

# Online nonnegative deconvolution for spike train inference from population calcium imaging

Joshua T. Vogelstein, Adam M. Packer, Tim A. Machado,  
Tanya Sippy, Baktash Babadi, Rafael Yuste, Liam Paninski

November 24, 2009

## Abstract

Calcium imaging technologies for observing spiking activity simultaneously from large populations of neurons are quickly gaining popularity. While the raw data are fluorescence movies often the underlying spike trains are of interest. This work presents an online non-negative deconvolution filter to infer the approximately most likely spike trains for each neuron, given the fluorescence observations. This algorithm outperforms optimal linear deconvolution (aka, Wiener filter) on both simulated and in vitro data. The performance gains come from restricting the inferred spike trains to be positive (using an interior-point method), unlike the Wiener filter. The algorithm is fast enough that even when imaging  $\approx 100$  neurons, inference can be performed on the set of all observed traces super-real-time (ie, faster than real-time). Performing optimal spatial filtering on the images further refines the estimates. Importantly, all the parameters required to perform the inference can be estimated using only the fluorescence data, obviating the need to perform joint electrophysiological and imaging calibration experiments.

## 1 Introduction

Simultaneously imaging large populations of neurons using calcium sensors is becoming increasingly popular [1], both in vitro [2, 3] and in vivo [4, 5, 6], especially as the signal-to-noise-ratio (SNR) of genetic sensors continues to improve [7, 8, 9]. Whereas the data from these experiments are movies of time-varying fluorescence signals, the desired signal is typically the spike trains of the observable neurons. Unfortunately, finding the most likely spike train is a challenging computational task, due to poor SNR, poor temporal resolution, unknown parameters, and computational intractability. One is therefore effectively forced to find an approximately most likely spike train, or hope that the inferred spike train is indeed most likely.

A number of groups have therefore proposed algorithms to infer spike trains from calcium fluorescence data. For instance, Greenberg et al. [10] developed a novel template matching algorithm. Holekamp et al. [11] took a very different strategy, by performing the optimal linear deconvolution (i.e., the Wiener filter) on the fluorescence data. This approach is natural from a signal processing standpoint, but does not utilize the knowledge that spikes are always positive. Previously, a sequential Monte Carlo method to efficiently estimate the probability of a spike in each image frame, given the entire fluorescence time-series, was proposed [12]. While effective, that approach is not suitable for online analyses of populations of neurons, as the computations run in about real-time per neuron (i.e., analyzing one minute of data requires about one minute of computational time on a standard laptop computer), and real-time for a whole population of neurons would be desirable.

The present work starts by building a generative model relating spiking activity to fluorescence traces. Unfortunately, inferring the most likely spike train given this model is computationally intractable. Making some reasonable approximations leads to an algorithm that infers the approximately most likely spike train, given the fluorescence data. This algorithm has a few particularly noteworthy features, relative to other approaches. First, spikes are assumed to be positive. This assumption often improves filtering results when the underlying signal has this property [13, 14, 15, 16, 17, 18, 19, 20]. Second, the algorithm is extremely fast: it can process a calcium trace from 50,000 images in about one second on a standard laptop computer. In fact, filtering the signals for an entire population of about 100 neurons runs in super-real-time (ie, *faster* than real-time). This speed facilitates using this filter online, as observations are being collected. In addition to these two features, the model may be generalized in a number of ways, including incorporating spatial filtering of the raw movie. The efficacy of the proposed filter is demonstrated

on several real data-sets, suggesting this algorithm is a powerful and robust tool for online spike train inference. The code (which is a simple Matlab script) is available from the authors upon request.

## 2 Methods

### 2.1 Data driven generative model

Figure 1 shows a typical in vitro, epifluorescence data-set (see section 2.7 for data collection details). The top panel shows a field-of-view, including 3 neurons, two of which are patched. To build the model, a region-of-interest (ROI) is found, which in this case is the circled neuron. Given the ROI, all the pixel intensities of each frame can be averaged, to get a one-dimensional fluorescence time-series, as shown in the bottom left panel (black line). By patching onto this neuron, the spike train can also be directly observed (black bars). Previous work suggests that this fluorescence signal might be well characterized by convolving the spike train with an exponential, and adding noise [1]. This model is confirmed by convolving the true spike train with an exponential (gray line, bottom left panel), and then looking at the distribution of the residuals. The bottom right panel shows a histogram of the residuals (black line), and the best fit Gaussian distribution (gray line).

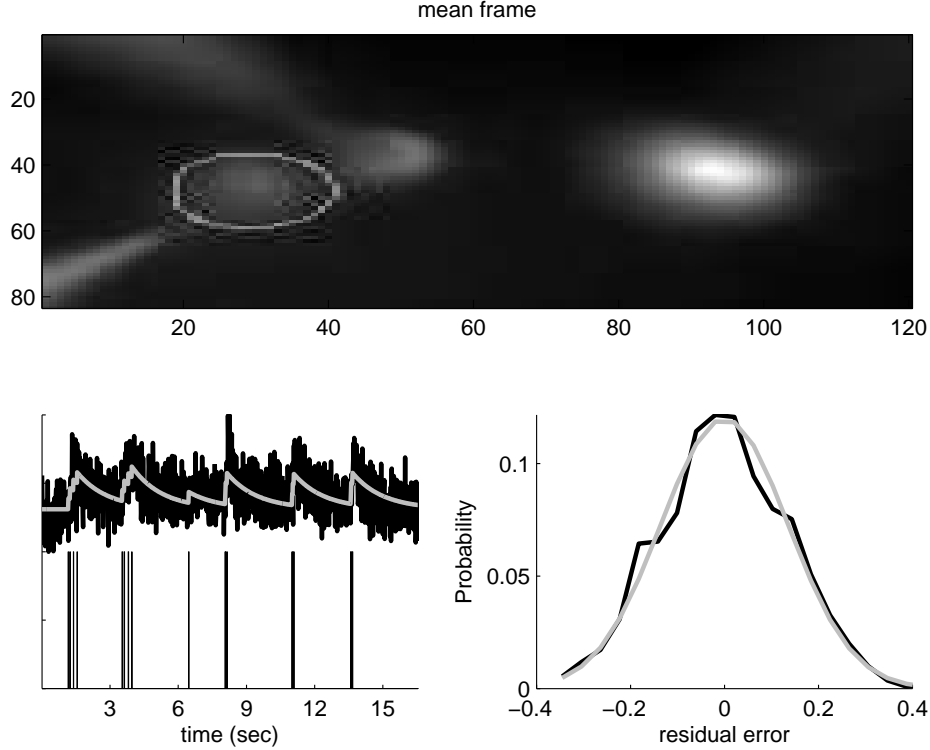


Figure 1: Typical in vitro data suggest that a reasonable first order model may be constructed by convolving the spike train with an exponential, and adding Gaussian noise. Top panel: the average (over frames) of a typical field-of-view. Bottom left: spike train (black bars), convolved with an exponential (gray line), superimposed on the one-dimensional fluorescence time-series (black line). Bottom right: a histogram of the residual error between the gray and black lines from the bottom left panel (black line), and the best fit Gaussian (gray line).

The above observations may be formalized as follows. Assume there is a one-dimensional fluorescence trace,  $F$  (throughout this text  $X$  indicates the vector  $(X_1, \dots, X_T)$ , where  $T$  is the index of the final frame), from a neuron. At time  $t$ , the fluorescence measurement,  $F_t$  is a linear-Gaussian function of the intracellular calcium concentration at

that time,  $[\text{Ca}^{2+}]_t$ :

$$F_t = \alpha([\text{Ca}^{2+}]_t + \beta) + \sigma\varepsilon_t, \quad \varepsilon_t \stackrel{iid}{\sim} \mathcal{N}(0, 1). \quad (1)$$

The scale,  $\alpha$ , absorbs all experimental variables impacting the scale of the signal, including number of sensors within the cell, photons per calcium ion, amplification of imaging system, etc. Similarly, the offset,  $\beta$ , absorbs baseline calcium concentration of the cell, background fluorescence of the fluorophore, imaging system offset, etc. The standard deviation,  $\sigma$ , results from calcium fluctuations independent of spiking activity, fluorescence fluctuations independent of calcium, and imaging noise. The noise at each time,  $\varepsilon_t$ , is independently and identically distributed according to a standard normal distribution (i.e., Gaussian with zero mean and unit variance), as indicated by the notation  $\stackrel{iid}{\sim} \mathcal{N}(0, 1)$ .

Then, assuming that the intracellular calcium concentration,  $[\text{Ca}^{2+}]_t$ , jumps by  $A \mu\text{M}$  after each spike, and subsequently decays back down to  $C_b \mu\text{M}$  with time constant,  $\tau$ , yields:

$$[[\text{Ca}^{2+}]_{t+\Delta} = (1 - \Delta/\tau)[\text{Ca}^{2+}]_t + C_b + An_t \quad (2)$$

where  $\Delta$  is the time step size — which is the frame duration, or  $1/(\text{frame rate})$  — and  $n_t$  indicates the number of times the neuron spiked in frame  $t$ . Note that because  $[\text{Ca}^{2+}]_t$  and  $F_t$  are linearly related to one another, the fluorescence scale,  $\alpha$ , and calcium scale,  $A$ , are not identifiable. In other words, either can be set to unity without loss of generality, as the other can absorb the scale entirely. Similarly, the fluorescence offset,  $\beta$ , and calcium baseline,  $C_b$  are not identifiable, so either can be set to zero without loss of generality. Finally, letting  $\gamma = (1 - \Delta/\tau)$ , Eq. (2) can be rewritten replacing  $[\text{Ca}^{2+}]_t$  with its non-dimensionalized counterpart,  $C_t$ :

$$C_{t+\Delta} = \gamma C_t + n_t. \quad (3)$$

Note that  $C_t$  does not refer to absolute intracellular concentration of calcium, but rather, a relative measure (but see [12] for a more general model). The gray line in the bottom left panel of Figure 1 corresponds to the putative  $C$  of the observed neuron.

