imaged at 7 Hz for just over 56 minutes, giving 23559 frames in
 144 total across [X] trials of the task.

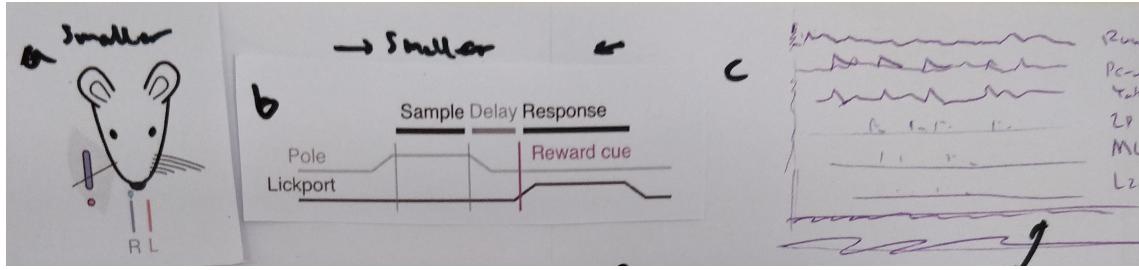


Figure 3: Experimental data from (Peron et al., 2015a).

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

145 Our primary goal is to understand how the choices of deconvolving these calcium-
 146 imaging data alter the scientific inferences we can draw. As our baseline, we use the
 147 “raw” $\Delta F/F$ time-series of changes in calcium indicator fluorescence. We use the above
 148 three discrete deconvolution methods to extract spike counts (MLSpike), event occurrence
 149 (LZero), or event magnitude (Suite2p) per frame. For comparison, we use (Peron et al.,
 150 2015a)’s own version of denoised calcium time-series, created by detecting calcium event
 151 onsets by a threshold and convolving events with a spike-response kernel. As an example
 152 of simpler methods, we use Yaksi and Friedrich (2006)’s simple deconvolution of the raw
 153 calcium with a fixed kernel of the calcium response to a single spike. And finally we create
 154 smoothed versions of the discrete-deconvolution methods, by convolving their recovered
 155 spikes/events with a spike-response kernel. Figure 3c show an example raw calcium time-
 156 series for one neuron, and the result of applying each of these 8 processing methods.
 157 We thus repeat all analyses on 9 different sets of time-series extracted from the same
 158 population recording.

159 We choose the algorithm parameters as follows. Simple deconvolution Yaksi and
 160 Friedrich (2006) is just kernel of the calcium response to a single spike; here is GCAMP6s[?],
 161 so parameterise kernel to that [Same as Perons kernel?]. For the three discrete deconvolu-
 162 tion methods, we choose the modal values of the best-fit parameters that optimised the
 163 Error Rate over the ground-truth recordings. This seems a reasonably consistent choice,
 164 of using the most consistently performing values obtained from comparable data: neurons
 165 in the same layer (L2/3) in the same species (mouse), in another primary sensory area
 166 (V1). Most importantly for our purposes, choosing the modal values means we avoid
 167 pathological regions of the parameter space. [FILL IN DETAILS AND TIDY]

168 2.4 Deconvolution methods disagree on estimates of simple neural statistics

170 We first check how well each approach recovers the basic statistics of neural activity event
 171 rates in L2/3 of barrel cortex. Electrophysiology has shown that the distribution of firing
 172 rates across neurons in a population is consistently long-tailed, and often log-normal, all
 173 across rodent cortex (Wohrer et al., 2013); and L2/3 neurons in barrel cortex are no
 174 different (O’Connor et al., 2010), with median firing rates less than 1 Hz, and a long right-
 175 hand tail of rarer high-firing neurons. We thus expect the calcium event rates or spike
 176 rates from our time-series would follow such a distribution. (Event rates for raw calcium,

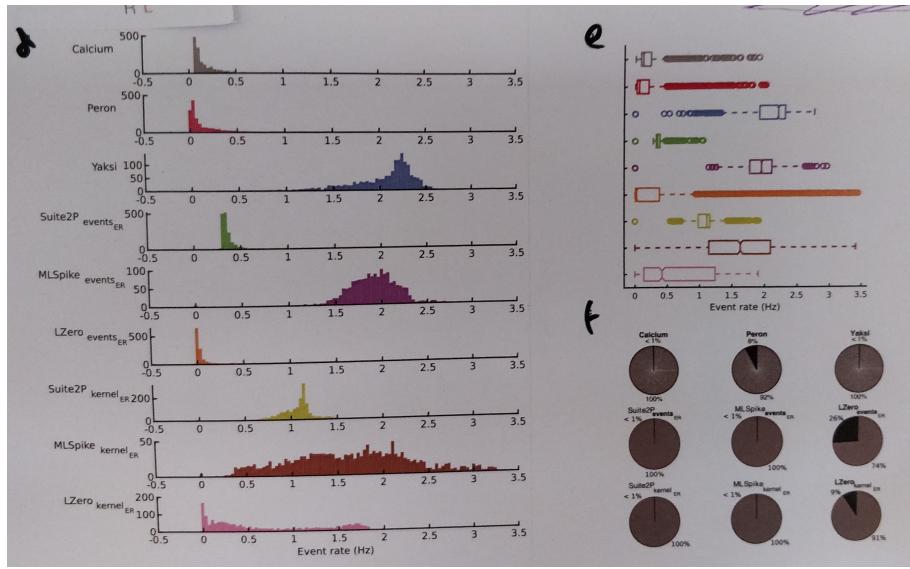


Figure 4: Estimates of population-wide event rates vary qualitatively across deconvolution methods.

