

Fast spike train inference from calcium imaging

Joshua T. Vogelstein, others, Liam Paninski

July 6, 2009

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 3 |
| 2 | Methods | 4 |
| 2.1 | Data driven generative model | 4 |
| 2.2 | Goal | 6 |
| 2.3 | Inference | 7 |
| 2.4 | Assessment | 9 |
| 2.5 | Experimental Methods | 9 |
| 3 | Results | 10 |
| 3.1 | Main Result | 10 |
| 3.2 | Spatial Filtering | 14 |
| 3.3 | Learning | 16 |
| 3.4 | in vitro data | 19 |
| 3.5 | Overlapping spatial filters | 20 |
| 3.6 | Population imaging | 22 |
| 3.7 | Slow rise time | 23 |
| 3.8 | Poisson observation model | 24 |
| 3.9 | in vivo data | 25 |
| 4 | Discussion | 26 |
| 4.1 | Incorporating a nonlinear observation model | 27 |
| 4.2 | Dynamic prior | 28 |
| | References | 29 |
| A | Wiener Filter | 29 |

Abstract

A fundamental desideratum in neuroscience is to simultaneously observe the spike trains from large populations of neurons. Calcium imaging technologies are bringing the field ever closer to achieving this goal, both in vitro and in vivo. To get the most information out of these preparations, one can increase the frame rate and image field, leading to corresponding increases in temporal resolution and number of observable cells. However, these increases come at the cost of reducing the dwell time per pixel, causing a decrease in the signal-to-noise ratio. Thus, to maximize the utility of these technologies, powerful computational tools must be built to compliment the experimental tools. In particular, by considering the statistics of the data — e.g., firing rates and photon counts are positive (or zero) — we develop an approximately optimal algorithm for inferring spike trains from fluorescence data. More specifically, we use an interior-point method to perform a non-negative deconvolution, inferring the approximately most likely spike train for each neuron, given their fluorescence signals. We demonstrate using simulations, in vitro, and in vivo data sets the improvement of our algorithm over other techniques. Moreover, because our inference is very fast — requiring only about 1 second of computational time on laptop to analyze a calcium trace from 50,000 image frames — we call this approach the Fast Non-negative deconvolution Spike Inference (FANSI) filter. We demonstrate that performing optimal spatial filtering on the images further refines the estimates. Importantly, all the

parameters required to perform our inference can be estimated using only the fluorescence data, obviating the need to perform simultaneous electrophysiological experiments. Finally, all the code written to perform the inference is freely available.

1 Introduction

Motivation Simultaneously imaging large populations of neurons using calcium sensors is becoming increasingly popular, both in vitro [1] and in vivo [2, 3, 4], especially as the signal-to-noise-ratio (SNR) of genetic sensors continues to improve [5, 6, 7]. Whereas, the data from these experiments are movies of time-varying fluorescence signals, the desired signal is typically the spike trains or firing rates of the observable neurons. Importantly, to a first approximation, somatic calcium concentration has a relatively simple relationship to spikes. Thus, in theory, one could infer the most likely spike train of each neuron, given the fluorescence data.

Limitations of calcium imaging Unfortunately, finding the *most* likely spike train is a challenging computational task, for a number of reasons. First, the signal-to-noise ratio (SNR) is often low, especially as one increases the image field and frame rate. Second, to find the most likely spike train for a given fluorescence signal, one would have to search over all possible spike trains, a search that would take far too long, pragmatically.

Computational tools as important as experimental tools One is therefore effectively forced to find an approximately most likely spike train, or guess that the inferred spike train is most likely (but not really be sure). The precise details of these approximations, however, are crucial, especially as one approaches shot-noise limited data. For in vitro studies, one often uses 1-photon imaging — either confocal or widefield — in which case the number of photons per neuron is a function of magnification and frame rate (obviously, other parameters, such as number of sensors per neuron, are also important, but typically more difficult to control). To maximize the amount of information one can extract from such preparations, one should increase the frame rate and image field until achieving the shot noise limited regime, assuming one has an inference algorithm that can operate in such a scenario. For in vivo studies, one often uses 2-photon laser scanning microscopy, for which the SNR is relatively reduced, because the dwell time per pixel is (frame duration)/(number of pixels in frame). To circumvent the low SNR for in vivo studies, some groups use long integration times (e.g., [8]), whereas others use small image fields (e.g., [9]). One would prefer to neither sacrifice temporal resolution nor number of observable neurons, to get sufficient signal quality to perform reliable inference. Therefore, it is of the utmost importance to extract as much information as possible from the signal; especially if one is asking quantitative questions about the statistics of the spike trains from the observable neurons — either spontaneously or in relation to some sensorimotor stimulus.

Previous approaches 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, which performed well on their data, but might not be quite as effective as the noise increased, because it relies on a relatively high SNR to not require unreasonable amounts of processing power. 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 carefully consider the statistics of the data. More specifically, their algorithm allows for both negative spikes and negative photon counts, neither of which are possible. In our previous work [12], we developed a sequential Monte Carlo method to efficiently compute the approximate probability of a spike in each image frame, given the entire fluorescence time series. While effective, that approach is not suitable for online analyses, as the computations run in approximately real-time (i.e., analyzing one minute of data requires about one minute of computational time).

Our approach The present work takes the following approach. First, we carefully consider the statistics of typical data sets, and then write down a generative model that accurately relates spiking to observations. Unfortunately, inferring the most likely spike train given this model is computationally intractable. We therefore make some well-justified approximations, which lead to an algorithm that infers the approximately most likely spike train, given the fluorescence data. Our algorithm has a few particularly noteworthy features, relative to other approaches. First, we assume that spikes are always non-negative (i.e., either positive or zero). This is often an important assumption when searching for non-negative signals [13, 14, 15]. Second, our algorithm is extremely fast: it can process a calcium trace from 50,000 images in about one second on a standard laptop computer. Because of these two features, we call our approach the FAst Non-negative deconvolution Spike Inference (FANSI) filter. In addition to these two features, we can generalize our model in a number of ways, to incorporate spatial filtering of the images, overlapping neurons, poisson observations (for use with in vivo data), and slow rise time (for use with genetic sensors).

2 Methods

As described above, to develop an algorithm to approximate the most likely spike train given fluorescence data, we first carefully analyze the statistics of typical data sets. We start by considering an in vitro experiment, for which the SNR is relatively high, and build an appropriate generative model (Section 2.1). Given this model, we can formally state our goal (Section 2.2). And given this goal, we derive an approximately optimal inference algorithm (Section 2.3). We then generalize our model in a number of ways, incorporating spatial filter (Section 3.2), overlapping spatial filters (Section 3.5), Poisson observations for shot-noise limited situations (Section 3.8), and slow rise time for genetic sensors (Section 3.7). Inferring the most likely spike train in all the above scenarios requires having an estimate of the parameters governing the relationship between the spikes and the movie. Thus, we also develop an approach to efficiently approximate the maximum likelihood estimate (MLE) of the parameters (Section 3.3). Finally, we describe several measures we use to assess and compare performance of our algorithm with others (Section 2.4).

2.1 Data driven generative model

Figure 1 shows a typical in vitro, epifluorescence data set (see Section 2.5 for data collection details). The top panel shows a field-of-view, including 3 neurons, two of which are patched. To build our model, we first define a region-of-interest (ROI), which in this case is the circled neuron. Given the ROI, we can average all the pixel intensities of each frame, to get a one-dimensional fluorescence time-series, shown in the bottom left panel (black line). Because we are patched onto this neuron, we also know when this neuron is spiking (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]. We confirmed that model for our data 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 (black line) a histogram of the residuals, and the best fit Gaussian distribution (gray line).