To complete the “generative model” (i.e., a model from which simulations can be generated), the distribution from which spikes are sampled must be defined. Perhaps the simplest first order description of spike trains is that at each time, spikes are sampled according to a Poisson distribution with some rate:

$$n_t \stackrel{iid}{\sim} \text{Poisson}(\lambda\Delta) \quad (4)$$

where  $\lambda\Delta$  is the expected firing rate, and  $\Delta$  is included to ensure that the expected firing rate is independent of the frame rate. Thus, Eqs. (1), (3), and (4) complete the generative model.

## 2.2 Goal

Given the above model, the goal is to find the maximum *a posteriori* (MAP) spike train, i.e., the most likely spike train,  $\hat{n}$ , given the fluorescence measurements,  $\mathbf{F}$ :

$$\hat{n} = \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmax}} P[\mathbf{n}|\mathbf{F}], \quad (5)$$

where  $P[\mathbf{n}|\mathbf{F}]$  is the posterior probability of a spike train,  $\mathbf{n}$ , given the fluorescent trace,  $\mathbf{F}$ , and  $n_t$  is constrained to be an integer ( $\mathbb{N}_0 = \{0, 1, 2, \dots\}$ ), because of the above assumed Poisson distribution. From Bayes’ Rule, the posterior can be rewritten:

$$P[\mathbf{n}|\mathbf{F}] = \frac{P[\mathbf{n}, \mathbf{F}]}{P[\mathbf{F}]} = \frac{1}{P[\mathbf{F}]} P[\mathbf{F}|\mathbf{n}] P[\mathbf{n}], \quad (6)$$

where  $P[\mathbf{F}]$  is the evidence of the data,  $P[\mathbf{F}|\mathbf{n}]$  is the likelihood of observing a particular fluorescence trace  $\mathbf{F}$ , given the spike train  $\mathbf{n}$ , and  $P[\mathbf{n}]$  is the prior probability of a spike train. Plugging Eq. (6) into Eq. (5), yields:

$$\hat{n} = \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmax}} \frac{1}{P[\mathbf{F}]} P[\mathbf{F}|\mathbf{n}] P[\mathbf{n}] = \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmax}} P[\mathbf{F}|\mathbf{n}] P[\mathbf{n}], \quad (7)$$

where the second equality follows because  $P[\mathbf{F}]$  merely scales the results, but does not change the relative quality of various spike trains. Fortunately, both  $P[\mathbf{F}|\mathbf{n}]$  and  $P[\mathbf{n}]$  are available from the above model:

$$P[\mathbf{F}|\mathbf{n}] = P[\mathbf{F}|\mathbf{C}] = \prod_{t=1}^T P[F_t|C_t], \quad (8a)$$

$$P[\mathbf{n}] = \prod_{t=1}^T P[n_t], \quad (8b)$$

where the first equality in Eq. (8a) follows because  $\mathbf{C}$  is deterministic given  $\mathbf{n}$ , and the second equality follows from Eq. (1). Further, Eq. (8b) follows from the Poisson process assumption, Eq. (4). Both  $P[F_t|C_t]$  and  $P[n_t]$  can be written explicitly:

$$P[F_t|C_t] = \mathcal{N}(\alpha(C_t + \beta), \sigma^2), \quad (9a)$$

$$P[n_t] = \text{Poisson}(\lambda\Delta), \quad (9b)$$

where both equations follow from the above model. Now, plugging Eq. (9) back into (8), and plugging that result into Eq. (7), yields:

$$\hat{\mathbf{n}} = \operatorname{argmax}_{n_t \in \mathbb{N}_0 \forall t} \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(F_t - \alpha(C_t + \beta))^2}{\sigma^2} \right\} \frac{\exp\{-\lambda\Delta\}(\lambda\Delta)^{n_t}}{n_t!} \quad (10a)$$

$$= \operatorname{argmax}_{n_t \in \mathbb{N}_0 \forall t} \sum_{t=1}^T \left( -\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + n_t \log \lambda\Delta - \log n_t! \right), \quad (10b)$$

where the second equality follows from taking the logarithm of the right-hand-side. Unfortunately, solving Eq. (10b) exactly is computationally intractable, as it requires a nonlinear search over an infinite number of possible spike trains. The search space could be restricted by imposing an upper bound,  $k$ , on the number of spikes within a frame. However, in that case, the computational complexity scales *exponentially* with the number of image frames — i.e., the number of computations required would scale with  $k^T$  — which for pragmatic reasons is intractable.

### 2.3 Inferring the most likely spike train, given a fluorescence trace

The goal here is to develop an algorithm to efficiently approximate  $\hat{\mathbf{n}}$ , the most likely spike train, given the fluorescence trace. Because of the computational intractability described above, Eq. (10) is approximated by modifying Eq. (4), replacing the Poisson distribution with an exponential distribution. Modifying Eq. (10) to reflect this approximation yields:

$$\hat{\mathbf{n}} \approx \operatorname{argmax}_{n_t > 0 \forall t} \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(F_t - \alpha(C_t + \beta))^2}{\sigma^2} \right\} (\lambda\Delta) \exp\{-\lambda\Delta n_t\} \quad (11a)$$

$$\hat{\mathbf{n}} \approx \operatorname{argmax}_{n_t > 0 \forall t} \sum_{t=1}^T -\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 - n_t \lambda\Delta \quad (11b)$$

where the constraint on  $n_t$  has been relaxed from  $n_t \in \mathbb{N}_0$  to  $n_t \geq 0$  (since the exponential distribution can yield any non-negative number). The advantage of this approximation is that the optimization problem becomes log-concave, meaning that any gradient ascent method guarantees achieving the global maximum (because there are no local maxima, other than the single global maximum). The disadvantage, however, is that the integer constraint is lost, i.e., the answer could include “partial” spikes. This disadvantage can be remedied by thresholding (ie, setting  $n_t = 1$  for all  $n_t$  greater than some threshold, and the rest setting to zero), or by considering the magnitude of a partial spike at time  $t$  as the probability of a spike occurring during frame  $t$ . Note that replacing a Poisson with an exponential is a common approximation technique in the machine learning literature [21, 20], as the exponential distribution is the closest convex relaxation to its non-convex counterpart, the Poisson distribution. More specifically, the probability mass function of a Poisson distributed random variable with low rate is very similar to the probability density function

of a random variable with an exponential distribution. While this convex relaxation makes the problem tractable, the “sharp” threshold imposed by the non-negativity constraint prohibits the use of standard gradient ascent techniques. This problem may be rectified by dropping the sharp threshold, and adding a barrier term, which must approach  $-\infty$  as  $n_t$  approaches zero (this approach is often called an “interior-point” method). Iteratively reducing the weight of the barrier term guarantees convergence to the correct solution. Thus, the goal is to efficiently solve:

$$\hat{n}_z = \operatorname{argmax}_{n_t \forall t} \sum_{t=1}^T \left( -\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 - n_t \lambda \Delta + z \log n_t \right). \quad (12)$$

where  $\log n_t$  is the “barrier term”, and  $z$  is the weight of the barrier term. Iteratively solving for  $\hat{n}_z$  for  $z$  going from one down to nearly zero, guarantees convergence to  $\hat{n}$  [21]. Since spikes and calcium are related to one another via a simple linear transformation, namely,  $n_t = C_t - \gamma C_{t-1}$ , Eq. (12) may be rewritten in terms of  $C$ :

$$\hat{C}_z = \operatorname{argmax}_{C_t - \gamma C_{t-1} \geq 0 \forall t} \sum_{t=1}^T \left( -\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 - (C_t - \gamma C_{t-1}) \lambda \Delta + z \log(C_t - \gamma C_{t-1}) \right). \quad (13)$$

The concavity of Eq. (13) facilitates utilizing any number of techniques guaranteed to find the global maximum. That the argument of Eq. (13) is twice analytically differentiable, allows for the use of the Newton-Raphson technique, which typically requires fewer steps than the gradient by itself [21]; and the special tridiagonal structure of the Hessian (as described below) enables each Newton-Raphson step to be very efficient. To proceed, Eq. (13) can be rewritten in matrix notation. Note that  $MC = n$ :

$$MC = \begin{bmatrix} 1 & 0 & 0 & \cdots & \cdots \\ 1 & -\gamma & 0 & \cdots & \cdots \\ 0 & 1 & -\gamma & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 & -\gamma \end{bmatrix} \begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ \vdots \\ C_T \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ \vdots \\ n_T \end{bmatrix} = n \quad (14)$$

where  $M \in \mathbb{R}^{T \times T}$  is a bidiagonal matrix. Then, letting  $\mathbf{1}$  be a  $T$  dimensional column vector and  $\lambda = \lambda \Delta \mathbf{1}^\top$ , yields:

$$\hat{C}_z = \operatorname{argmax}_{MC \geq 0} -\frac{1}{2\sigma^2} \|F - \alpha(C + \beta)\|^2 - (MC)^\top \lambda + z \log(MC)^\top \mathbf{1}, \quad (15)$$

where  $MC \geq 0$  indicates that every element of  $MC$  is greater than or equal to zero,  $^\top$  indicates the transpose operation, and  $\log(\cdot)$  indicates an element-wise logarithm. Now, when using Newton-Raphson to ascend a gradient, one iteratively computes both the gradient (first derivative) and Hessian (second derivative) of the argument to be minimized, with respect to the variables of interest ( $C_z$  here). Then, the estimate is updated using  $C_z \leftarrow C_z + s\mathbf{d}$ , where  $s$  is the step size. Solving  $H\mathbf{d} = \mathbf{g}$  provides  $\mathbf{d}$ , the step direction. The gradient,  $\mathbf{g}$ , and Hessian,  $H$ , are given by:

$$\mathbf{g} = -\frac{\alpha}{\sigma^2} (F - \alpha(\hat{C}_z^\top + \beta)) + M^\top \lambda - z M^\top (M \hat{C}_z)^{-1} \quad (16a)$$

$$H = \frac{\alpha^2}{\sigma^2} I + z M^\top (M \hat{C}_z)^{-2} M \quad (16b)$$

where the exponents indicate element-wise operations. The step size,  $s$ , is found using “backtracking linesearches”, which finds the maximal  $s$  that decreases the likelihood and is between zero and one.