(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) cells for each method. Silent cells are defined following (Peron et al., 2015b) as those with an event rate less than 0.0083Hz.

177 Peron, Yaksi and the convolved obtained by thresholding the calcium time-series)
178 Figure 4a shows that the raw calcium and two of the discrete deconvolution methods
179 (Suite2p, LZero) have qualitatively correct distributions of event rates (median near zero,
180 long right-hand tails). The Peron time-series also have the correct distribution of event
181 rates, which is unsurprising as it was tuned to do so. All other methods give qualitatively
182 wrong distributions of spike rates (MLSpike) or event rates (all other methods). There is
183 also little overlap in the distributions of spike rates between the three discrete deconvolu-
184 tion methods. Applying a kernel to their inferred spikes/events shifts rather than smooths
185 the firing rate distributions (Suite2P_{kernel} , MLSpike_{kernel} , LZero_{kernel}), suggesting noise
186 in the deconvolution process is amplified through the additional steps of convolving with
187 a kernel and thresholding.

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

197 2.5 Inferences of single cell tuning differ widely between raw calcium 198 and deconvolved methods

199 We turn now to what we can infer about simple properties of neural coding, and how our
200 choice of deconvolution method can alter those inferences. The decision task facing the
201 mouse (Fig. 3a) requires that it moves its whisker back-and-forth to detect the position
202 of the pole, delay for a second after the pole is withdrawn, and then make a choice of
203 the left or right lick-port based on the pole's position (Fig. 3b). As the imaged barrel
204 corresponds to the single spared whisker (on the contralateral side of the face), so the
205 captured population activity during each trial likely contains neurons tuned to different
206 aspects of the task. We show here that the number and identity of such task-tuned neurons
207 in the population differ widely between deconvolution methods.

208 Following Peron et al. 2015a, we define a task-tuned cell as one for which the peak
209 in its trial-averaged histogram of activity exceeds the predicted upper limit from shuffled
210 data (Fig. 5a0; see Methods). When applied to the raw calcium time-series, close to half
211 the neurons are tuned (Fig. 5a). This is more than double the proportion found for the
212 next nearest method (Yaksi's simple deconvolution), and at least a factor of 5 greater than
213 the proportion of tuned neurons resulting from any discrete deconvolution method, which
214 each report less than 10% of the neurons are tuned.

215 Worse, few neurons are detected as tuned in time-series resulting from multiple meth-
216 ods (Fig. 5b). Only XXX neurons are labelled as tuned in at least two sets of time-series,
217 and just 21 (X %) are labelled as tuned in all nine. Even separately considering the con-
218 tinuous and discrete time-series, we find only 38 cells are tuned across all six continuous
219 methods, and 25 neurons for all three discrete deconvolution methods (Fig. 5c). Figure
220 5d illustrates the diversity of detected tuning even amongst the neurons with the greatest
221 agreement between methods.

222 These results suggest that raw calcium alone over-estimates tuning in the population,
223 but also that there can be substantial disagreement between deconvolution methods. One
224 solution for robust detection of tuned neurons is to find those agreed between the raw

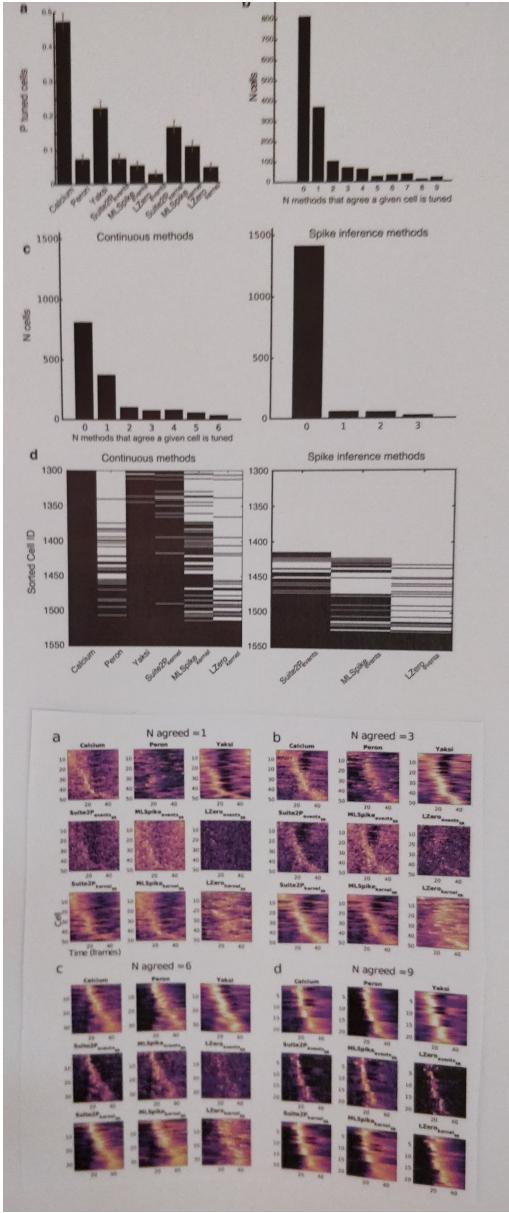


Figure 5: Inferences of single cell tuning show poor agreement between raw calcium and deconvolution methods, and between methods.

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

(a) Number of tuned cells 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 cells classified as tuned by exactly N methods.

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

(d) Comparison of cell tuning across methods. Each row shows whether that cell 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 cell as tuned.