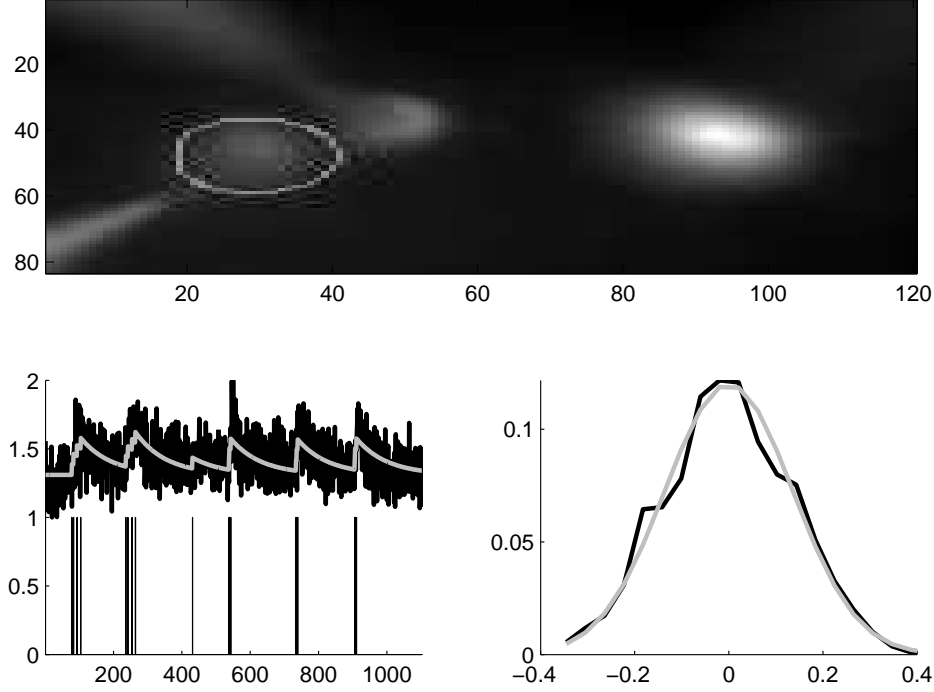


Figure 1: A typical in vitro data set suggests that a reasonable first order model may be constructed by convolving the spike train with an exponential, and adding Gaussian noise. Top panel: the average (over frames) of a typical field-of-view. Bottom left: spike train (black bars), convolved with an 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 we have a 1-dimensional fluorescence trace, $\mathbf{F} = (F_1, \dots, F_T)$ from a neuron. At time t , the fluorescence measurement, F_t is a linear-Gaussian function of the intracellular calcium concentration at that time, C_t :

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

The scale, α , absorbs all experimental variables impacting the scale of the signal, including number of sensors within the cell, photons per change in intracellular calcium concentration per sensor, amplification of imaging system, etc. Similarly, the offset, β , absorbs baseline calcium concentration of the cell, background fluorescence of the fluorophore, imaging system offset, etc. The standard deviation, σ , results from calcium fluctuations independent of spiking activity, fluorescence fluctuations independent of calcium, and imaging noise. The noise at each time, ε_t , is independently and identically distributed according to a standard normal (i.e., Gaussian) distribution. These three parameters, and the Gaussianity of the noise, correspond to a number of simplifying assumptions, that we will relax in Section 3.

We further assume that the intracellular calcium concentration, C_t , jumps after each spike, and subsequently decays back down to rest with time constant, τ , yielding:

$$\tau \frac{C_t - C_{t-1}}{\Delta} = -C_{t-1} + n_t \quad (2)$$

where Δ is the time step size — which in our case, is the frame duration, or $1/(\text{frame rate})$ — and n_t indicates the number of times the neuron spiked at time t . Note that C_t does not refer to absolute intracellular concentration of

calcium but rather, a relative measure. The assumed linearity of our model precludes the possibility of determining calcium in absolute terms (but see Section 4.1 for a modified 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”, we must also define the distribution from which spikes are sampled. 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:

$$n_t \stackrel{iid}{\sim} \text{Poisson}(\lambda\Delta) \quad (3)$$

where $\lambda\Delta$ is the rate, and we have included Δ to make λ be independent of the frame rate. Thus, Eqs. (1) – (3) complete our generative model.

2.2 Goal

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

$$\mathbf{n}_{MAP} = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} P[\mathbf{n}|\mathbf{F}] = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} \frac{1}{P[\mathbf{F}]} P[\mathbf{F}|\mathbf{n}]P[\mathbf{n}] = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} P[\mathbf{F}|\mathbf{n}]P[\mathbf{n}] \quad (4)$$

where n_t is constrained to be an integer ($\mathbb{N}_0 = \{0, 1, 2, \dots\}$), the second equality follows from two applications of Bayes rule, and the third equality follows because $P[\mathbf{F}]$ merely scales the results, but does not change the relative quality of various spike trains. A careful analysis of the above model, Eqs. (1)–(3), suggests that we have a *state-space model* [?]:

$$F_t = \alpha(C_t + \beta) + \sigma\varepsilon_t, \quad \varepsilon_t \stackrel{iid}{\sim} \mathcal{N}(0, 1) \quad (5a)$$

$$C_t = \gamma C_{t-1} + n_t, \quad n_t \stackrel{iid}{\sim} \text{Poiss}(\lambda\Delta) \quad (5b)$$

where $\gamma = 1 - \tau/\Delta$, and n_t has been rescaled appropriately. We use the above formulation to simplify Eq. (4):

$$\mathbf{n}_{MAP} = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} \prod_{t=1}^T P[F_t|C_t]P[C_t|C_{t-1}] \quad (6a)$$

$$= \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} \prod_{t=1}^T P[F_t|C_t]P[n_t] = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmax}} \sum_{t=1}^T (\log P[F_t|C_t] + \log P[n_t]) \quad (6b)$$

where Eq. (4) follows from Eq. (5), and Eq. (6b) follows from Eq. (2). Plugging Eqs. (1) and (3) into Eq. (6b) yields:

$$\mathbf{n}_{MAP} = \underset{\mathbf{n} \in \mathbb{N}_0^T}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 - n_t \log \lambda\Delta + \log n_t! \right), \quad (7)$$

Unfortunately, solving Eq. (7) exactly is computationally intractable, as it requires a nonlinear search over all possible spike trains. We could restrict our search space 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 — more specifically, requires $O(k^T)$ time — which for pragmatic reasons is intractable. Thus, we approximate Eq. (7), by modifying Eq. (2), replacing the Poisson distribution with one that yields a log-concave problem. This reduces the algorithmic complexity from requiring a search over $O(k^T)$ possible spike trains, to some kind of gradient ascent solution.

Two distributions naturally arise as possible approximations to a Poisson: (i) exponential, and (ii) Gaussian. As depicted in Figure 2 an exponential distribution (dashed gray line) is an excellent approximation to the Poisson distribution (solid black line), when λ is small (left panel). On the other hand, when λ is large, a Gaussian distribution

(dash-dotted gray line) closely approximates a Poisson when λ is large (right panel). Thus, naively, it seems as though it may be desirable to use a Gaussian approximation when the neuron is firing with a high firing rate, and approximate the Poisson with an exponential when the neuron is firing sparsely. The Wiener filter is, in fact, the optimal filter given the Gaussian approximation [16] (see [11] for an application of the Wiener filter to this problem). Below, we develop an algorithm to perform inference given the exponential distribution.