Typically, implementing Newton-Raphson requires inverting the Hessian, i.e.,  $\mathbf{d} = H^{-1}\mathbf{g}$ , a computation that scales *cubically* with  $T$ , i.e., requires on the order of  $T^3$  operations. Already, this would be a drastic improvement over the most efficient algorithm assuming Poisson spikes, which require  $k^T$  operations (where  $k$  is the maximum number of spikes per frame). Here, because  $M$  is bidiagonal, the Hessian is tridiagonal, the solution may be found in about  $T$  operations, via standard banded Gaussian elimination techniques (which can be implemented efficiently in Matlab using  $H \setminus \mathbf{g}$ , assuming  $H$  is represented as a sparse matrix). In other words, the above approximation and inference algorithm reduces computations from *exponential* time to *linear* time. Appendix A contains pseudocode for this algorithm, including learning the parameters, as described below.

## 2.4 Learning the parameters

In the above, the parameters governing the model,  $\theta = \{\alpha, \beta, \sigma, \gamma, \lambda\}$ , were assumed to be known. An algorithm to estimate the most likely parameters,  $\hat{n}$ , proceeds as follows: (i) initialize some estimate of the parameters,  $\hat{\theta}$ , then (ii) recursively compute  $\hat{n}$  using those parameters and update  $\hat{\theta}$  given  $\hat{n}$ , and (iii) stop recursing when some convergence criteria is met.

### 2.4.1 Initializing the parameters

Because the above model is linear, the scale of  $F$  relative to  $n$  is arbitrary, so  $\alpha$  can be fixed at one without loss of generality. The offset, however, is relative but is not arbitrary. Because spiking is assumed to be sparse,  $F$  tends to be around baseline,  $\beta$  is initialized to be the median of  $F$ , and  $\sigma$  is initialized as median absolute deviation of  $F$ , ie,  $\sigma = \text{median}_t(|F_t - \text{median}(F)|)/K$ , where  $K$  is the correction factor when using median absolute deviation as a robust estimator of the standard deviation. Because previous work has shown that results are somewhat robust to minor variations in time constant [22],  $\gamma$  is initialized at  $1 - \Delta/(1\text{sec})$ , which is fairly typical [23]. Finally,  $\lambda$  is initialized at 1 Hz, which is between typical baseline firing rates and evoked spike rate activity, for these data-sets.

### 2.4.2 Estimating the parameters given $\hat{n}$

Ideally, one could integrate out the hidden variables, to find the most likely parameters:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \iint P[F, C, n | \theta] dC dn = \underset{\theta}{\operatorname{argmax}} \int P[F | C; \theta] dC \int P[n | \theta] dn. \quad (17)$$

However, evaluating those integrals is not currently tractable (to our knowledge). Therefore, Eq. (17) is approximated by simply maximizing the parameters given the MAP estimate of the hidden variables:

$$\hat{\theta} \approx \underset{\theta}{\operatorname{argmax}} P[F, \hat{C}, \hat{n} | \theta] = \underset{\theta}{\operatorname{argmax}} P[F | \hat{C}; \theta] P[\hat{n} | \theta] \quad (18)$$

where  $\hat{C}$  and  $\hat{n}$  are determined using the above described inference algorithm. The approximation in Eq. (18) is good whenever the likelihood is very peaky, meaning that most of the mass is around the MAP sequence.<sup>1</sup> The argument from the right-hand-side of Eq. (18) may be expanded:

$$P[F | \hat{C}; \theta] P[\hat{n} | \theta] = \prod_{t=1}^T P[F_t | \hat{C}_t; \beta, \sigma] P[\hat{n}_t | \lambda]. \quad (19)$$

where  $\gamma$  is not present because its posterior is relatively flat, ie, large changes in  $\gamma$  result in relatively small changes in the posterior, so it is not estimated recursively.  $\beta$  is estimated using:

$$\hat{\beta} = \underset{\beta > 0}{\operatorname{argmax}} \prod_{t=1}^T P[F_t | \hat{C}_t; \beta, \sigma] = \underset{\beta > 0}{\operatorname{argmax}} \sum_{t=1}^T (F_t - (C_t + \beta))^2, \quad (20)$$

which is solved by letting  $\hat{\beta} = \frac{1}{T} \sum_t (F_t - C_t)$ .  $\sigma$  is then the root-mean-square of the residuals, and  $\hat{\lambda}$  is the mean of  $\hat{n}$ .

### 2.4.3 Convergence criteria

Iterations stop whenever (i) iteration number exceeds some upper bound, or (ii) relative change in likelihood does not exceed some lower bound. In practice, parameters tend to converge after several iterations, given the above initializations.

<sup>1</sup>Eq. (18) may be considered a first-order Laplace approximation

## 2.5 Spatial filtering

In the above, the raw movie of fluorescence measurements collected by the experimenter had undergone two stages of preprocessing before filtering. First, the movie was segmented, to determine regions-of-interest (ROIs), yielding a vector,  $\vec{F}_t = (F_{1,t}, \dots, F_{N_p,t})$ , which corresponded to the fluorescence intensity at time  $t$  for each of the  $N_p$  pixels in the ROI. Second, at each time  $t$ , that vector was projected into a scalar, yielding  $F_t$ , the assumed input. In this section, the optimal projection is determined by considering a more general model:

$$F_{x,t} = \alpha_x(C_t + \beta) + \sigma \varepsilon_{x,t}, \quad \varepsilon_{x,t} \stackrel{iid}{\sim} \mathcal{N}(0, 1) \quad (21)$$

where  $\alpha_x$  scales each pixel, from which some number of photons are contributed due to calcium fluctuations,  $C_t$ , and others due to baseline fluorescence,  $\beta$ . Further, the noise is assumed to be both spatially and temporally white, with variance,  $\sigma^2$ , in each pixel (this assumption can easily be relaxed, by making the covariance matrix diagonal, for example). Performing inference in this more general model proceeds nearly identical manner as before. Rewriting in vector notation:

$$\hat{C}_z = \operatorname{argmax}_{\mathbf{MC} \geq \mathbf{0}} -\frac{1}{2\sigma^2} \left\| \vec{F} - \vec{\alpha}(\mathbf{C}^\top + \beta \mathbf{1}^\top) \right\|^2 - (\mathbf{MC})^\top \boldsymbol{\lambda} + z \log(\mathbf{MC})^\top \mathbf{1}, \quad (22)$$

$$\mathbf{g} = \frac{\vec{\alpha}}{\sigma^2} (\vec{F} - \vec{\alpha}(\hat{C}_z^\top + \beta)) - \mathbf{M}^\top \boldsymbol{\lambda} + z \mathbf{M}^\top (\mathbf{M} \hat{C}_z)^{-1} \quad (23)$$

$$\mathbf{H} = -\frac{\vec{\alpha}^\top \vec{\alpha}}{\sigma^2} \mathbf{I} - z \mathbf{M}^\top (\mathbf{M} \hat{C}_z)^{-2} \mathbf{M} \quad (24)$$

where  $\vec{F}$  is an  $N_p \times T$  element matrix,  $\vec{\alpha}$  is column vectors of length  $N_p$ , and  $\mathbf{I}$  is an  $N_p \times N_p$  identity matrix. Typically, the spatial filter,  $\vec{\alpha}$  is unknown, and therefore must be estimated from the data. In practice, letting  $\vec{\alpha} = \frac{1}{T} \sum_t \vec{F}_t$  was improved results over a naïve averaging of the whole field. Given this alpha,  $\beta$  could be estimated using  $\hat{\beta} = \frac{1}{T} \sum_t \sum_x \frac{F_{x,t}}{\alpha_x} - C_t$ .

## 2.6 Overlapping spatial filters

In the above, the image was segmented before filtering, such that only a single neuron was within each ROI. However, segmentation is itself a difficult problem [24]. Therefore, using a crude segmentation technique, that might not actually produce ROIs with only a single cell, and then building spatial filters for each neuron in the ROI, might be more efficient and improve SNR. This requires a minor modification to Eq. (21). Specifically, letting the superscript  $i$  index the  $N_c$  neurons in this ROI, yields:

$$\vec{F}_t = \sum_{i=1}^{N_c} \vec{\alpha}^i (C_t^i + \beta^i) + \sigma \vec{\varepsilon}_t, \quad \vec{\varepsilon}_t \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25)$$

$$C_t^i = \gamma^i C_{t-1}^i + n_t^i, \quad n_t^i \stackrel{iid}{\sim} \text{Poisson}(n_t^i; \lambda_i \Delta) \quad (26)$$

where each neuron is implicitly assumed to be independent, and each pixel is independent and identically distributed with variance  $\sigma^2$ . To perform inference in this more general model, let:

$$\mathbf{n} = [n_1^1, n_1^2, \dots, n_1^{N_c}, n_2^1, \dots, n_T^{N_c}]^\top \quad (27)$$

$$\mathbf{C} = [C_1^1, C_1^2, \dots, C_1^{N_c}, n_2^1, \dots, C_T^{N_c}]^\top \quad (28)$$

$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 0 & 0 & \dots & \dots \\ 1 & -\gamma^1 & 1 & -\gamma^2 & \dots & 1 & -\gamma_{N_c} & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 \dots & 1 & -\gamma^{N_c-1} & 1 & -\gamma^{N_c} & \dots & \dots \end{bmatrix} \quad (29)$$

and proceed as above, making minor algorithmic adjustments to deal with dimensionality issues. Since the parameters will be unknown, they must be estimated. Define  $\boldsymbol{\alpha}_x = [\alpha_x^1, \dots, \alpha_x^{N_c}]^\top$  and  $\boldsymbol{\beta} = [\beta^1, \dots, \beta^{N_c}]^\top$ . To initialize, let  $\boldsymbol{\beta} = \mathbf{0}$ ; then, perform principal component analysis on the movie,  $\vec{F}$ . Because PCA finds a low-dimensional,

orthogonal representation of the data, if a neuron in the ROI is relatively uncorrelated from the others, then a principal component is expected to correspond to its spatial filter. If PCA is insufficient to separate the spatial filters for each neuron, then the first  $N_c$  principal components are used to initialize the  $N_c$  spatial filters. Then, the following two steps are iterated. Given  $\beta$ , each  $\hat{\alpha}_x$  can be updated using:

$$\hat{\alpha}_x = \underset{\alpha_x}{\operatorname{argmin}} \sum_{t=1}^T \left( F_{x,t} \sum_{i=1}^{N_c} \alpha_x^i (C_t^i + \beta^i) \right)^2 \quad (30)$$

$$\hat{\beta} = \quad (31)$$

which is solved efficiently in Matlab using  $\alpha_x = (C + \tilde{\beta}) \backslash F_x$ , where  $\tilde{\beta}$  is  $\beta$  reparameterized to be the same size as  $C$ . Given this estimate of  $\alpha_x$  for all  $x$ , update  $\hat{\beta}^i = \frac{1}{T} \sum_t (\vec{F}_t / \vec{\alpha}^i - C_t^i)$  for each  $i$ , where  $/$  indicates an element-wise division. In practice, iterating these two steps converged after several iterations, assuming enough spikes were present in the two neurons, and they were sufficiently uncorrelated. Note that this approach — iteratively updating each  $\alpha_x$  — is more efficient than solving for the entire vector, due to the quadratic structure of the problem.