(e-h) Identifying robust cell tuning. Panel groups (e) to (h) show cells classed as tuned by increasing numbers of deconvolution methods. Each panel within a group plots one cell'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 cells, but following processing by each of the eight deconvolution methods. [Q: How were the neurons chosen?]

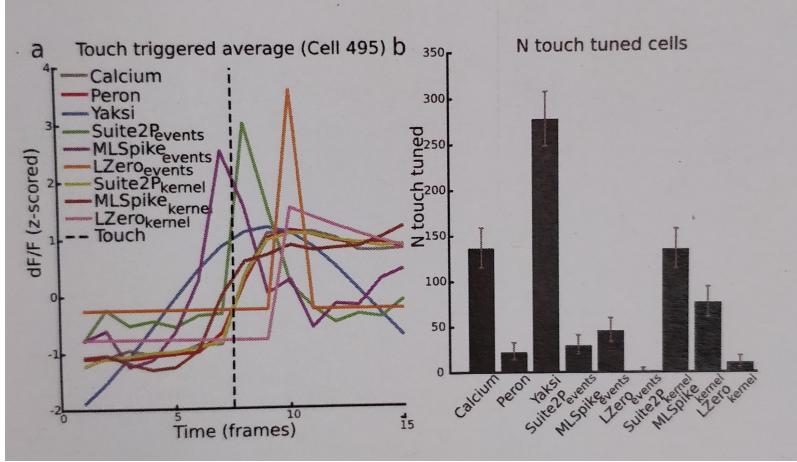


Figure 6: Touch-triggered neuron responses.

(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 cells across deconvolution methods. A cell is classed as touch-tuned if its peak touch-triggered activity is significantly greater than shuffled data. Error bars are Jeffreys intervals

calcium time-series and more than one deconvolution method. In Figure 5e-h, we show how increasing the number of methods required to agree on a neuron’s tuned status creates clear agreement between time-series processed with all methods, even if a particular method did not reach significance for that cell. Even requiring agreement between the raw calcium and just two other methods is enough to see tuning of many cells. The identification of unambiguously task-tuned cells could thus be achieved by triangulating the raw calcium with the output of multiple deconvolution methods.

In the pole detection task considered here, neurons tuned to pole contact are potentially crucial to understanding the sensory information used to make a decision. Touch onset is known to drive a subset of neurons to spike with short latency and low jitter (O’Connor et al., 2010; Hires et al., 2015). Detecting such rapid, precise responses in the slow kinetics of calcium imaging is challenging, suggesting discrete-deconvolution methods might be necessary to detect touch-tuned neurons. To test this, in each of the 9 sets of time-series we identify touch-tuned neurons by a significant peak in their touch-triggered activity (Fig ??a). Figure 6b shows that, while all data-sets have touch-tuned neurons, the number of such neurons differs substantially between them. And rather than being essential, discrete deconvolution methods disagree strongly on touch-tuning, with MLSpike (events) finding 45 touch-tuned neurons and LZero (events) finding one. Thus our inferences of the coding of task-wide or specific sensory events crucially depends on our choice of calcium imaging time-series.

2.6 Inconsistent recovery of population correlation structure across deconvolution approaches

The high yield of neurons from calcium imaging is ideal for studying the dynamics and coding of neural populations [cite some examples here]. Many analyses of populations start from pairwise correlations between cells, whether as measures of a population’s synchrony or joint activity, or as a basis for further analyses like clustering and dimension reduction [Cunningham]. We now show how our inferences of population correlation structure also

252 depend strongly on the choice of deconvolution method.

253 Figure 7a shows that the distributions of pairwise correlations qualitatively differ be-
254 tween the sets of time-series we derived from the same calcium imaging data. The con-
255 siderably narrower distributions from the discrete deconvolution time-series compared to
256 the others is expected, as these time-series are sparse. Nonetheless, there are qualitative
257 differences within the sets of discrete and continuous time-series. Some distributions are
258 approximately symmetric, with broad tails; some asymmetric with narrow tails; the corre-
259 lation distribution from the Peron method time-series is the only one with a median below
260 zero. These qualitative differences are not due to noisy estimates of the pairwise corre-
261 lations: for all our sets of time-series the correlations computed on a sub-set of time-points
262 in the session agree well with the correlations computed on the whole session (Figure 7b).
263 Thus pairwise correlation estimates for each method are stable, but their distributions
264 differ between methods.

265 Looking in detail at the full correlation matrix shows that even for methods with simi-
266 lar distributions, their agreement on correlation structure is poor. Some neuron pairs that
267 appear correlated from time-series processed by one deconvolution method are uncorre-
268 lated when processed with another method (Figure 7c). Over the whole population, the
269 correlation structure obtained from the raw calcium, Yaksi and Suite2p (kernel) time-series
270 all closely agree, but nothing else does (Figure 7d): the correlation structure obtained from
271 LZero agrees with nothing else; and the discrete deconvolution methods all generate dis-
272 similar correlation structures (Figure 7e). Our inferences about the extent and identity
273 of correlations within the population will differ qualitatively depending on our choice of
274 imaging time-series.

275 **2.7 Deconvolution methods show the same population activity is both 276 low and high dimensional**

277 Dimensionality reduction techniques, like principal components analysis (PCA), allow re-
278 searchers to make sense of large scale neuroscience data ([Cunningham and Yu, 2014](#)) [more
279 here], by reducing the data from N neurons to $d < N$ dimensions. Key to such analyses
280 is the choice of d , a choice guided by how much of the original data we can capture. To
281 assess such inferences of population dimensionality, we apply PCA to our 9 sets of imag-
282 ing time-series to estimate the dimensionality of the imaging data (which for PCA is the
283 variance explained by each eigenvector of the data's covariance matrix).