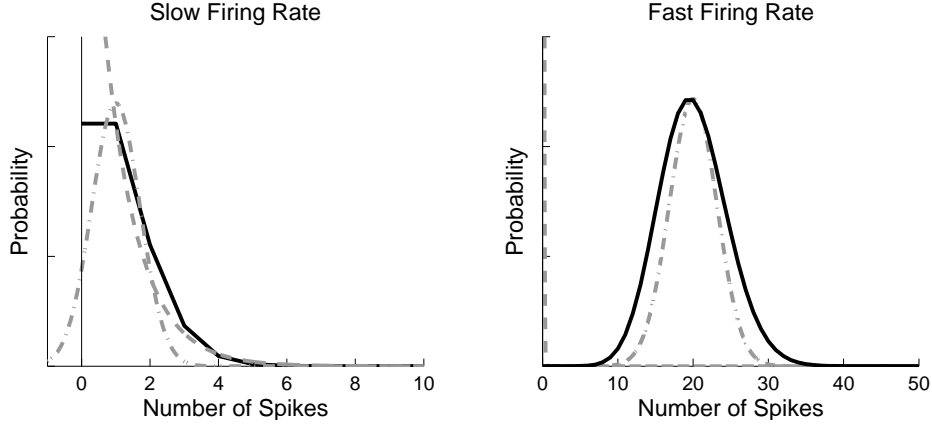


Figure 2: Approximating a Poisson distribution. Left panel: when λ is small (e.g., ≈ 1), an exponential distribution (dashed gray line) approximates the Poisson distribution (solid black line) well, but a Gaussian distribution (dash-dotted gray line) does not. Right panel: when λ is large (e.g., ≈ 20), the inverse is true.

2.3 Inference

Fast Non-negative Deconvolution (FANSI) Filter Our goal here is to develop an algorithm to efficiently approximate $\hat{\mathbf{n}}_{MAP}$. The main contribution of this work demonstrating that we can substitute the Poisson distribution with an exponential distribution, reducing computational complexity from $O(k^T)$ to $O(T)$. Formally, we replace Eq. (2) with:

$$C = \gamma C_{t-1} + n_t, \quad n_t \stackrel{iid}{\sim} \text{Exponential}(\lambda\Delta) \quad (8)$$

where the only difference between Eq. (2) and Eq. (8) is the distribution on n_t . Plugging this change into Eq. (??) yields

$$\mathbf{n}^{FANSI} = \underset{n_t > 0 \forall t}{\operatorname{argmax}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + n_t \lambda \Delta \right) \quad (9)$$

where the constraint on \mathbf{n} has been relaxed from $n_t \in \mathbb{N}_0$ to $n_t \geq 0$ (since the support of an exponential distribution is the non-negative number line). Note that this is a common approximation technique in the machine learning literature [17], as it is the closest convex relaxation to its non-convex counterpart. 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 [18]. We therefore take an “interior-point” (or “barrier”) approach, in which we drop the sharp threshold, and add a barrier term, which must approach $-\infty$ as n_t approaches zero (e.g., $-\log n_t$) [18]. By iteratively reducing the weight of the barrier term, we are guaranteed to converge to the correct solution [18]. Thus, our goal is to efficiently solve:

$$\hat{\mathbf{n}}_z = \underset{n_t \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + n_t \lambda \Delta - z \log(n_t) \right), \quad (10)$$

Since spikes and calcium are related to one another via a simple linear transformation, namely, $n_t = C_t - \gamma C_{t-1}$, we may rewrite Eq. (10) in terms of \mathbf{C} :

$$\hat{\mathbf{C}}_z = \underset{C_t - \gamma C_{t-1} \geq 0 \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + (C_t - \gamma C_{t-1})\lambda\Delta - z \log(C_t - \gamma C_{t-1}) \right). \quad (11)$$

The concavity of Eq. (11) facilitates utilizing any number of techniques guaranteed to find the global optimum. The fact that the argument of Eq. (11) is twice differentiable, suggests that we use the Newton-Raphson technique. Importantly, the state-space nature of this problem yields a particularly efficient approach. Specifically, note that the Hessian is *tridiagonal*, which is clear upon rewriting (11) in matrix notation:

$$\hat{\mathbf{C}}_z = \underset{\mathbf{MC} \geq \mathbf{0}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{F} - \alpha(\mathbf{C} + \beta)\|^2 + (\mathbf{MC})^\top \boldsymbol{\lambda} - z \log(\mathbf{MC})^\top \mathbf{1}, \quad (12)$$

where $\mathbf{M} \in \mathbb{R}^{T \times T}$ is a bidiagonal matrix, $\mathbf{MC} \geq \mathbf{0}$ indicates that every element of \mathbf{MC} is greater than or equal to zero, $^\top$ indicates transpose, $\mathbf{1}$ is a T dimensional column vector, $\boldsymbol{\lambda} = \lambda\Delta\mathbf{1}^\top$, and $\log(\cdot)$ indicates an element-wise logarithm. Note that Eq. (12) follows from writing \mathbf{n} in terms of \mathbf{M} and \mathbf{C} :

$$\mathbf{MC} = \begin{bmatrix} 1 & 0 & 0 & \dots & \dots \\ 1 & -\gamma & 0 & \dots & \dots \\ 0 & 1 & -\gamma & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 & -\gamma \end{bmatrix} \begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ \vdots \\ C_T \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ \vdots \\ n_T \end{bmatrix} = \mathbf{n} \quad (13)$$

Thus, a little bit of calculus yields our update algorithm for \mathbf{C}_z :

$$\hat{\mathbf{C}}_z \leftarrow \hat{\mathbf{C}}_z + s\mathbf{d} \quad (14a)$$

$$\mathbf{H}\mathbf{d} = \mathbf{g} \quad (14b)$$

$$\mathbf{g} = -\frac{\alpha}{\sigma^2} (\mathbf{F} - \alpha(\hat{\mathbf{C}}_z^\top + \beta)) + \mathbf{M}^\top \boldsymbol{\lambda} - z\mathbf{M}^\top (\mathbf{M}\hat{\mathbf{C}}_z)^{-1} \quad (14c)$$

$$\mathbf{H} = \frac{\alpha^2}{\sigma^2} \mathbf{I} + z\mathbf{M}^\top (\mathbf{M}\hat{\mathbf{C}}_z)^{-2} \mathbf{M} \quad (14d)$$

where s is the step size, \mathbf{d} is the step direction, and \mathbf{g} and \mathbf{H} are the gradient (first derivative) and Hessian (second derivative) of the argument in Eq. (12) with respect to \mathbf{C} , respectively, and the exponents indicate element-wise operations. Note that we use “backtracking line searches”, meaning that for each iteration, we find the maximal s that is (i) between 0 and 1 and (ii) decreases the likelihood.

Typically, implementing Newton-Raphson requires inverting the Hessian, i.e., $\mathbf{d} = \mathbf{H}^{-1}\mathbf{g}$, a computation consuming $O(T^3)$ time. Already, this would be a drastic improvement over the most efficient algorithm assuming Poisson spikes, which require $O(k^T)$ time (where k is the maximum number of spikes per frame). Here, because \mathbf{M} is bidiagonal, the Hessian is tridiagonal, the solution may be found in $O(T)$ time via standard banded Gaussian elimination techniques (which can be implemented efficiently in Matlab using $\mathbf{H} \setminus \mathbf{g}$). In other words, the above approximation and inference algorithm reduces computations from exponential time to *linear* time. We refer to this fast algorithm for solving (9) the FASt Nonnegative Deconvolution (FANSI) filter.

Fast Wiener Filter Instead of replacing the Poisson distribution on spikes with an exponential, we can replace it with a Gaussian:

$$C = \gamma C_{t-1} + n_t, \quad n_t \stackrel{iid}{\sim} \mathcal{N}(\lambda\Delta, \lambda\Delta) \quad (15)$$

which, when plugged into Eq. (??) yields

$$\mathbf{n}^{Wiener} = \underset{n_t}{\operatorname{argmax}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - \alpha(C_t + \beta))^2 + \frac{1}{2\lambda\Delta} (n_t - \lambda\Delta)^2 \right) \quad (16)$$

Using the same tridiagonal trick as above, we can solve Eq. (16) using Newton-Raphson once (since we have a quadratic problem here, see Appendix A for details). Because we know only positive spikes are possible, at times we will also consider, $[\text{Wiener}]_+$, which is the Wiener filter half-wave rectified, i.e., all sub-zero values are set to zero.