## 2.7 Experimental Methods

### 2.7.1 Slice Preparation and Imaging

All animal handling and experimentation was done according to the National Institutes of Health and local Institutional Animal Care and Use Committee guidelines. Somatosensory thalamocortical or coronal slices 350 – 400  $\mu\text{m}$  thick were prepared from C57BL/6 mice at age P14 as described [25]. Neurons were filled with 50  $\mu\text{M}$  Oregon Green Bapta 1 hexapotassium salt (OGB-1; Invitrogen, Carlsbad, CA) through the recording pipette or bulk loaded with Fura-2 AM (Invitrogen, Carlsbad, CA). Pipette solution contained 130 K-methylsulfate, 2  $\text{MgCl}_2$ , 0.6 EGTA, 10 HEPES, 4 ATP-Mg, and 0.3 GTP-Tris, pH 7.2 (295 mOsm). After cells were fully loaded with dye, imaging was done by using a modified BX50-WI upright microscope (Olympus, Melville, NY). Image acquisition was performed with the C9100-12 CCD camera from Hamamatsu Photonics (Shizuoka, Japan) with arclamp illumination at 480 – 500 nm and 510 – 550 nm collection filters (Chroma, Rockingham, VT) for Oregon Green. Imaging of Fura-2 loaded slices was performed with a confocal spinning disk (Solamere Technology Group, Salt Lake City, UT) and an Orca CCD camera from Hamamatsu Photonics (Shizuoka, Japan). Images were saved and analyzed using custom software written in Matlab (Mathworks, Natick, MA).

### 2.7.2 Electrophysiology

All recordings were made using the Multiclamp 700B amplifier (Molecular Devices, Sunnyvale, CA), digitized with National Instruments 6259 multichannel cards and recorded using custom software written using the LabView platform (National Instruments, Austin, TX). Square pulses of sufficient amplitude to yield the desired number of action potentials were given as current commands to the amplifier using the LabView and National Instruments system.

### 2.7.3 Fluorescence preprocessing

Traces were extracted using custom Matlab scripts to first segment the mean image into ROIs, and then all pixels within the ROI were averaged. The Fura 2 fluorescence traces were inverted. As some slow drift was sometimes present in the traces, the lowest ten frequency components were set to zero (ten was chosen by eye), and the resulting fluorescence trace was then normalized to be between zero and one.

## 3 Results

### 3.1 Main Result

The main result of this paper is that the fast filter can approximate  $\hat{n}$  very efficiently, and that this approach outperforms a more naïve approach of a typical deconvolution filter, half-wave rectified (i.e., setting everything below zero equal to zero). Fig. 2 depicts a simulation showing this result. Clearly, the fast filter is outperforming the optimal linear



deconvolution filter (also called a Wiener filter). The Wiener filter implicitly approximates the Poisson spike rate with a Gaussian spike rate (see Appendix B for details). While a Gaussian well approximates a Poisson distribution when rates are about 10 spikes per frame, this example is very far from that regime, and so the Gaussian approximation performs relatively poorly. Furthermore, the Gaussian approximation allows for the inferred spike train to include negative numbers, which is undesirable, as spike trains are non-negative entities. To counteract the negative values, the Wiener filter then infers large positive values, contributing to a “ringing” effect. The non-negative constraint imposed by the fast filter ensures that such ringing does not take place. Finally, by utilizing Gaussian elimination and interior-point methods, as described in the Methods section, the computational complexity of fast filter is the same as an efficient implementation of the Wiener filter.

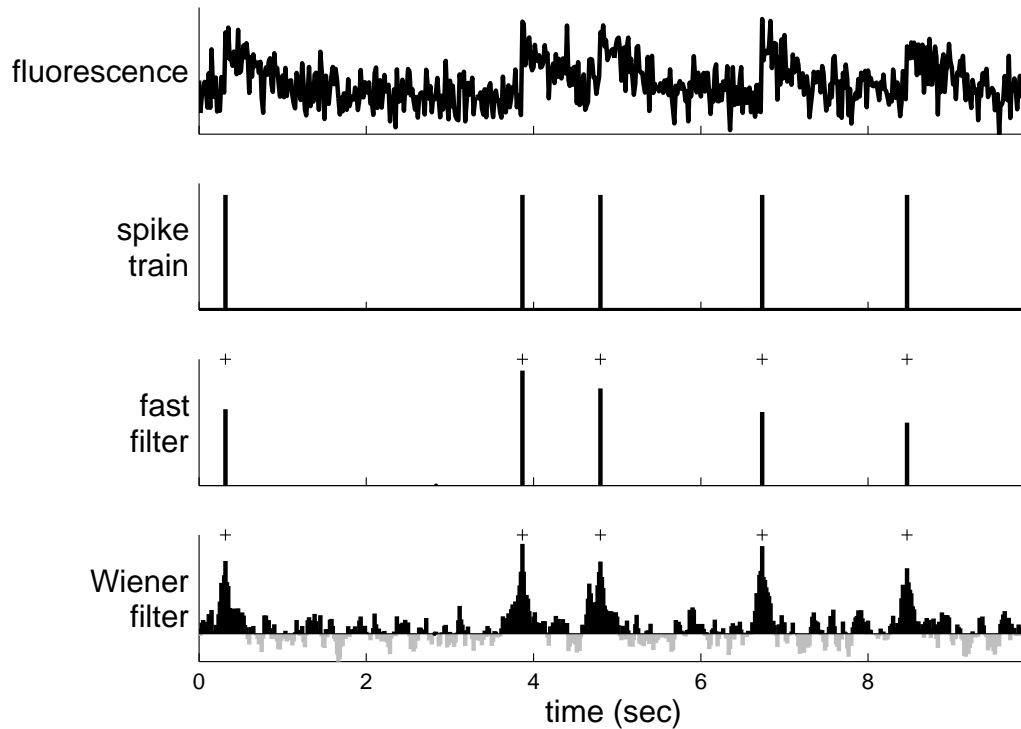


Figure 2: The fast filter significantly outperforms the optimal linear deconvolution (aka, Wiener filter) on typical simulated data. Top panel: fluorescence trace. Second panel: spike train. Third panel: fast filter inference. Bottom panel: Wiener filter inference. Note that the gray bars indicate *negative* spikes, which, of course, are meaningless. Black '+'s in bottom two panels indicate true spike times. Simulation details:  $T = 2930$  time steps,  $\Delta = 5$  msec,  $\alpha = 1$ ,  $\beta = 0$ ,  $\sigma = 0.3$ ,  $\tau = 1$  sec,  $\lambda = 1$  Hz.

Although in Figure 2 the model parameters were provided, in the general case, the parameters are unknown, and must therefore be estimated from the data. section 2.4 describes how the parameters of the model may be estimated directly from the observations. Importantly, this obviates the need to conduct joint imaging and electrophysiological experiments to obtain “training” data, as the developed approach is fully unsupervised. Figure 3 shows another simulated example; in this example, however, the parameters are estimated from the observed fluorescence trace alone. Again, it is clear that the fast filter far outperforms the Wiener filter.

Given the above two results, the fast filter was applied to real data. More specifically, by jointly recording electrophysiologically and imaging, the true spike times are known, and the accuracy of the two filters can be compared. Figure 4 shows a result typical of the 12 joint electrophysiological and imaging experiments conducted. Note that the first few “events” are actually pairs of spikes, which is reflected in the inferred spike trains. This suggests that although the scale is arbitrary, the fast filter can correctly ascertain the number of spikes within spike events.

Figure 5 further evaluates this claim. The cell was forced to spike in “doublets” for each spiking event, meaning

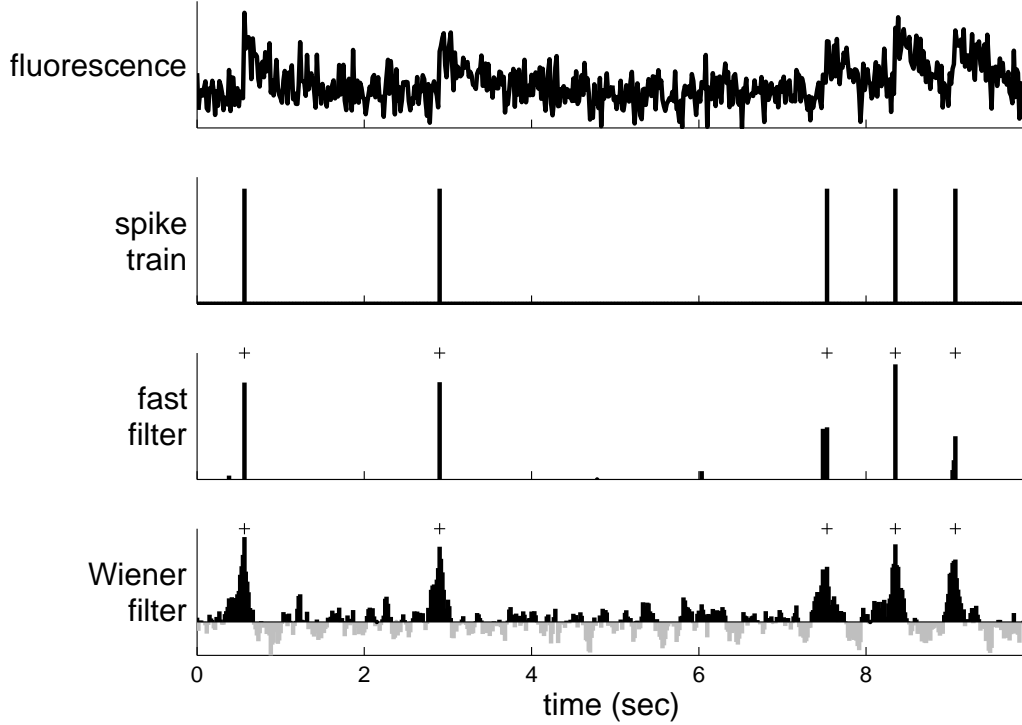


Figure 3: The fast filter significantly outperforms the Wiener filter, even when estimating the parameters only from the observed simulated data. Simulation details as in Figure 2.

that the algorithm never saw a single spike event. Naïvely, this would suggest that a purely linear approach would struggle to resolve spike fidelity in this high frequency spiking regime. However, the fast filter correctly identifies nearly all spiking events as “doublets”. On the contrary, there is no obvious way to count the number of events upon using the Wiener filter.