284 Figure 8a plots for each deconvolution method the cumulative variance explained when
285 increasing the number of retained dimensions. Most deconvolution methods qualitatively
286 disagree with the raw calcium data-set on the relationship between dimensions and vari-
287 ance. This relationship is also inconsistent across deconvolution methods; indeed the
288 discrete deconvolution methods result in the shallowest (MLSpike_{events}) and amongst the
289 steepest (LZero_{events}) relationships between increasing dimensions and variance explained.
290 The number of dimensions required to explain 80% of the variance in the data ranges from
291 $d = 125$ (Peron) to $d = 1081$ (MLSpike_{events}), a jump from 8% to 70% of all possible di-
292 mensions (Fig 8b). Thus we could equally infer that the same L2/3 population activity
293 is low dimensional (<10% dimensions required to explain 80% of the variance) or high-
294 dimensional (>50% of dimensions required) depending on our choice of imaging time-series.

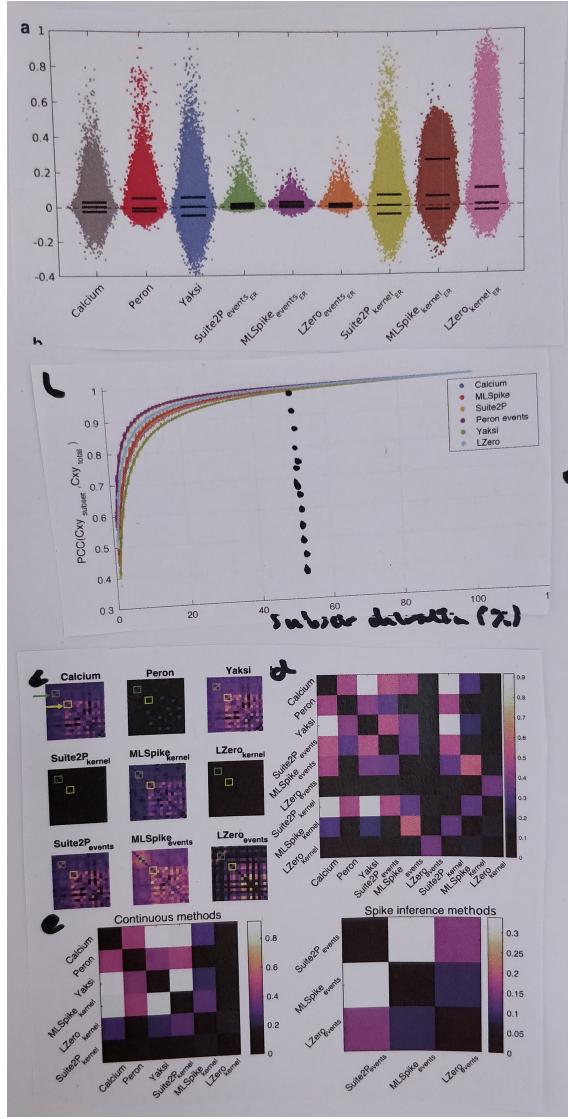


Figure 7: Effects of deconvolution on pairwise correlations between neurons.

- (a) Distributions of pairwise correlations between all cells, for each deconvolution method (one dot per cell 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 ($C_{xy\text{total}}$) to the same matrix constructed on a subset of the data ($C_{xy\text{subset}}$). Each data-point is the correlation between $C_{xy\text{total}}$ and $C_{xy\text{subset}}$; one line per deconvolution method.
- (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 cells: 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 correlation between the full-data correlation matrix for that pair of methods. [REDO: compare correlation matrices using Spearmans to test rank ordering of pairs]
- (e) as in (d), but split to show continuous methods (left) or discrete deconvolution methods (right). [Check all labels: clearly something inconsistent between (d) and (e): see Suite2pkernel results in both]

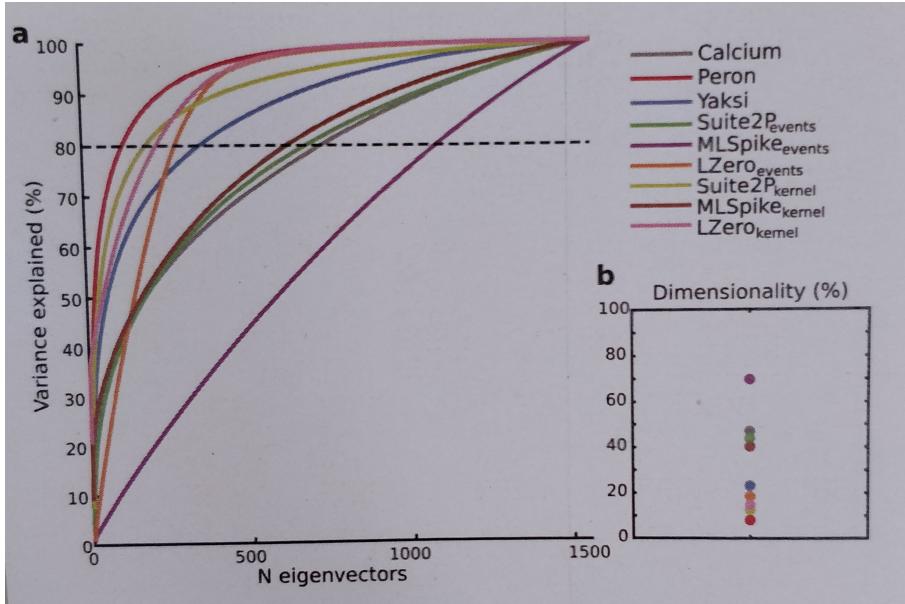


Figure 8: Dimensionality of population activity.