2.4 Assessment

Let \hat{n} be the inferred spike train, using some algorithm, and n be the true spike train. When spiking is sparse (i.e., at most, one spike per frame), a reasonable measure of inference quality is the effective signal-to-noise ratio, which we define as the squared size of \hat{n} when there is a spike, divided by the squared size of \hat{n} when there is no spike:

$$\text{eSNR}_{\hat{n}} = \frac{\sum_{t|n_t=1} \hat{n}_t}{\sum_{t|n_t=0} \hat{n}_t} \quad (17)$$

When the neuron emits many spikes per frame, eSNR is not meaningful. Instead, we compute the mean squared error between the magnitude of inferred spikes and the true spikes:

$$\text{MSE}_{\hat{n}} = \frac{1}{T} \sum_t (n_t - \hat{n}_t)^2 \quad (18)$$

2.5 Experimental Methods

3 Results

3.1 Main Result

The main result of this paper is that we can approximate \mathbf{n}^{MAP} in *linear* time, whereas an exact solution would require exponential time. Fig. 3 shows two examples of running the FANSI filter on simulated data. On the left, the data is simulated according to Eqs. (1) and (2), with a low expected firing rate (5 Hz). The FANSI filter performs very well (middle left panel): when spikes occurred (gray downward facing triangles), the FANSI filter’s output is relatively high, and in the absence of a spike, the FANSI filter’s output is relatively low. The height of the “spikes” in \mathbf{n}^{FANSI} can be thought of as the probability of a spike having occurred. The Wiener filter does not perform as well on this data. Specifically, \mathbf{n}^{Wiener} exhibits a “ringing” effect, in which the inference oscillates around zero in the absence of true spikes. This occurs because negative spikes improve the inference, if negative spikes are allowed. By constraining our inference to be non-negative (for the FANSI filter), we completely circumvent this problem.

It is unsurprising that the FANSI filter significantly outperforms the Wiener filter when spiking is sparse, given that the exponential is a much better approximation to the Poisson in this regime (c.f. Figure 2). However, even in the fast firing rate scenario, where the Gaussian approximation is far more accurate, the FANSI filter performs relatively well, as depicted in the right panels of Figure 3. Note however that the computational time for computing the FANSI filter scales linearly with T , i.e. is $O(T)$, whereas the naïve implementation of the Wiener filter requires $O(T \log T)$ (but see Appendix ?? for an implementation of the Wiener for that only requires $O(T)$).

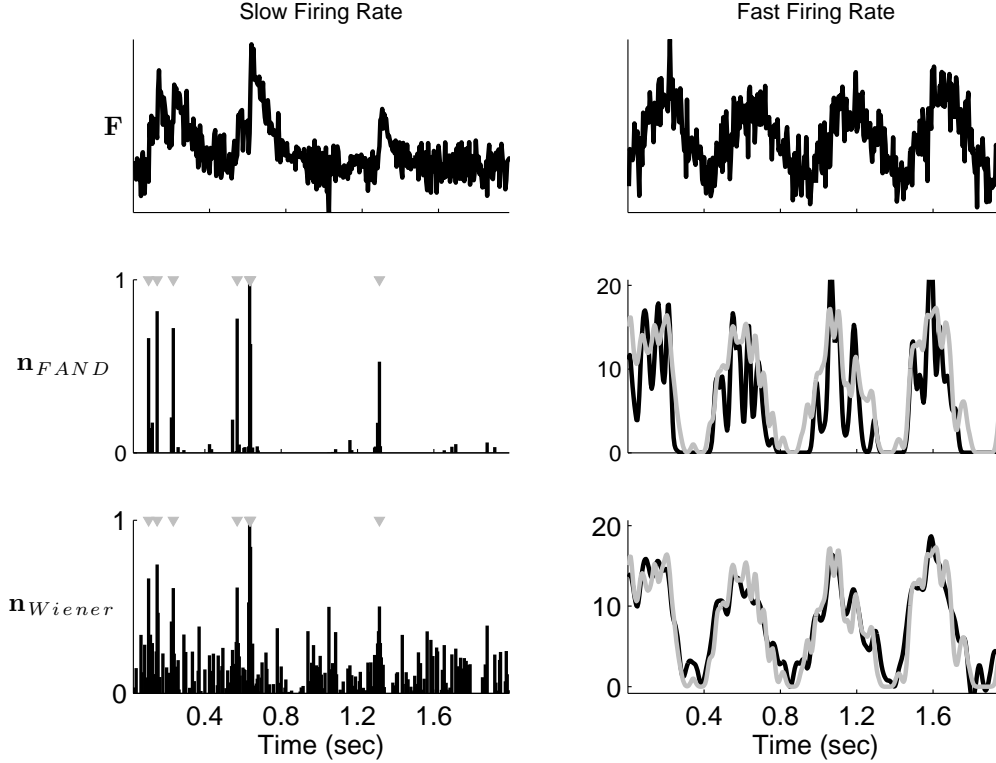


Figure 3: A simulation demonstrating the performance of the FANSI filter in different firing regimes. The left panels show that in the sparse firing regime, the FANSI filter outperforms the Wiener filter in terms of SNR. This follows because an exponential is a closer approximation to a Poisson than a Gaussian, when spiking is sparse. The right panels show that both approximations are good in the fast firing regime. Top left panel: fluorescence time series for a neuron with a slow firing rate. Middle left panel: the FANSI filter's inferred spike train. Bottom left panel: Wiener filter's inferred spike train. Note that (i) the Wiener filter does not impose a non-negativity constraint, and (ii) the effective SNR of the Wiener filter in this example is worse than the FANSI filter's. Top right panel: same as top left panel, for a neuron with a high firing rate. Middle right panel: the FANSI filter's inferred spike train smoothed with a Gaussian kernel for visualization purposes (black line), and the true spike train smoothed with the same Gaussian kernel (gray line). Bottom right panel: same as middle right panel, but with the Wiener filter. Parameters for left panels: same as in Figure ???. Parameters for right panels: same as Figure ??, except: $\sigma = 8$ photons, $\lambda = 500$ Hz.

To quantify the relative quality of these inference algorithms in the sparsely spiking regime, we compute the effective signal-to-noise ratio (eSNR), defined as the average squared magnitude of the inferred spikes during frames in which there was a spike, divided by the same quantity computed in frames lacking a spike. The left panels of Figure 4 show how the eSNR varies as a function of the magnitude of the noise on observations (top left panel), and the expected firing rate (bottom left panel). The FANSI filter's inference (blue line) dominates the Wiener filter (red line), as well as the post-hoc half wave rectifier Wiener filter, $[\text{Wiener}]_+$ (green line), and a simple dF/F (turquoise line).

When the neuron is exhibiting a high firing rate, we compute the mean-squared error (MSE) between the inferred number of spikes, and the true number of spikes. The top right panel of Figure 4 shows that when noise is relatively low, the FANSI filter performs about as well as the Wiener filter. However, in the high noise regime, the Wiener filter clearly outperforms the FANSI filter.

To verify that both our implementations of the FANSI filter and the Wiener filter scale linear with respect to the number of image frames, the bottom right panel of Figure 4 shows such a linear relationship (on a log-log scale).

The take home message from Figure 4 is that when we expect < 1 spike per frame, the FANSI filter significantly outperforms the Wiener filter, regardless of variance of the observation noise. Furthermore, because of our linear time algorithm, filtering around 50,000 image frames requires only about 1 second on a standard Apple laptop. Below, we improve on our inference quality by generalizing our model in a number of ways.