### 3.2 Online analysis of spike trains using the fast filter

A central aim for this work was the development of an algorithm that infers spikes fast enough to use online while imaging a large population (eg,  $\approx 100$ ) of neurons. Figure 6 shows a segment of the results of running the fast filter on 136 neurons, recorded simultaneously, as described in section 2.7. Note that the filtered fluorescence signals much more clearly show fluctuations in spiking. These spike trains were inferred in super-real-time, meaning that one could infer spike trains for the past experiment while conducting the subsequent experiment. More specifically, a trace with 50,000 frames can be analyzed in about one second on a standard desktop computer.

### 3.3 Extensions

section 2.1 describes a simple principled first-order model relating the spike trains to the fluorescence trace. A number of the simplifying assumptions can be straightforwardly relaxed, as described below.

#### 3.3.1 Replacing Gaussian observations with Poisson

In the above, observations had a Gaussian distribution. The statistics of photon emission and counting, however, suggest that a Poisson distribution would be more natural [26], yielding:

$$F_t \stackrel{iid}{\sim} \text{Poisson}(\alpha C_t + \beta). \quad (32)$$

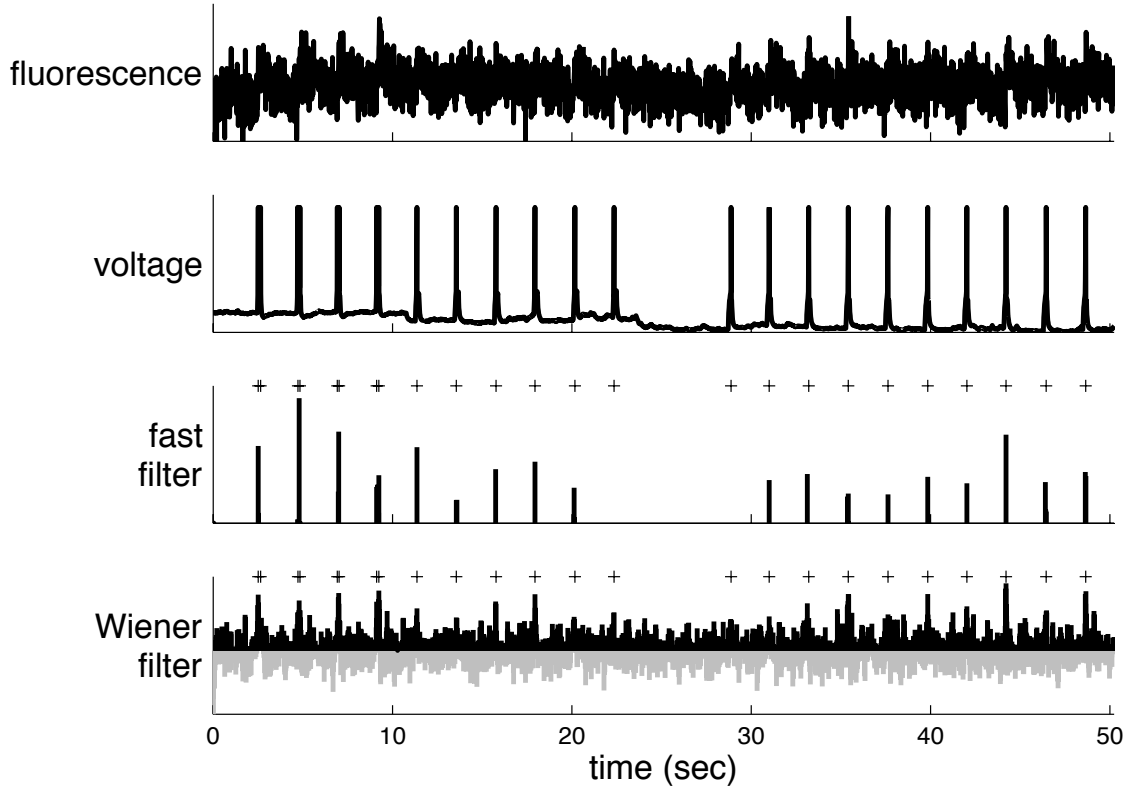


Figure 4: The fast filter significantly outperforms the Wiener filter on typical in vitro data, using Fura 2. Note that all the parameters for both filters were estimated from the fluorescence data alone (ie, not considering the voltage data at all).

One advantage to this model, over the Gaussian one, besides accuracy, is that the variance parameter no longer exists, which might make learning the parameters simpler. Importantly, the posterior is still concave in  $C$ , so the same techniques can be used for this model (which requires analytically calculating the gradient and Hessian for the likelihood term implied by Eq. (32)). In practice, modifying the filter for this model extension did not seem to improve inference results in any data (not shown).

### 3.3.2 Allowing for a time-varying prior

In Eq. (4), the rate of spiking is a constant. Often, additional knowledge about the experiment, including external stimuli, or other neurons spiking, can provide strong time-varying prior information. A simple model modification can incorporate that feature:

$$n_t \stackrel{iid}{\sim} \text{Poisson}(\lambda_t \Delta), \quad (33)$$

where  $\lambda_t$  is now a function of time. Approximating this time-varying Poisson with a time-varying exponential (as in Eq. (10)) yields a log-concave problem again, so the same techniques can be used to find the most likely spike train, give this prior. However, as above, this model extension did not yield any improved filtering results (not shown).

### 3.3.3 Saturating fluorescence

Although all the above models assumed a *linear* relationship between  $F_t$  and  $C_t$ , the relationship between fluorescence and calcium is typically better approximated by certain nonlinear functions, like the Hill equation [23];

$$F_t = \alpha \frac{C_t}{C_t + k_d} + \beta + \sigma \varepsilon_t, \quad \varepsilon_t \stackrel{iid}{\sim} \mathcal{N}(0, 1). \quad (34)$$

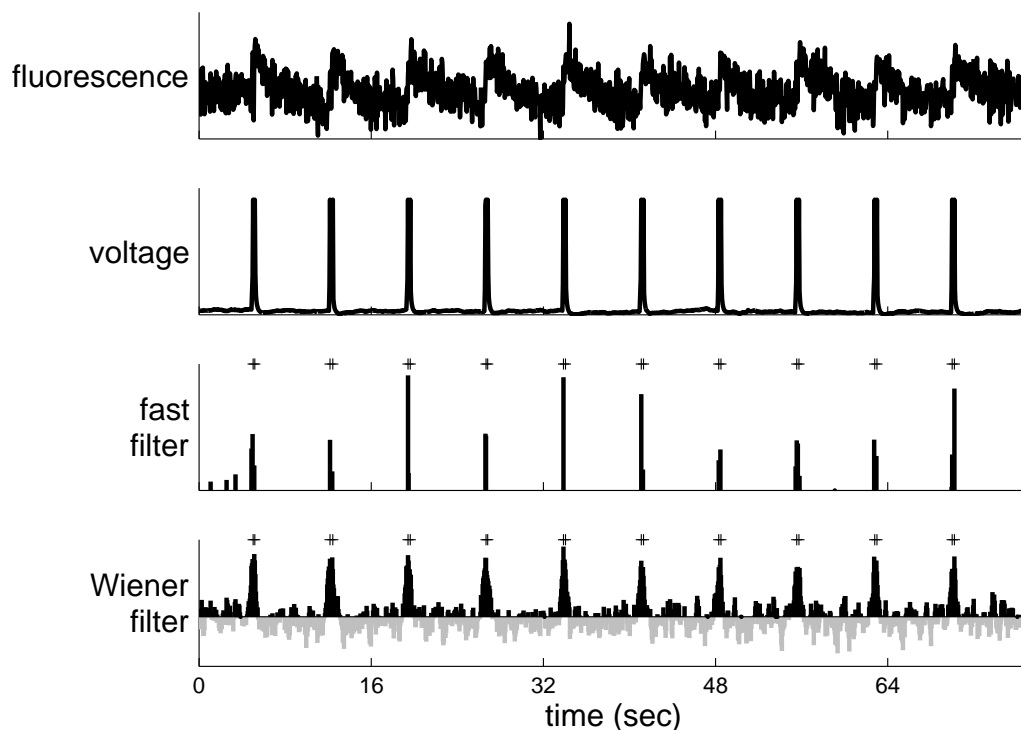


Figure 5: The fast filter can resolve a sequence of doublets, even in the absence of ever seeing an individual spike. It is difficult, if not impossible, to count the number of spikes given the Wiener filter output. Recording and fitting parameters as in Figure 4

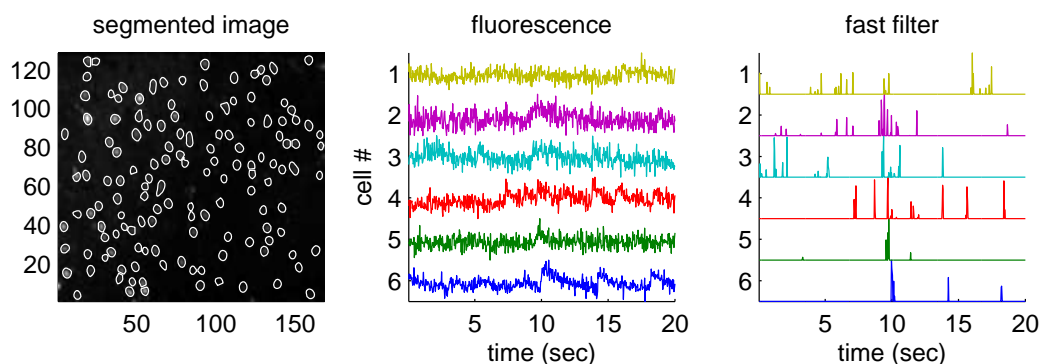


Figure 6: The fast filter infers spike trains from a large population of neurons imaged simultaneously in vitro, using OBG-1, in super-real-time. The inferred spike trains convey much more clearly the neural activity. Left panel: Mean segmented image field. Middle panel: example fluorescence traces. Right panel: fast filter output corresponding to each associated trace.