(a) Cumulative variance explained by each dimension of the data's 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.

295 3 Discussion

296 Imaging of somatic calcium is a remarkable tool for capturing the simultaneous activity of
 297 hundreds to thousands of neurons. But the time-series of each neuron's calcium fluorescence
 298 is inherently noisy and non-linearly related to its spiking. We sought here to address
 299 how our choice of corrections to these time-series – to use them raw, deconvolve them into
 300 continuous time-series, or deconvolve them into discrete events – affect the quality and
 301 reliability of the scientific inferences drawn.

302 Our results show the choice qualitatively changes the potential scientific inferences we
 303 draw about neural activity, coding, and correlation structure. We consistently observe
 304 that the analysis results differ sharply between the raw calcium and most, if not all, of the
 305 processed time-series. However, the deconvolved time-series also consistently disagreed
 306 with each other, even between methods of the same broad class (continuous or discrete
 307 time-series).

308 3.1 Accurate discrete deconvolution is possible, but sensitive

309 We find much that is encouraging. In fitting discrete deconvolution methods to ground-
 310 truth data, we found they can in principle accurately recover neural activity. A caveat
 311 here is that the choice of metric for evaluation and fitting of parameters is of critical
 312 importance. The widely-used Pearson correlation coefficient is a poor choice of metric
 313 as it returns inconsistent results with small changes in algorithm parameters, and leads
 314 to poor estimates of simple measures such as firing rate when used across methods and
 315 sampling rates. By contrast, the Error Rate metric (Deneux et al., 2016; Victor and
 316 Purpura, 1996) resulted in excellent recovery of ground-truth spike trains. Other recently
 317 developed methods for comparing spike-trains based on information theory (Theis et al.,

318 2016) or fuzzy set theory (Reynolds et al., 2017), may also be appropriate.

319 However, while good estimates of neural activity can be achieved with modern dis-
320 crete deconvolution methods [cf Spikefinder, Pitcharriu 2018 JNS], the best parameters
321 vary substantially between cells, and small changes in analysis parameters result in poor
322 performance. This variation and sensitivity of parameters played out as widely-differing
323 results between the three discrete deconvolution methods in analyses of neural activity,
324 coding, and correlation structure.

325 3.2 Choosing parameters for deconvolution methods

326 A potential limitation of our study is that we use a single set of parameter values for
327 each discrete deconvolution method applied to the population imaging data from barrel
328 cortex. But then our situation is the same as that facing any experimentalist: in the
329 absence of ground-truth, how do we set the parameters? Our solution here was to use
330 the modal parameter values from ground-truth fitting, as these values are candidates for
331 the most general solutions. We also felt these were a reasonable choice for the population
332 imaging data from barrel cortex, given that the ground-truth recordings came from the
333 same species (mouse) in the same layer (2/3) of a different bit of primary sensory cortex
334 (V1).

335 Rather than use the most general parameters values, another solution would be to
336 tune the parameters to obtain known gross statistics of the neural activity. This was the
337 approach used by Peron and colleagues [cite] to obtain the denoised Peron time-series we
338 included here. But as we've seen, this approach can lead to its own problems: for example,
339 in the Peron time-series, it created a distributions of correlations that differed from any
340 other set of time-series. Indeed, finding good parameter values may be an intractable
341 problem, as it is possible each neuron requires individual fitting, to reflect the combination
342 of its expression of fluorescent protein, and its particular non-linearity between voltage and
343 calcium.

344 3.3 Ways forward

345 The simplest solution to the inconsistencies between different forms of time-series is to
346 triangulate them, and take the consensus across their results. For example, our finding of
347 a set of tuned neurons across multiple methods is strong evidence that neurons in L2/3
348 of barrel cortex are responsive across the stages of the decision task. Further examples
349 of such triangulation in the literature are rare; Klaus and colleagues (Klaus et al., 2017)
350 used two different pipelines to derive raw $\Delta F/F$ of individual neurons from one-photon
351 fibre-optic recordings in the striatum, and replicated all analyses using the output of both
352 pipelines. Our results encourage the further use of triangulation to create robust inference:
353 obtaining the same result in the face of wide variation increases our belief in its reliability
354 (Munaf and Davey Smith, 2018).

355 There are caveats to using triangulation. For single neuron analyses, triangulation
356 inevitably comes at the price of reducing the yield of neurons to which we can confidently
357 assign roles. A further problem for triangulation is how to combine more complex analyses,
358 such as pairwise correlations; the alternative is to rely on qualitative comparisons.

359 Many studies use the raw calcium signal as the basis for all their analyses [cite], perhaps
360 assuming this is the least biased approach. Our result show this is not so: the discrepancy
361 between raw and deconvolved calcium on single neuron coding suggests an extraordinary
362 range of possible results, from about half of all neurons tuned to the task down to less
363 5 percent. The qualitative conclusion – there is coding – is not satisfactory. Thus our

364 results should not be interpreted as a call to abandon deconvolution methods; rather they
365 serve to delimit how we can interpret their outputs.

366 Instead, we need deconvolution solved: as sensors with faster kinetics (though funda-
367 mentally limited by kinetics of calcium release itself) and higher signal-to-noise ratios are
368 developed [cite GCAMP7; red-shifted], so the accuracy and robustness of de-noising and
369 deconvolution should improve; and as the neuron yield continues to increase [Ahrens], so
370 the potential for insights from inferred spikes or spike-driven events grows. Developing
371 further advanced deconvolution algorithms will harness these advances, but are potentially
372 always limited by the lack of ground-truth to fit their parameters. Our results may pro-
373 vide impetus for a different direction of research, focussing on how we can get consensus
374 among the output of different algorithms, and thus provide robust scientific inferences
375 about neural populations.