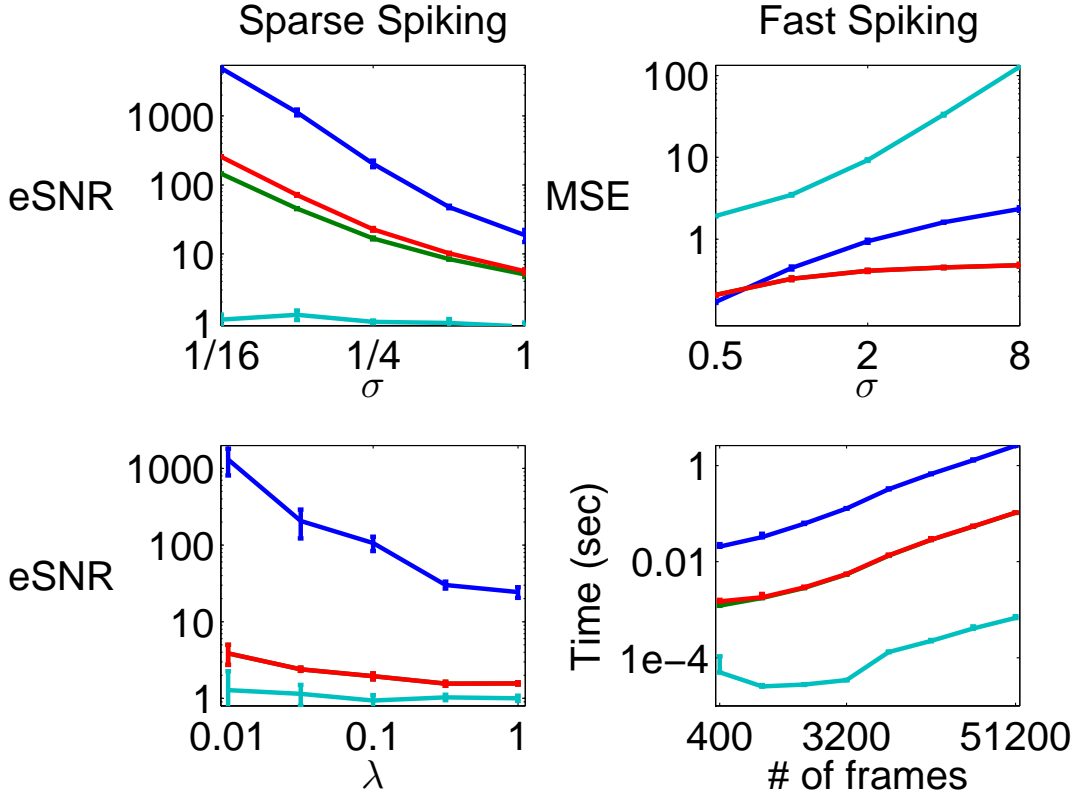


Figure 4: Quantitative assessment of the FANSI filter’s inference quality and speed. The left columns show that the FANSI filter outperforms the Wiener filter in terms of eSNR — Eq. (17) — regardless of variance of observation noise (top) and expected firing rate (bottom). The top right panel shows that when noise is low, even when firing rates are relatively (e.g. $\lambda\Delta = 20$ Hz), the FANSI filter and Wiener filter both perform well. Both algorithms scale linearly with the number of image frames, meaning that it takes about 1 second of computational time per 50,000 image frames, using either algorithm. Simulation details: 5 simulations for each point on all the plots, mean (solid line) and standard deviation (bars) are shown for each. Parameters: $\Delta = 0.005$ sec, $\alpha = 1$, $\beta = 0$. Top left: $\lambda = 1$ Hz, $\tau = 1$ sec. Top right: $\lambda = 10$ Hz, $\tau = 1$ sec. Bottom left: $\sigma = 0.25$, $\tau = 0.5$ sec. Bottom right: $\sigma = 0.25$, $\tau = 0.1$ sec, $\lambda = 1$ Hz.

Finally, often one is interested in understanding the relationship between spike trains and the environment. Therefore, we simulated a neuron whose spiking activity was a function of a 5-dimensional external stimulus. More specifically, we let $\lambda_t = \mathbf{k}^T \mathbf{x}_t$, where \mathbf{k} is a 5-dimensional linear kernel, \mathbf{x}_t is the stimulus at time t , and $P(n_t) = \text{Poisson}(\lambda_t \Delta)$. We then computed the maximum likelihood estimate of the linear kernel, $\hat{\mathbf{k}} = \hat{\mathbf{n}} \mathbf{x}^T (\mathbf{x} \mathbf{x}^T)^{-1}$. Importantly, the simulation was constructed to incorporate both sparse spiking and fast spiking periods, much like sensory neurons have periods of quiescence, followed by stimulus driven bursts. Thus, neither the FANSI filter nor the Wiener filter’s assumptions are entirely appropriate. Figure 5 shows the results of this simulation. The true kernel, kernel estimated using the true spike times, and kernel estimated using the FANSI filter, are nearly overlapping (black, gray, and blue lines, respectively). The kernel estimated using the Wiener filter output, however, is effectively flat (red line). Post-hoc half-wave rectification of the Wiener filter output does not improve its ability to estimate this kernel (green line — completely obscured by the unrectified Wiener filter output). Finally, dF/F does not yield anything useful at all (turquoise line). This simulation provides further evidence of the quantitative advantage of utilizing the FANSI filter before performing an analysis on calcium fluorescence data.

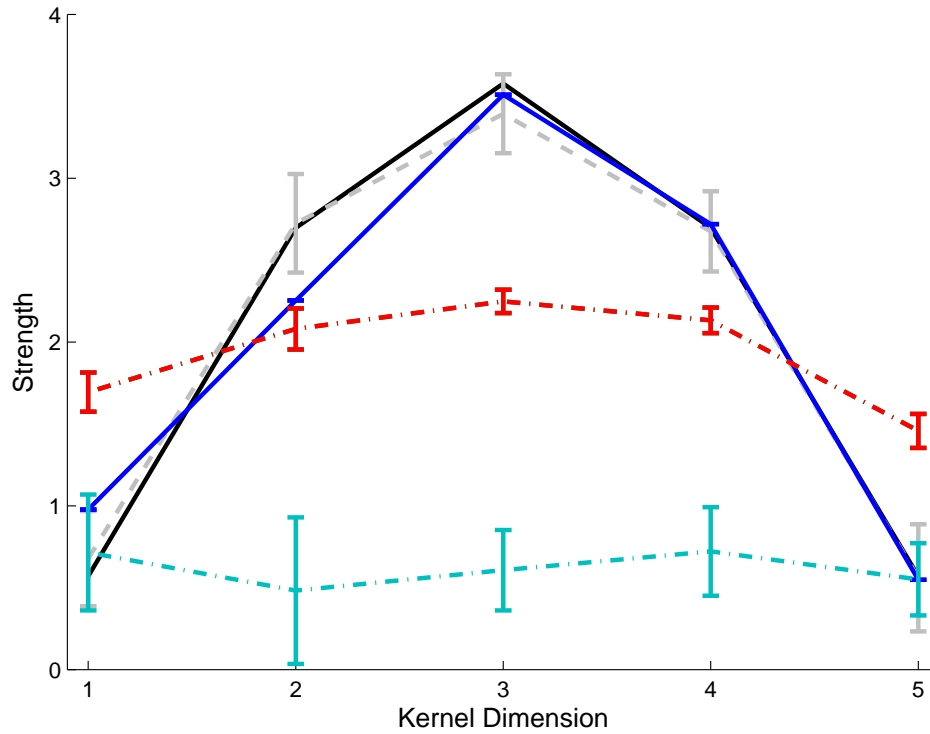


Figure 5: The tuning curve of a neuron (black line), estimating from the true spike train (gray line), and the inferred spike train using the (i) FANSI filter (blue line), (ii) Wiener filter (red line), and (iii) dF/F. Clearly, the FANSI filter performs nearly as well as the true spike train, whereas the Wiener filter and dF/F do not. Note that half-wave rectification of the Wiener did not change the results at all. Simulation details: mean (solid lines) and standard deviation (bars) of 5 simulations, $T = 800$, $\Delta = 0.005$ sec, $\alpha = 1$, $\beta = 0$, $\tau = 0.1$ sec, $x_{i,t} \sim \mathcal{U}(0, 0.2)$, $\sigma = 0.25$.