Incorporating this nonlinearity breaks the log-concavity of the posterior, meaning that converging to the global maximum is no longer guaranteed. Assuming a good initialization can be found, however, if this model is more accurate, then ascending the gradient for this model might yield improved inference results. In practice, initializing with the linear solution often resulted in nearly equally good inference, but this nonlinear inference was far less robust than the inference assuming the linear model (not shown).

### 3.3.4 Using the fast filter to initialize the SMC filter

A previously proposed sequential Monte Carlo (SMC) method to infer spike trains [12] can incorporate this nonlinearity, as well as the other model extensions discussed above. However, this SMC filter is not nearly as efficient as the fast filter proposed here. The SMC filter estimates the model parameters in a completely unsupervised fashion, ie, from the fluorescence observations, using an expectation-maximization algorithm (which requires iterating between computing the expected value of the hidden variables —  $C$  and  $n$  — and updating the parameters). In [12], parameters for the SMC filter were initialized based on other data. While effective, this initialization was often far from the final estimates, and therefore, required a relatively large number of iterations (eg, 20–25) before converging. Thus, it seemed that the fast filter could be used to obtain an improvement to the initial parameter estimates, reducing the required number of iterations. Indeed, Figure 7 shows how the SMC filter outperforms the fast filter on in vitro data, and only required 3–5 iterations to converge. Note that the first few events are individual spikes, resulting in relatively small fluorescence fluctuations, whereas the next events are actually spike doublets, causing a much larger fluorescence fluctuation. Only the SMC filter picks up the individual spikes in this trace, a result typical when the effective signal-to-noise ratio (SNR) is so poor. Thus, these two inference algorithms are complementary: the fast filter can be used for rapid, online inference, and for initializing the SMC filter, which can then be used to further refine the spike train estimate. Importantly, although the SMC filter often outperforms the fast filter, the fast filter is more robust, meaning that it more often works “out-of-the-box”.

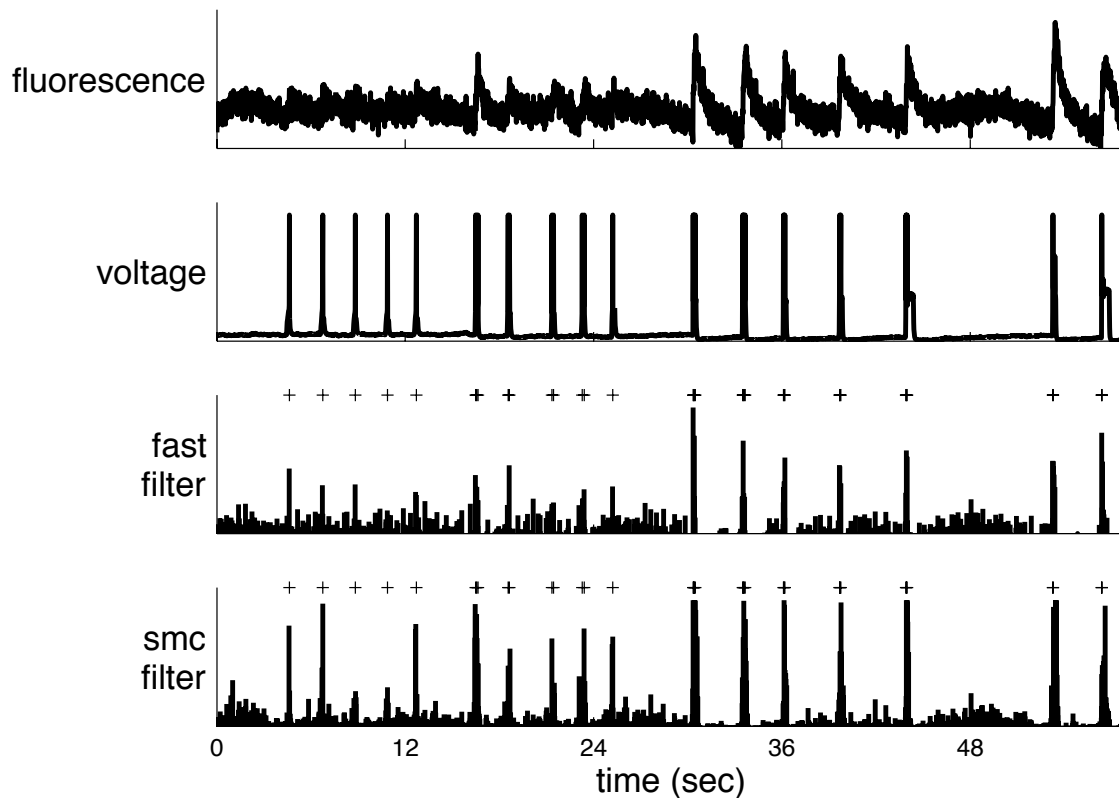


Figure 7: The fast filter effectively initializes the parameters for the SMC filter (which outperforms the fast filter), significantly reducing the number of expectation-maximization iterations to convergence, on typical in vitro data, using Fura 2. Note that the ordinate on the bottom panel corresponds to the inferred probability of a spike having occurred in each frame.

## 3.4 Spatial filter

In the above, the filter operated on one-dimensional fluorescence traces. Typically the data are a time-series of images which are first segmented into regions-of-interest (ROI), and then (usually) averaged to obtain  $F_t$ . In theory, one could

improve the effective SNR of the fluorescence trace by scaling each pixel relative to one another. In particular, pixels not containing any information about calcium fluctuations can be ignored, and pixels that are partially anti-correlated with one another could have weights with opposing signs.

Figure 8 demonstrates the potential utility of this approach. The top row shows different depictions of an ROI containing a single neuron. On the far left panel is the true spatial filter for this neuron. This particular spatial filter was chosen based on experience analyzing both in vitro and in vivo movies; often, it seems that the pixels immediately around the soma are anti-correlated with those in the soma. This effect is possibly due to the influx of calcium from extracellular space immediately around the soma. This simulated movie is relatively noisy, as indicated by the second panel, which depicts an exemplary image frame. The standard approach, given such a noisy movie, would be to first segment the movie to find an ROI corresponding to the soma of this cell, and then spatially average all the pixels found to be within this ROI. The third panel shows this “typical spatial filter”. The fourth panel shows the mean frame. Clearly, this mean frame is very similar to the true spatial filter. When the true spatial filter is a single Gaussian, the typical spatial filter works about as well as the true one (not shown).

The bottom panels of Figure 8 depict the effect of using the true spatial filter, versus the typical one. The left side shows the fluorescence trace and its associated spike inference obtained from using the typical spatial filter. The right side shows the same when using the true spatial filter. Clearly, the true spatial filter results in a much cleaner fluorescence trace and spike inference.