376 **4 Methods**

377 [FINISH all methods: add all details from deconvolution debugging doc; fix all issues
378 flagged below; turn notes into prose.]

379 **Ground truth data**

380 Ground truth data was accessed from crcns.org ([Svoboda, 2015](#)), and the experiments have
381 been described previously [Chen et al. \(2013\)](#). Briefly, mouse visual cortical neurons ex-
382 pressing the fluorescent calcium reporter protein GCamP6s were imaged with two-photon
383 microscopy at 60Hz. Loose-seal cell-attached recordings were performed simultaneously
384 at 10kHz. The data-set contains twenty one recordings from nine cells.

385 **Spike train metrics**

386 Pearson correlation coefficient was computed between the ground truth and inferred spike
387 trains following convolution of both with a gaussian kernel (61 samples wide, 1.02 seconds).
388 [WHAT were used as inferred spike trains]

389 Error Rate was computed between the ground truth and inferred spike trains using the
390 [Deneux et al. \(2016\)](#) implementation of normalised error rate, derived from [Victor and](#)
391 [Purpura \(1996\)](#) Error Rate (code available <https://github.com/MLspike>). Briefly, the
392 error rate is 1 - F1-score, where the F1-score is the harmonic mean of sensitivity (number
393 of missed spikes divided by total spikes) and precision (number of falsely detected spikes
394 divided by total detected spikes),

395 [write this in full, defining sensitivity and precision]

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

396 Hits, misses and false detections were counted with a temporal precision of 0.5 seconds.
397 [Firing rate: how computed for true and inferred?]

398 For normalised estimation of errors in firing rate, we computed (estimated FR - true
FR / true FR)

399 **Parameter fitting**

400 For each method the best parameters for each cell were determined by brute force search
401 over an appropriate range. [Give the ranges. How was “appropriate” defined?]

402 The modal best parameters, as determined using Error Rate on downsampled data,
403 were then fixed for the population imaging data analysis. These were: [Complete with
404 parameter values]

405 **Downsampling**

406 Ground truth calcium data was downsampled from 60Hz to 7Hz in Matlab by up-sampling
407 by 7 - `interp(ca, 7)` and downsampling the resultant signal by 60, as Matlab’s downsam-
408 pling must be done in integer steps [??? how was the downsampling done then? By taking
409 every Nth frame from the interpolated data? Or averaging over every approx 8.5 frames
410 of the 60Hz signal?].

411 **Population imaging data description**

412 Population imaging data was accessed from crcns.org and have been described previously
413 (Peron et al., 2015b). Briefly, volumetric two photon calcium imaging of primary
414 somatosensory cortex (S1) was performed in awake head-fixed mice performing a whisker-
415 based object localisation task. In the task a metal pole was presented on one of two loca-
416 tions and mice were motivated with fluid reward to lick at one of two lick ports depending
417 on the location of the pole following a brief delay. Two photon imaging of GCaMP6s
418 expressing neurons in superficial S1 was performed at 7Hz. Images were motion corrected
419 and aligned, before regions of interest were manually set and neuropil-subtracted. A single
420 recording from this dataset was used for population analysis. The example session had
421 1552 neurons recorded for a total of 23559 frames (56 minutes).

422 **4.1 Event rate estimation**

423 Spike inference methods (Suite2P_{events} , MLSpike_{events} , LZero_{events}) return estimated spikes
424 per frame which were converted into mean event rates (Hz) per cell [FIX THESE descrip-
425 tions]. The event rate for continuous methods (Calcium, Peron, Yaksi, Suite2P_{kernel} , MLSpike_{kernel} ,
426 LZero_{kernel}) for each cell was determined by counting activity/fluorescence events greater
427 than three standard deviations of the background noise. Background noise was calculated
428 by taking a four-point moving average [missing object in sentence: average of what?] and
429 subtracting this [what is "this"?] from the activity/fluorescence trace. Event rate was
430 then computed in Hz.

431 [what does this all mean? That each data-point for the background noise was an
432 average over 4 adjacent frames, and the standard deviation of the noise was computed
433 from those data-points? Why smooth the noise? And what segment of the data was
434 treated as "background"? (i.e. how many data-points)? Or does this mean that the whole
435 Ca²⁺ trace for each cell was smoothed using a 4-frame average (shifting 1 frame?), and
436 the SD of that smoothed signal was used as an estimate of background noise?]

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

439 **4.2 Task-tuned cells**

440 Task-tuning was determined for each neuron using the model-free approach of Peron et al.
441 (2015b). Neurons were classed as task-tuned if their peak trial-average activity exceeded
442 the 95th percentile of a distribution of trial-average peaks from shuffled data (10000 shuf-
443 fles). The shuffle test was done separately for correct lick-left and lick-right trials and cells
444 satisfying the tuning criteria in either case were counted as task-tuned.

445 Tuned cell agreement was calculated as the number of methods that agreed to the
446 tuning status of a given cell, for all methods and separately for continuous and spike
447 inference methods.