3.2 Spatial Filtering

In the previous sections, we implicitly assumed that the raw movie of fluorescence measurements collected by the experimenter had undergone two stages of pre-processing. First, the movie was segmented, to determine regions-of-interest (ROIs). This yields a vector, \vec{F}_t , corresponding to the fluorescence intensity at time t for each of the N_p pixels in the ROI. Second, we projected that vector into a scalar, yielding F_t , the assumed input. In this section, we still assume that somebody has gone through our movies and performed some segmentation, but we do not assume that they have projected the vector \vec{F}_t into a scalar F_t^{proj} . Formally, we posit a more general model:

$$\vec{F}_t = \vec{\alpha}(C_t + \beta) + \sigma\vec{\varepsilon}_t, \quad \vec{\varepsilon}_t \sim \mathcal{N}(\vec{0}, \mathbf{I}) \quad (19)$$

where \vec{F}_t , $\vec{\alpha}$, $\vec{\varepsilon}_t$, and $\vec{0}$ are all column vectors of length N_p , and \mathbf{I} is an $N_p \times N_p$ identity matrix. This model follows because the fluorescence at any individual pixel may be composed of a static element, β , and a dynamic element, that we assume is purely due to calcium fluctuations, C_t . Further, we have assumed that the noise is uncorrelated and has the same variance, σ^2 , in each pixel (an assumption that can be relaxed quite easily). Performing inference in this more general model proceeds nearly identical as before:

$$\hat{C}_z = \underset{MC \geq 0}{\operatorname{argmin}} \frac{1}{2\sigma^2} \left\| \vec{F} - \vec{\alpha}(C^\top + \beta \mathbf{1}^\top) \right\|^2 + (MC)^\top \lambda - z \log(MC)^\top \mathbf{1}, \quad (20)$$

$$\mathbf{g} = -\frac{\alpha}{\sigma^2} (\mathbf{F} - \alpha(\hat{C}_z^\top + \beta)) + M^\top \lambda - z M^\top (M \hat{C}_z)^{-1} \quad (21)$$

$$\mathbf{H} = \frac{\alpha^\top \alpha}{\sigma^2} \mathbf{I} + z M^\top (M \hat{C}_z)^{-2} M \quad (22)$$

Figure 6 demonstrates the utility of this generalization. The top row shows different depictions of an ROI containing a single neuron. On the far left panel is the “true” spatial filter. We modeled the true spatial filter as a sum of Gaussians: a positively weighted small variance Gaussian, and a negatively weighted large variance Gaussian. We chose this model based on our empirical observations that often pixels immediately around a neuron exhibit calcium sensitive fluctuations that are anti-correlated with the pixels on the neuron, probably due to the calcium influx from the vicinity of the neuron. The mean frame (second panel from left) looks very similar to the true spatial filter, as individual frames are effectively just modulating the magnitude of the spatial filter, and adding noise. Typically, to obtain an ROI, one would simply identify the pixels with high positive values from the mean frame, and average them together (third panel from left). Such an approach yields the 1-dimensional fluorescence projection depicted on the left, the typical projection, F_t^{typ} , with its associated inferred spike train beneath. Using the true spatial filter to project \vec{F} onto a 1-dimensional fluorescence time series results on the right, the optimal projection, F_t^{opt} , with its associated inferred spike train beneath. It should be clear that using the true spatial filter improves the SNR of the fluorescence signal, and therefore, the inferred spike train accuracy.

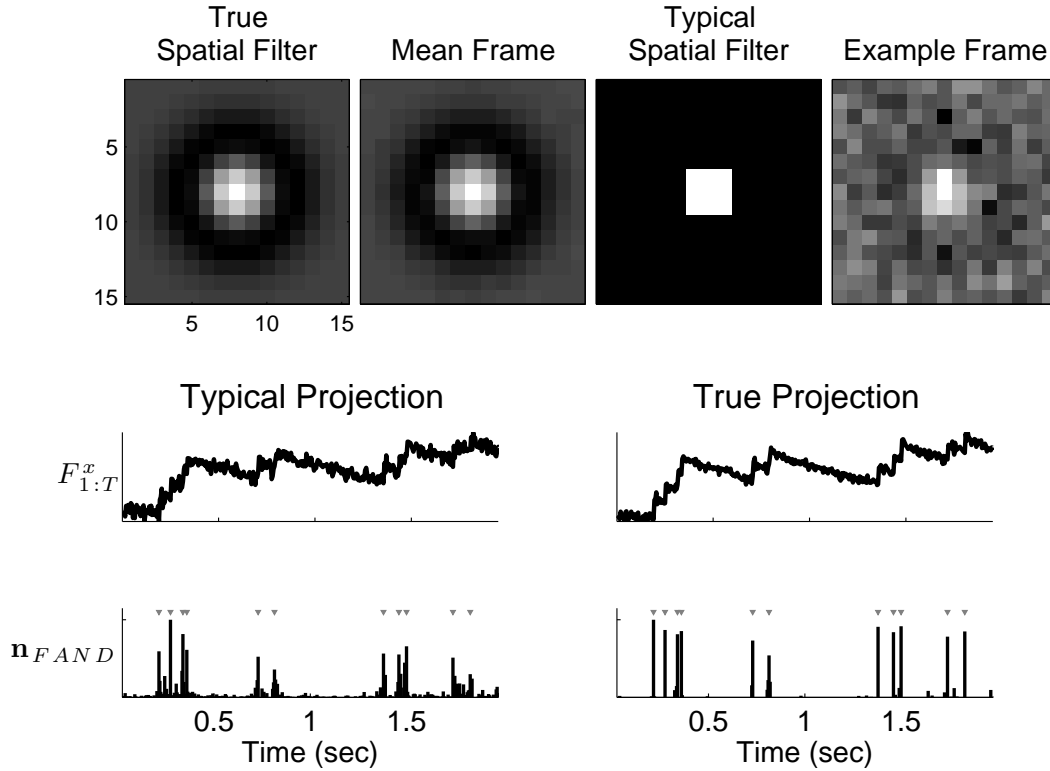


Figure 6: A simulation demonstrating that using a better spatial filter can significantly enhance the effective SNR (see Supplementary Movie 1 for the full movie associated with this simulation). Top left: true spatial filter. Top second from left: mean frame. Top second from right: typical spatial filter. Top right: example frame from movie (frame number 100). Middle left: 1-dimensional fluorescence projection using typical spatial filter. Bottom left: \mathbf{n}^{FANDSI} using typical spatial filter. Middle right: 1-dimensional fluorescence projection using true spatial filter. Bottom right: \mathbf{n}^{FANDSI} using true spatial filter. Simulation details: $\alpha = \mathcal{N}(\mathbf{0}, 2\mathbf{I}) - 1.1\mathcal{N}(\mathbf{0}, 2.5\mathbf{I})$ where $\mathcal{N}(\mu, \Sigma)$ indicates a Gaussian with mean μ and covariance matrix Σ , $\beta = 1$, $\tau = 0.85$ sec, $\lambda = 5$ Hz.

3.3 Learning

In the above, we assumed that the parameters governing our model, $\theta = \{\alpha, \beta, \sigma, \gamma, \lambda\} \in \Theta$, were known. In general, however, these parameters must be estimated from the data. Therefore, we take an approximate expectation-maximization approach for computing n^{FANSI} : initialize some estimate of the parameters, $\hat{\theta}$, recursively compute n^{FANSI} using those parameters and update $\hat{\theta}$ given n^{FANSI} , stop recursing when some convergence criteria is met. Below, we provide details for each of the above steps.