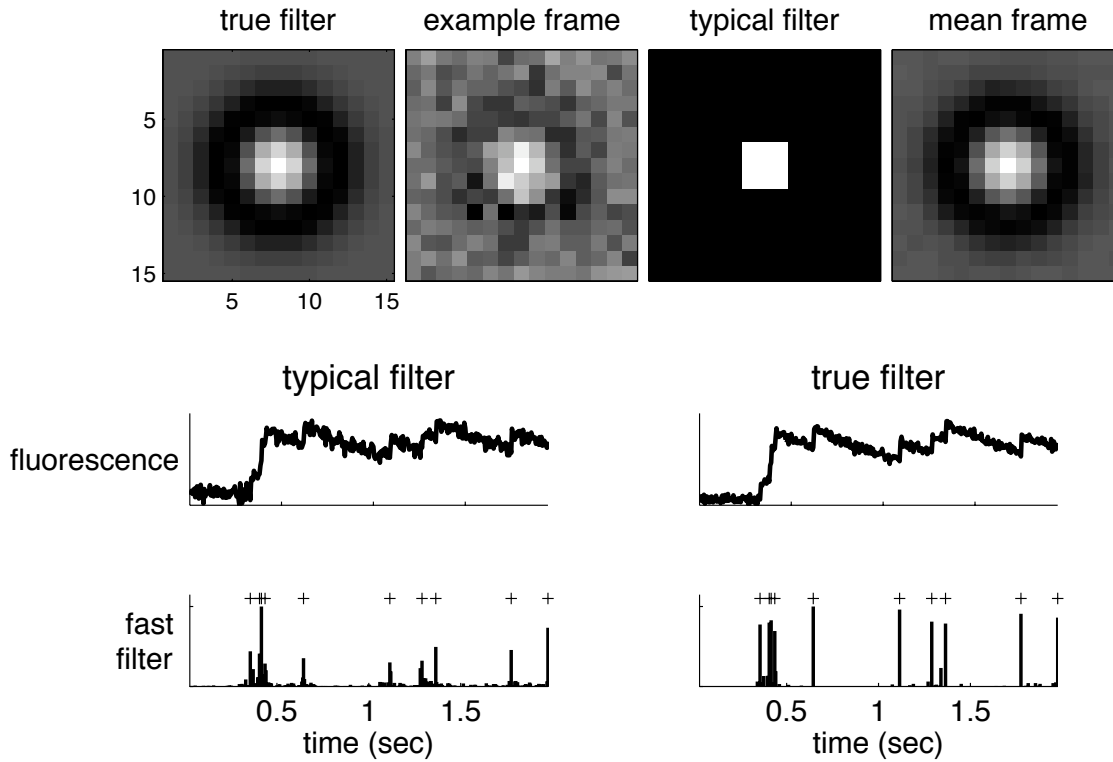


Figure 8: A simulation demonstrating that using a better spatial filter can significantly enhance the effective SNR (see Supplementary Movie 1 for the full movie associated with this simulation). The true spatial filter was a sum of Gaussians: a positively weighted small variance Gaussian, and a negatively weighted large variance Gaussian (both with the same mean). Top row far left: true spatial filter. Top row second from left: example frame (frame number 100). Top row second from right: typical spatial filter. Top row far right: mean frame. Middle row left: fluorescence trace using typical spatial filter. Bottom row left: fast filter output using typical spatial filter. Middle row right: fluorescence trace using true spatial filter. Bottom right: fast filter output using true spatial filter. Simulation details:  $\vec{\alpha} = \mathcal{N}(\mathbf{0}, 2\mathbf{I}) - 1.1\mathcal{N}(\mathbf{0}, 2.5\mathbf{I})$  where  $\mathcal{N}(\mathbf{mu}, \mathbf{\Sigma})$  indicates a Gaussian with mean  $\mathbf{\mu}$  and covariance matrix  $\mathbf{\Sigma}$ ,  $\beta = 1$ ,  $\tau = 0.85$  sec,  $\lambda = 5$  Hz.

### 3.5 Overlapping spatial filters

The above shows that if a ROI contains only a single neuron, the effective SNR can be enhanced by spatially filtering. However, this analysis assumes that only a single neuron is in the ROI. Often, neural spatial filters are overlapping, or nearly overlapping, making the segmentation problem even more difficult. Therefore, it is desirable to have an ability to crudely segment, yielding only a few neurons in each ROI, and then spatially filter within each ROI to pick out the spike trains from each neuron. This may be achieved in a principled manner by generalizing the model as described in section 2.6. Figure 9 shows an example of this approach on simulated data.

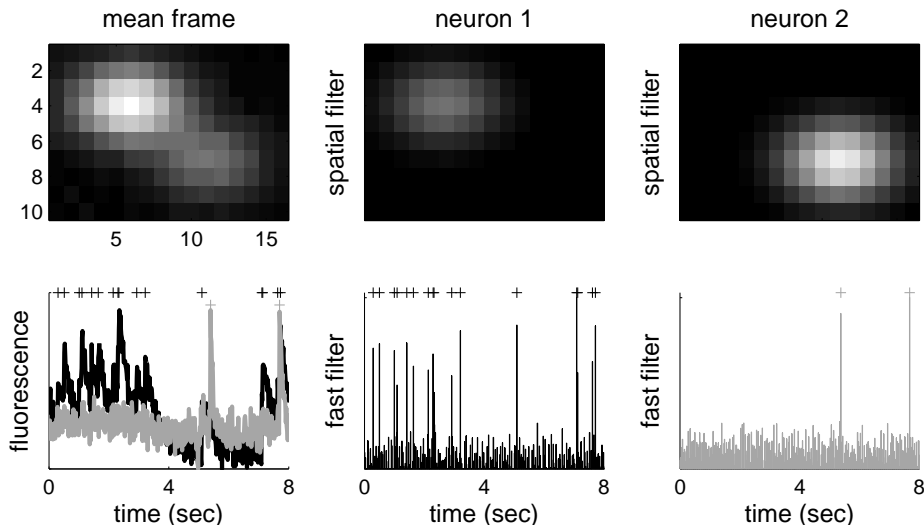


Figure 9: Simulation showing that even when two neurons’ spatial filters are overlapping, one can separate the two signals by spatial filtering. Simulation details:  $\vec{\alpha}^1 = \mathcal{N}([-1.81.8], 2\mathbf{I})$ ,  $\vec{\alpha}^2 = \mathcal{N}([1.8 - 1.8], 5\mathbf{I})$ ,  $\beta = [11]^T$ ,  $\tau = [0.50.5]^T$  sec,  $\lambda = [1.51.5]$  Hz.

While Figure 9 shows that one could separate the signals if the spatial filters of the neurons were known, Figure 10 shows that the spatial filters can be estimated using only the fluorescence movie, by using the approach described in section 2.6.

## 4 Discussion

This work describes an algorithm that approximates the *maximum a posteriori* (MAP) spike train, given a fluorescence movie. The approximation is required because finding the actual MAP estimate is not currently computationally tractable. Replacing the assumed Poisson distribution on spikes with an exponential distribution yields a log-concave optimization problem, which can be solved using standard gradient ascent techniques (such as Newton-Raphson). This exponential distribution has an advantage over a Gaussian distribution by restricting spikes to be positive, which improves inference quality (c.f. Figure 2), and is a better approximation to a Poisson distribution with low rate. Furthermore, by utilizing the special structure of the Hessian matrix (ie, it is tridiagonal), this approximate MAP spike train can be inferred fast enough on standard computers to use it for online analyses. Finally, all the parameters can be estimated from only the fluorescence observations, obviating the need for joint electrophysiology and imaging (c.f. Figure 3). This approach is robust, in that it works “out-of-the-box” on all the in vivo and in vitro data analyzed (c.f. Figure 4).

Ideally, one could compute the full joint posterior of entire spike trains, conditioned on the fluorescence data. This distribution is analytically intractable, due to the Poisson assumption on spike trains. A Bayesian approach could use Markov Chain Monte Carlo to recursively sample spikes until a whole sample spike train is obtained [27, 28]. Because a central aim here was computational expediency, a “greedy” approach is natural: ie, recursively sample the most likely spike, update the posterior, and repeat until the posterior stops increasing. Template matching, projection pursuit regression [29], and matching pursuit [30] are examples of such a greedy approach (Greenberg et al’s algorithm could

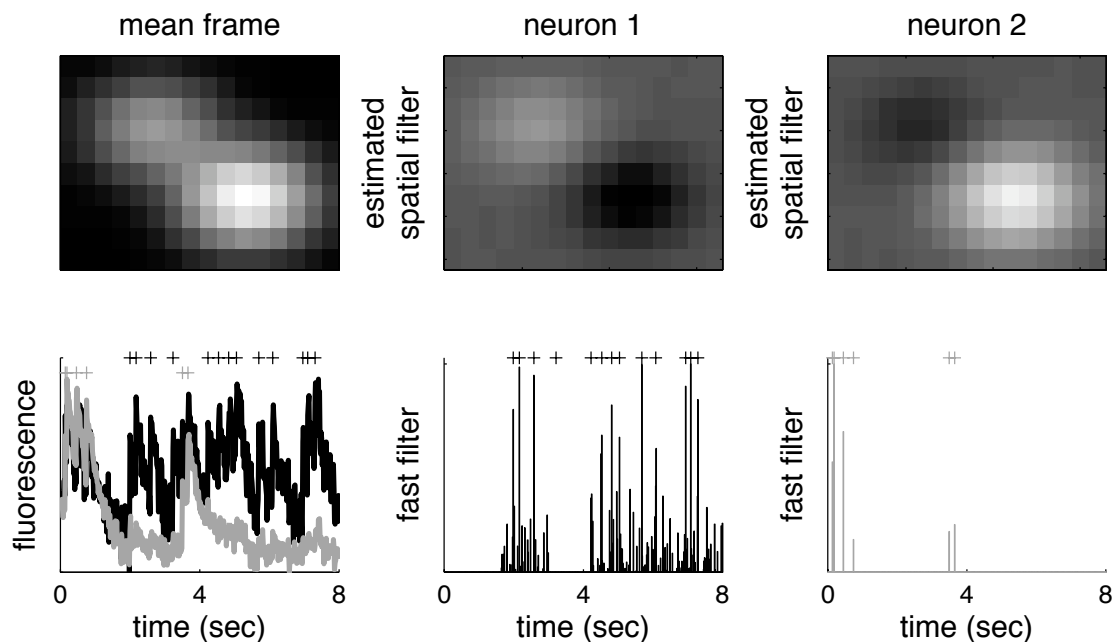


Figure 10: Simulation showing that even when two neuron’s spatial filters are largely overlapping, linearly scaled estimates of the spatial filter of each neuron can be inferred, to separate the two signals. Simulation details as above.

also be considered a special case of such a greedy approach). Both the greedy methods, and the one developed here, aim to optimize a similar objective function. While greedy methods reduce the computational burden by restricting the search space of spike trains, here analytic approximations are made. The advantage of the greedy approaches relative to this one is that they result in a spike train (ie, a binary sequence), whereas the approach developed herein is guaranteed to be approximately optimal, given the model.

Because this filter is model based, it can be generalized in several ways to improve accuracy. Unfortunately, some of these generalizations do not improve inference accuracy, probably because of the exponential approximation. Instead, the fast filter output can be used to initialize the SMC filter [12], to further improve inference quality (c.f. Figure 7). Another model generalization allows incorporation of spatial filtering of the raw movie into this approach (c.f. Figure 8). The parameters of the spatial filter can be estimated from the data, even when spatial filters are overlapping (c.f. Figure 10).

A number of extensions follow from this work. First, further development on some of the model generalizations may improve inference results. Second, putting this filter with a crude but automatic segmentation tool to obtain ROIs would create a completely automatic algorithm that converts raw movies of populations of neurons into populations of spike trains. Third, combining this algorithm with recently developed connectivity inference algorithms on this kind of data [28], could yield very efficient connectivity inference.

**Acknowledgments** The authors would like to express appreciation for helpful discussions with Vincent Bonin. Support for JTV was provided by NIDCD DC00109. LP is supported by an NSF CAREER award, by an Alfred P. Sloan Research Fellowship, and the McKnight Scholar Award. RY’s laboratory was supported by grant EY11787 from the National Eye Institute and by the Kavli Institute for Brain Science at Columbia University. LP and RY share NSF IIS-0904353.