448 **4.3 Touch-related responses**

449 Touch-tuned cells were determined by computing touch-triggered average activity for each
450 cell, before calculating whether the data distribution of peak touch-induced activity ex-
451 ceeds the expected activity of shuffled data. In more detail, the time of first touch -
452 between the mouse's whisker and the metal pole - on each trial was recorded. For each
453 touch time, one second of activity (seven data samples) was extracted before and after the
454 frame closest to touch (15 samples total); taking the mean of these gave the average touch

455 response for the cell. The time of peak touch-triggered average activity was calculated,
456 and a ranksum test (bonferroni corrected) between the true data distribution at peak time
457 and a matched random sample of data from the same cell.

458 [Figure legend said: (Mann-Whitney U test, Benjamini Hochberg corrected); which is
459 it? And what is α here?]

460 [unclear what the pairwise tests were between; the mean activity at the peak time, and
461 a N-length vector of mean activity as the same time obtained from N shuffled datasets?]

462 [Shuffled how many times? Shuffled how?].

463 4.4 Pairwise correlations

464 Pairwise correlations were calculated for all pairs of neurons in Matlab (corrcoef) at the
465 data sampling rate (7Hz). [State what forms of time-series were correlated - from debug-
466 ging doc]

467 Correlations between correlation matrices (Fig. ??) were computed between the unique
468 pairwise correlations from each method (i.e. CXY(find(triu(CXY)))). [Add Spearman's
469 rank here]

470 Stability of correlation. For each deconvolution method, we computed the pairwise
471 correlation matrix using the entire sessions data, as above. We then sampled a subset of
472 time-points (1%-100%) of the full dataset at random without replacement and compute
473 its matrix of pairwise correlations. We compute the similarity between the [Total] and
474 [Subset] matrices using Pearsons correlation coefficient. [Finish this]

475 4.5 Dimensionality

476 To determine the dimensionality of each dataset we performed eigendecomposition of the
477 covariance matrix of each dataset. [Computed as per the pairwise correlations above?] The
478 resultant eigenvalues were sorted into descending order, and the variance explained
479 (`cumsum(egs)/sum(egs)`) plotted.

480 List of deconvolution methods

481 Suite2P

482 Suite2P (<https://github.com/cortex-lab/Suite2P>) is actively developed by Marius Pachitariu
483 (HHMI Janelia) and members of the cortexlab (Kenneth Harris and Matteo Carandini)
484 at UCL. Suite2P's USP is it's application to large scale 2-photon imaging analysis,
485 with an emphasis on end-to-end processing (images to neural event time series) and speed.
486 A preprint describing the toolbox is available here ([Pachitariu et al.](#))

487

488 [Explain: the basic method, and its free parameters]

489 <http://biorxiv.org/content/early/2016/06/30/061507>,

490

491 and our own notes on the spike detection algorithm are here:

492

493 <https://drive.google.com/open?id=1NeQhmoRpS-x8R0e84w3TqkUR1PNMXiem6ZIjJta-U7A>.

494 ML Spike

495 ML Spike (<https://github.com/mlspike>) was developed by Thomas Deneux at INT, CRNS
496 Marseille, France. A model-based probabilistic approach, ML Spike was developed to re-

497 cover spike trains in calcium imaging data by taking baseline fluctuations and cellular
498 properties into account. A comprehensive explanation of the algorithm and its benefits
499 can be found in the paper ([Deneux et al., 2016](#)).
500 [Explain: the basic method, and its free parameters]

501 MLSpike can return a maximum a posteriori spike train, or a spike probability per
502 time step. (*TO DO: We show results for both denoted $MLSpike_{events}$ and $MLSpike_{pspike}$
503 in Supplement*)

504 **LZero**

505 The method we refer to as **LZero** was developed by Sean Jewell and Daniela Witten
506 from U.Washington, Seattle, USA. The goal for this implementation was to cast spike
507 detection as a change-point detection problem, which could be solved with an existing
508 l_0 optimization algorithm. In their paper Jewell and Witten show that the l_0 solution is
509 better than previously implemented l_2 solutions, with results much closer to the real spike
510 train (l_2 solutions tend to overestimate the true firing rate). Details can be found in the
511 paper ([Jewell and Witten, 2017](#)).

512 [Explain: the basic method, and its free parameters]

513 Link: <https://arxiv.org/abs/1703.08644>

514 **Yaksi**

515 **Yaksi** refers to the ‘vanilla’ deconvolution of Yaksi and Friedrich (2006). This is to be used
516 as a baseline for comparison with more sophisticated methods. The method is detailed in
517 the paper: ([Yaksi and Friedrich, 2006](#)).

518 [Explain: the basic method]

519 **Peron events**

520 **Peron events** refer to the extracted events detailed in the original [Peron et al. \(2015b\)](#)
521 paper. It is a version of the ‘peeling’ algorithm ([Lütcke et al., 2013](#)) tuned to generate a
522 low number of false positive detections (a rate of 0.01Hz) on ground truth data, leading
523 to a hit rate of 54%.

524 [Explain: the basic method, and its free parameters]

525 **Events + kernel versions**

526 Where a spike inference method returns spike rates per time point, these are plotted as
527 Method_{events}. To compare to other methods that return a de-noised dF/F or firing rate
528 estimates, these events are convolved with a calcium kernel and plotted as Method_{kernel}.