Initializing the parameters We initialize our estimate α using Principal Component Analysis (PCA). More specifically, we perform PCA on F (the whole movie), and let α be the PC with the highest eigenvalue. If the movie is very long, we either only use part of it, or simply take the mean frame. Then, we let $\beta = 0$. σ is estimated by finding a sequence of the time series lacking any obvious spikes, and computing the root mean square of that segment (i.e., $\hat{\sigma} = \sqrt{\vec{F}_{s:t}^2 / (t - s)}$). We typically set γ and λ based on our previous experience with these cells. For instance, $\gamma \approx 0.95$ is often reasonable, and λ is somewhere between 1 and 10 Hz.

Estimating the parameters given \hat{n} To find the maximum likelihood estimator for the parameters, $\hat{\theta}$, we must integrate over all possible spike trains. Unfortunately, it is not currently known how to perform this integral exactly, and approximating this integral using Monte Carlo methods is relatively time consuming (see [12] for details). Thus, we resort to a more drastic approximation, commonly used in state-space models. For specifically, instead of integrating over all possible spike trains, we only consider the most likely sequence (often referred to as the Viterbi path [19]):

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \int P[F|C, \theta] P[C|\theta] dC \approx \operatorname{argmax}_{\theta \in \Theta} P[F|\hat{C}, \theta] P[\hat{C}|\theta] \quad (23)$$

where \hat{C} is determined using the above described inference algorithm. The approximation in (23) is good whenever the likelihood is very peaky, meaning that most of the mass is around the MAP sequence.¹

Due to the state space nature of the above model (Eqs (19) and (2)), the optimization in Eq (23) simplifies significantly. More specifically, we can write the argument from the right-hand-side of Eq (23) as a product of terms that we have defined in our model:

$$P[F|\hat{C}, \theta] P[\hat{C}|\theta] = \prod_{t=1}^T P[F_t|\hat{C}_t; \alpha, \beta, \sigma] P[\hat{C}_t|\hat{C}_{t-1}, \hat{n}_t; \gamma] P[\hat{n}_t|\lambda]. \quad (24)$$

This optimization simplifies into several separable problems. First, we solve for α and β . Because solving for them jointly is non-concave, we solve for each separately. Consider α only:

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \prod_{t,x=1}^{T,N_p} P_{\theta}[F_{t,x}|C_t] = \operatorname{argmax}_{\alpha} \prod_{t,x=1}^{T,N_p} \mathcal{N}(F_{t,x}; \alpha_x(C_t + \beta), \sigma^2) \quad (25a)$$

$$= \operatorname{argmax}_{\alpha} \sum_{t,x=1}^{T,N_p} \log \mathcal{N}(F_{t,x}; \alpha_x(C_t + \beta), \sigma^2) \quad (25b)$$

$$= \operatorname{argmax}_{\alpha} -\frac{N_p T}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t,x=1}^{T,N} (F_{t,x} - \alpha_x(C_t + \beta))^2 \quad (25c)$$

$$= \operatorname{argmin}_{\alpha} \sum_{t,x=1}^{T,N_p} (F_{t,x} - \alpha_x(C_t + \beta))^2. \quad (25d)$$

Therefore, we can solve for each of the $\hat{\alpha}_x$'s separately and efficiently using Matlab's `mldivide`: $\hat{\alpha}_x = (C + \beta \mathbf{1}) \backslash F_x$, where $F_x = [F_{1,x}, \dots, F_{T,x}]^T$. Given $\hat{\alpha}$, we can estimate β :

¹The approximation in (23) may be considered a first-order Laplace approximation

$$\hat{\beta} = \operatorname{argmax}_{\beta > 0} \sum_{t,x=1}^{T,N} (F_{t,x} - \hat{\alpha}_x(C_t + \beta))^2 = \operatorname{argmax}_{\beta > 0} \sum_{t,x=1}^{T,N} (F_{t,x} - \hat{\alpha}_x C_t + \beta \hat{\alpha}_x)^2 \quad (26)$$

which can again be solved efficiently again using Matlab's `mldivide`: $\tilde{\beta} = \hat{\alpha}_x \backslash \left(\sum_{t=1}^T F_{t,x} - \hat{\alpha}_x C_t \right)$. If $\tilde{\beta} < 0$, we simply let $\hat{\beta} = 0$, else, $\hat{\beta} = \tilde{\beta}$.

Given $\hat{\alpha}$ and $\hat{\beta}$, we can now estimate σ using the residuals, i.e.,

$$\hat{\sigma}^2 = \frac{1}{TN_p} \left\| \vec{F} - \hat{\alpha}(C^T + \beta \mathbf{1}) \right\|^2 \quad (27)$$

Estimating λ is also totally straightforward: $\hat{\lambda} = \hat{n}' \mathbf{1} / (T\Delta)$.

We don't update γ as we have found that it does not improve inference quality (not shown), in agreement with previous work [20].

We stop iterating the parameter update whenever (i) iteration number exceeds some upper bound, or (ii) relative change in likelihood does not exceed some lower bound. In practice, we find that parameters tend to converge after several iterations, given our initialization. Figure 7 shows a simulation demonstrating the efficacy of this procedure. The left panel shows the true spatial filter (top), the 1-dimensional fluorescence projection using the true spatial filter, F^{opt} (middle), and the inferred spike train using the true parameters (bottom). The right panel shows the estimated spatial filter, the 1-dimensional fluorescence projection using the estimated filter, F^{est} , and the inferred spike train using the estimated filter. Note that we seeded the algorithm with the movie, and parameter estimates that were at least a factor of 2 from the true values.

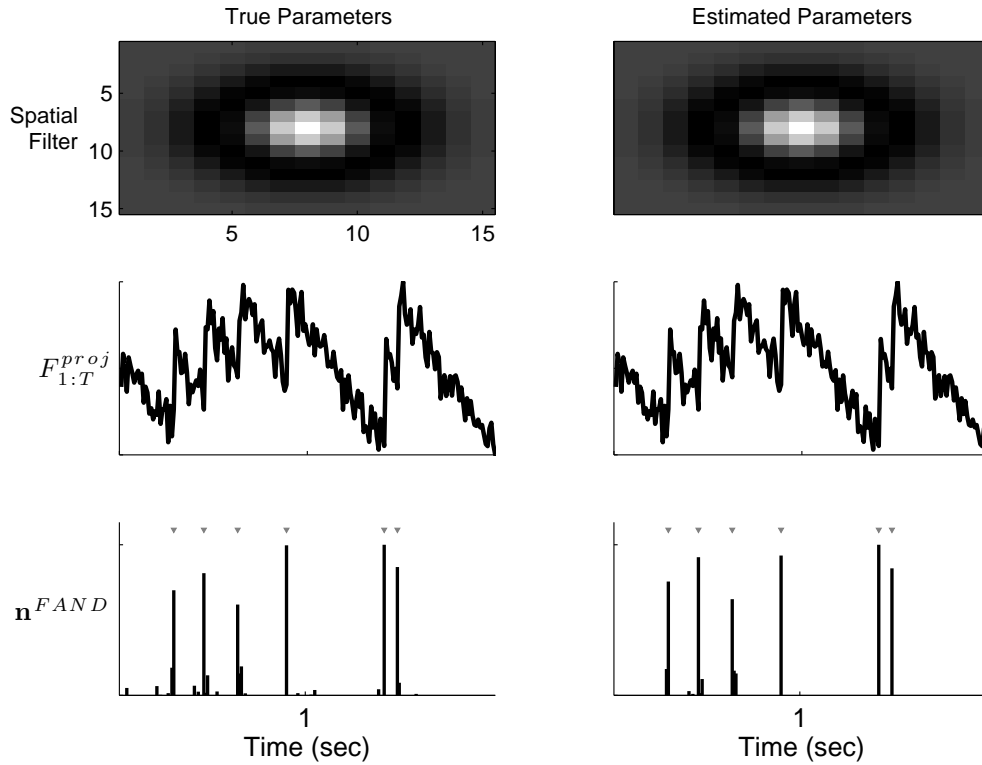


Figure 7: A simulation demonstrating that given only the fluorescence movie, the parameters may be estimated, and the spike train inferred (c.f. Supplementary Movie 2). Top left panel: true spatial filter. Middle left panel: projection of movie onto true spatial filter. Bottom left panel: inferred spike train using true parameters. Right panels: same as left except estimating parameters. All parameters estimated other than γ , which was assumed known. Parameters converged within 7 iterations. Simulation details: $T = 1000$, $\Delta = 5$ msec, α is the same as in Figure 6, $\beta = 0$, $\tau = 500$ msec, $\lambda = 10$ Hz.