## A Pseudocode



---

**Algorithm 1** pseudocode for inferring the approximately most likely spike train, given fluorescence data.

---

```

1: initialize parameters,  $\theta$  (section 2.4.1)
2: while convergence criteria not met do
3:   for  $z = 1, 0.1, 0.01, \dots, \eta_1$  do {interior point method to find  $\hat{C}$ }
4:     Initialize  $n_t = \eta_2$  for all  $t = 1, \dots, T$ ,  $C_1 = 0$  and  $C_t = \gamma C_{t-1} + n_t$  for all  $t = 2, \dots, T$ .
5:     let  $C_z$  be the initialized calcium, and  $\mathcal{P}_z$ , be the posterior given this initialization
6:     while  $\mathcal{P}_{z'} > \mathcal{P}_z + \eta_3$  do {Newton-Raphson with backtracking line searches}
7:       compute  $\mathbf{g}$  using Eq. (23)
8:       compute  $\mathbf{H}$  using Eq. (24)
9:       compute  $\mathbf{d}$  using  $\mathbf{H} \backslash \mathbf{g}$  {Gaussian elimination}
10:      let  $\mathbf{C}_{z'} = \mathbf{C}_z + s\mathbf{d}$ , where  $s$  is between 0 and 1, and  $\mathcal{P}_{z'} > \mathcal{P}_z$  {backtracking line search}
11:    end while
12:  end for
13:  check convergence criteria
14:  update  $\vec{\alpha}$  using Eq. (30) {only if spatial filtering}
15:  update  $\vec{\beta}$  using Eq. (31)
16:  let  $\sigma$  be the root-mean square of the residual
17:  let  $\lambda = \hat{n}^\top \mathbf{1} / T$ 
18: end while

```

---

## B Wiener Filter

The Poisson distribution can be replaced with a Gaussian instead of a Poisson distribution, ie,  $n_t \stackrel{iid}{\sim} \mathcal{N}(\lambda\Delta, \lambda\Delta)$ , which, when plugged into Eq. (7) yields:

$$\hat{n} = \operatorname{argmax}_{n_t} \sum_{t=1}^T \left( \frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + \frac{1}{2\lambda\Delta} (n_t - \lambda\Delta)^2 \right). \quad (35)$$

Note that since fluorescence integrates over  $\Delta$ , it makes sense that the mean scales with  $\Delta$ . Further, since the Gaussian here is approximating a Poisson with high rate [26], the variance should scale with the mean. Using the same tridiagonal trick as above, Eq. (35) can be solved using Newton-Raphson once (because its quadratic). Writing the above in matrix notation, substituting  $C_t - \gamma C_{t-1}$  for  $n_t$ , yields:

$$\hat{C} = \operatorname{argmax}_C \frac{1}{2\sigma^2} \|\mathbf{F} - \mathbf{C}\|^2 - \frac{1}{2\lambda\Delta} \|\mathbf{M}\mathbf{C} - \lambda\Delta\mathbf{1}\|^2, \quad (36)$$

which is quadratic in  $\mathbf{C}$ . The gradient and Hessian are given by:

$$\mathbf{g} = -\frac{1}{\sigma^2} (\mathbf{C} - \mathbf{F}) - \frac{1}{\lambda\Delta} ((\mathbf{M}\hat{C})^\top \mathbf{M} + \lambda\Delta \mathbf{M}^\top \mathbf{1}), \quad (37)$$

$$\mathbf{H} = \frac{1}{\sigma^2} \mathbf{I} + \frac{1}{\lambda\Delta} \mathbf{M}^\top \mathbf{M}. \quad (38)$$

Note that this solution is the optimal linear solution, under the assumption that spikes follow a Gaussian distribution, and if often referred to as the Wiener filter, regression with a smoothing prior, or ridge regression. Estimating the parameters for this model follows similarly as above.

## References

- [1] R. Yuste, A. Konnerth, B.R. Masters, et al. *Imaging in Neuroscience and Development, A Laboratory Manual*, 2006.
- [2] D Smetters, A Majewska, and R Yuste. Detecting action potentials in neuronal populations with calcium imaging. *Methods*, 18(2):215–221, Jun 1999.
- [3] Yuji Ikegaya, Gloster Aaron, Rosa Cossart, Dmitriy Aronov, Ilan Lampl, David Ferster, and Rafael Yuste. Synfire chains and cortical songs: temporal modules of cortical activity. *Science*, 304(5670):559–564, Apr 2004.
- [4] Shin Nagayama, Shaoqun Zeng, Wenhui Xiong, Max L Fletcher, Arjun V Masurkar, Douglas J Davis, Vincent A Pieribone, and Wei R Chen. In vivo simultaneous tracing and  $\text{Ca}^{2+}$  imaging of local neuronal circuits. *Neuron*, 53(6):789–803, Mar 2007.
- [5] Werner Göbel and Fritjof Helmchen. In vivo calcium imaging of neural network function. *Physiology (Bethesda)*, 22:358–365, Dec 2007.
- [6] L. Luo, E. M. Callaway, and K. Svoboda. Genetic dissection of neural circuits. *Neuron*, 57:634–60, —2008—.
- [7] Olga Garaschuk, Oliver Griesbeck, and Arthur Konnerth. Troponin c-based biosensors: a new family of genetically encoded indicators for in vivo calcium imaging in the nervous system. *Cell Calcium*, 42(4-5):351–361, 2007.
- [8] Marco Mank and Oliver Griesbeck. Genetically encoded calcium indicators. *Chem Rev*, 108(5):1550–1564, May 2008.
- [9] D.J. Wallace, S.M. zum Alten Borgloh, S. Astori, Y. Yang, M. Bausen, S. Kugler, A.E. Palmer, R.Y. Tsien, R. Sprengel, J.N.D. Kerr, W. Denk, and M.T. Hasan. Single-spike detection in vitro and in vivo with a genetic  $\text{Ca}^{2+}$  sensor. *Nature methods*, 5(9):797–804, 2008.
- [10] David S Greenberg, Arthur R Houweling, and Jason N D Kerr. Population imaging of ongoing neuronal activity in the visual cortex of awake rats. *Nat Neurosci*, Jun 2008.
- [11] T. Holekamp, D. Turaga, and T. Holy. Fast three-dimensional fluorescence imaging of activity in neural populations by objective-coupled planar illumination microscopy. *Neuron*, 57:661–672, 2008.
- [12] Joshua T Vogelstein, Brendon O Watson, Adam M Packer, Rafael Yuste, Bruno Jedynek, and Liam Paninski. Spike inference from calcium imaging using sequential monte carlo methods. *Biophys J*, 97(2):636–655, Jul 2009.
- [13] L.F. Portugal, J.J. Judice, and L.N. Vicente. A comparison of block pivoting and interior-point algorithms for linear least squares problems with nonnegative variables. *Mathematics of Computation*, 63(208):625–643, 1994.
- [14] J. Markham and J. A. Conchello. Parametric blind deconvolution: a robust method for the simultaneous estimation of image and blur. *Journal of The Optical Society Of America A. Optics, Image Science, and Vision*, 16(10):2377–2391, Oct 1999.
- [15] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, Oct 1999.
- [16] Y. Lin, D. Lee, and L. Saul. Nonnegative deconvolution for time of arrival estimation. *International Conference on Acoustics, Speech, and Signal Processing*, 2004.
- [17] O’Grady, P.D. and Pearlmutter, B.A. Convolutional non-negative matrix factorisation with a sparseness constraint. *Machine Learning for Signal Processing, 2006. Proceedings of the 2006 16th IEEE Signal Processing Society Workshop on*, pages 427–432, 2006.
- [18] Quentin J M Huys, Misha B Ahrens, and Liam Paninski. Efficient estimation of detailed single-neuron models. *J Neurophysiol*, 96(2):872–890, Aug 2006.

- [19] John P. Cunningham, Krishna V. Shenoy, and Maneesh Sahani. Fast Gaussian process methods for point process intensity estimation. *ICML*, pages 192–199, 2008.
- [20] Liam Paninski, Yashar Ahmadian, Daniel Ferreira, Shinsuke Koyama, Kamiar Rahnama Rad, Michael Vidne, Joshua Vogelstein, and Wei Wu. A new look at state-space models for neural data. *J Comput Neurosci*, Aug 2009.
- [21] S. Boyd and L. Vandenberghe. *Convex Optimization*. Oxford University Press, 2004.
- [22] Emre Yaksi and Rainer W Friedrich. Reconstruction of firing rate changes across neuronal populations by temporally deconvolved  $\text{Ca}^{2+}$  imaging. *Nature Methods*, 3(5):377–383, May 2006.
- [23] Thomas A Polgruto, Ryohei Yasuda, and Karel Svoboda. Monitoring neural activity and  $[\text{Ca}^{2+}]$  with genetically encoded  $\text{Ca}^{2+}$  indicators. *J Neurosci*, 24(43):9572–9579, Oct 2004.
- [24] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- [25] Jason N MacLean, Brendon O Watson, Gloster B Aaron, and Rafael Yuste. Internal dynamics determine the cortical response to thalamic stimulation. *Neuron*, 48(5):811–823, Dec 2005.
- [26] Lucas Sjulson and Gero Miesenböck. Optical recording of action potentials and other discrete physiological events: a perspective from signal detection theory. *Physiology (Bethesda)*, 22:47–55, Feb 2007.
- [27] C. Andrieu, É. Barat, and A. Doucet. Bayesian deconvolution of noisy filtered point processes. *IEEE Transactions on Signal Processing*, 49(1):134–146, 2001.
- [28] Mishchenko Y, Vogelstein JT, and Paninski L. A Bayesian approach for inferring neuronal connectivity from calcium fluorescent imaging data. *Annals of Applied Statistics*, in press, 2009.
- [29] JH Friedman and W. Stuetzle. Projection Pursuit Regression. *J. AM. STAT. ASSOC.*, 76(376):817–823, 1981.
- [30] S. Mallat and Z. Zhang. Matching pursuit with time-frequency dictionaries: *IEEE Trans. Signal Processing*, 41:3397–3415, 1993.