529 [Explain: what are the parameters of the kernel]

530 **References**

531 Philipp Berens, Jeremy Freeman, Thomas Deneux, Nicolay Chenkov, Thomas McColgan,
532 Artur Speiser, Jakob H Macke, Srinivas C Turaga, Patrick Mineault, Peter Rupprecht,
533 Stephan Gerhard, Rainer W Friedrich, Johannes Friedrich, Liam Paninski, Marius Pa-
534 chitariu, Kenneth D Harris, Ben Bolte, Timothy A Machado, Dario Ringach, Jasmine
535 Stone, Luke E Rogerson, Nicolas J Sofroniew, Jacob Reimer, Emmanouil Froudarakis,
536 Thomas Euler, Miroslav Román Rosón, Lucas Theis, Andreas S Tolias, and Matthias

- 537 Bethge. Community-based benchmarking improves spike rate inference from two-photon
538 calcium imaging data. *PLoS Comput. Biol.*, 14(5):e1006157, May 2018.
- 539 Emery N Brown, Robert E Kass, and Partha P Mitra. Multiple neural spike train data
540 analysis: state-of-the-art and future challenges. *Nat. Neurosci.*, 7(5):456–461, May 2004.
- 541 Tsai-Wen Chen, Trevor J Wardill, Yi Sun, Stefan R Pulver, Sabine L Renninger, Amy
542 Baohan, Eric R Schreiter, Rex A Kerr, Michael B Orger, Vivek Jayaraman, Loren L
543 Looger, Karel Svoboda, and Douglas S Kim. Ultrasensitive fluorescent proteins for
544 imaging neuronal activity. *Nature*, 499(7458):295–300, July 2013.
- 545 John P Cunningham and Byron M Yu. Dimensionality reduction for large-scale neural
546 recordings. *Nat. Neurosci.*, 17(11):1500–1509, November 2014.
- 547 Thomas Deneux, Attila Kaszas, Gergely Szalay, Gergely Katona, Tamás Lakner, Amiram
548 Grinvald, Balázs Rózsa, and Ivo Vanzetta. Accurate spike estimation from noisy calcium
549 signals for ultrafast three-dimensional imaging of large neuronal populations *in vivo*.
550 *Nat. Commun.*, 7:12190, July 2016.
- 551 Samuel Andrew Hires, Diego A Gutnisky, Jianing Yu, Daniel H O’Connor, and Karel
552 Svoboda. Low-noise encoding of active touch by layer 4 in the somatosensory cortex.
553 *Elife*, 4, August 2015.
- 554 Sean Jewell and Daniela Witten. Exact spike train inference via ℓ_0 optimization. March
555 2017.
- 556 Andreas Klaus, Gabriela J Martins, Vitor B Paixao, Pengcheng Zhou, Liam Paninski, and
557 Rui M Costa. The spatiotemporal organization of the striatum encodes action space.
558 *Neuron*, 95:1171–1180.e7, Aug 2017. ISSN 1097-4199. doi: 10.1016/j.neuron.2017.08.
559 015.
- 560 Henry Lütcke, Felipe Gerhard, Friedemann Zenke, Wulfram Gerstner, and Fritjof Helm-
561 chen. Inference of neuronal network spike dynamics and topology from calcium imaging
562 data. *Front. Neural Circuits*, 7:201, December 2013.
- 563 Marcus R Munaf and George Davey Smith. Robust research needs many lines of evidence.
564 *Nature*, 553:399–401, January 2018. ISSN 1476-4687. doi: 10.1038/d41586-018-01023-3.
- 565 Daniel H O’Connor, Simon P Peron, Daniel Huber, and Karel Svoboda. Neural activity
566 in barrel cortex underlying vibrissa-based object localization in mice. *Neuron*, 67(6):
567 1048–1061, September 2010.
- 568 Marius Pachitariu, Carsen Stringer, Sylvia Schröder, Mario Dipoppa, L Federico Rossi,
569 Matteo Carandini, and Kenneth D Harris. Suite2p: beyond 10, 000 neurons with stan-
570 dard two-photon microscopy.
- 571 António R C Paiva, Il Park, and José C Príncipe. A comparison of binless spike train
572 measures. *Neural Comput. Appl.*, 19(3):405–419, April 2010.
- 573 Simon Peron, Tsai-Wen Chen, and Karel Svoboda. Comprehensive imaging of cortical
574 networks. *Curr. Opin. Neurobiol.*, 32:115–123, June 2015a.
- 575 Simon P Peron, Jeremy Freeman, Vijay Iyer, Caiying Guo, and Karel Svoboda. A cellular
576 resolution map of barrel cortex activity during tactile behavior. *Neuron*, 86(3):783–799,
577 May 2015b.

- 578 Stephanie Reynolds, Simon R Schultz, and Pier Luigi Dragotti. CosMIC: A consistent
579 metric for spike inference from calcium imaging. December 2017.
- 580 K Svoboda. Simultaneous imaging and loose-seal cell-attached electrical recordings from
581 neurons expressing a variety of genetically encoded calcium indicators. *GENIE project,*
582 *Janelia Farm Campus, HHMI; CRCNS.org*, 2015.
- 583 Lucas Theis, Philipp Berens, Emmanouil Froudarakis, Jacob Reimer, Miroslav
584 Román Rosón, Tom Baden, Thomas Euler, Andreas S Tolias, and Matthias Bethge.
585 Benchmarking spike rate inference in population calcium imaging. *Neuron*, 90(3):471–
586 482, May 2016.
- 587 J D Victor and K P Purpura. Nature and precision of temporal coding in visual cortex:
588 a metric-space analysis. *J. Neurophysiol.*, 76(2):1310–1326, August 1996.
- 589 Adrien Wohrer, Mark D Humphries, and Christian K Machens. Population-wide distri-
590 butions of neural activity during perceptual decision-making. *Prog. Neurobiol.*, 103:
591 156–193, April 2013.
- 592 Emre Yaksi and Rainer W Friedrich. Reconstruction of firing rate changes across neuronal
593 populations by temporally deconvolved ca2+ imaging. *Nat. Methods*, 3(5):377–383, May
594 2006.