3.4 in vitro data

3.5 Overlapping spatial filters

Model

$$\mathbf{F}_t = \sum_{i=1}^{N_c} \boldsymbol{\alpha}_i (C_{i,t} + \beta_i) + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (28)$$

$$C_{i,t} = \gamma_i C_{i,t-1} + n_{i,t}, \quad n_{i,t} \sim \text{Poisson}(n_{i,t}; \lambda_i \Delta) \quad (29)$$

implicit assumption that $\mathbf{n}_i \perp \mathbf{n}_j, \forall i \neq j$

Inference let $\mathbf{n} = (n_{1,1}, n_{2,1}, \dots, n_{N_c,1}, n_{1,2}, \dots, n_{N_c,T})^\top$. similar def of \mathbf{C} .

$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 0 & 0 & \dots & \dots \\ 1 & -\gamma_1 & 1 & -\gamma_2 & \dots & 1 & -\gamma_{N_c} & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & & & & \\ 0 & 0 & 0 \dots & 1 & -\gamma_{N_c-1} & 1 & -\gamma_{N_c} & & \end{bmatrix} \quad (30)$$

inference as before, but replacing the scalar β with $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{N_c})^\top$, and making minor adjustments to deal with dimensionality issues.

Learning estimating $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{N_c})^\top$ is similar. now, we compute $\hat{\boldsymbol{\alpha}}_x = (\hat{\alpha}_{1,x}, \dots, \hat{\alpha}_{N_c,x})^\top$ using

$$\hat{\boldsymbol{\alpha}}_x = (\mathbf{C} + \tilde{\boldsymbol{\beta}}) \setminus \mathbf{F}_x, \quad (31)$$

where $\tilde{\boldsymbol{\beta}}$ is $\boldsymbol{\beta}$ reparameterized to be the same size as \mathbf{C} .

estimating $\boldsymbol{\beta}$ proceeds as before, but since it is a vector, we use Matlab's `quadprog`, imposing the constraint that $\beta_i > 0 \forall i$.

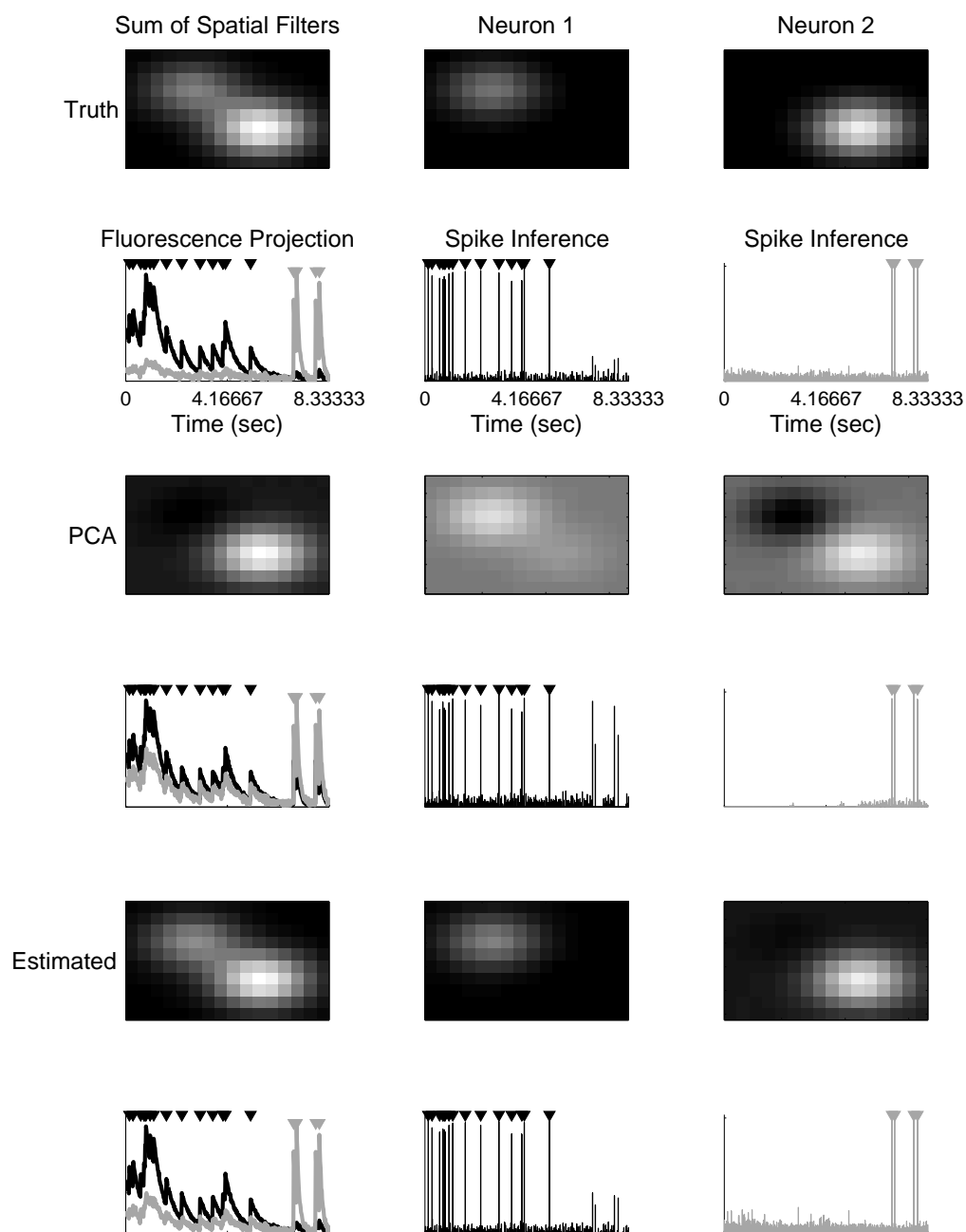


Figure 8: Simulation showing that even when two neuron's spatial filters are largely overlapping

3.6 Population imaging

Figure 9: full movie, in vitro data

3.7 Slow rise time

Figure 10: Slow rise time

3.8 Poisson observation model

Model

$$\mathbf{F}_{x,t} \sim \text{Poisson}(\alpha_x(C_t + \beta)) \quad (32)$$

Inference

$$\mathcal{L}_{x,t} = -\alpha_x(C_t + \beta) + F_{x,t} \log(\alpha_x(C_t + \beta)) - \log(F_{x,t}!) \quad (33a)$$

$$g_{x,t} = -\alpha_x + F_{x,t}(C_t + \beta)^{-1} \quad (33b)$$

$$H_{x,t} = -F_{x,t}(C_t + \beta)^{-2} \quad (33c)$$

where $\mathcal{L} = \sum_{x,t} \mathcal{L}_{x,t}$, $\mathbf{g} = \sum_x (g_{x,1}, \dots, g_{x,T})^\top$, and $\mathbf{H} = \text{diag}(\sum_x H_{x,t})$.

Figure 11: Poisson

3.9 in vivo data

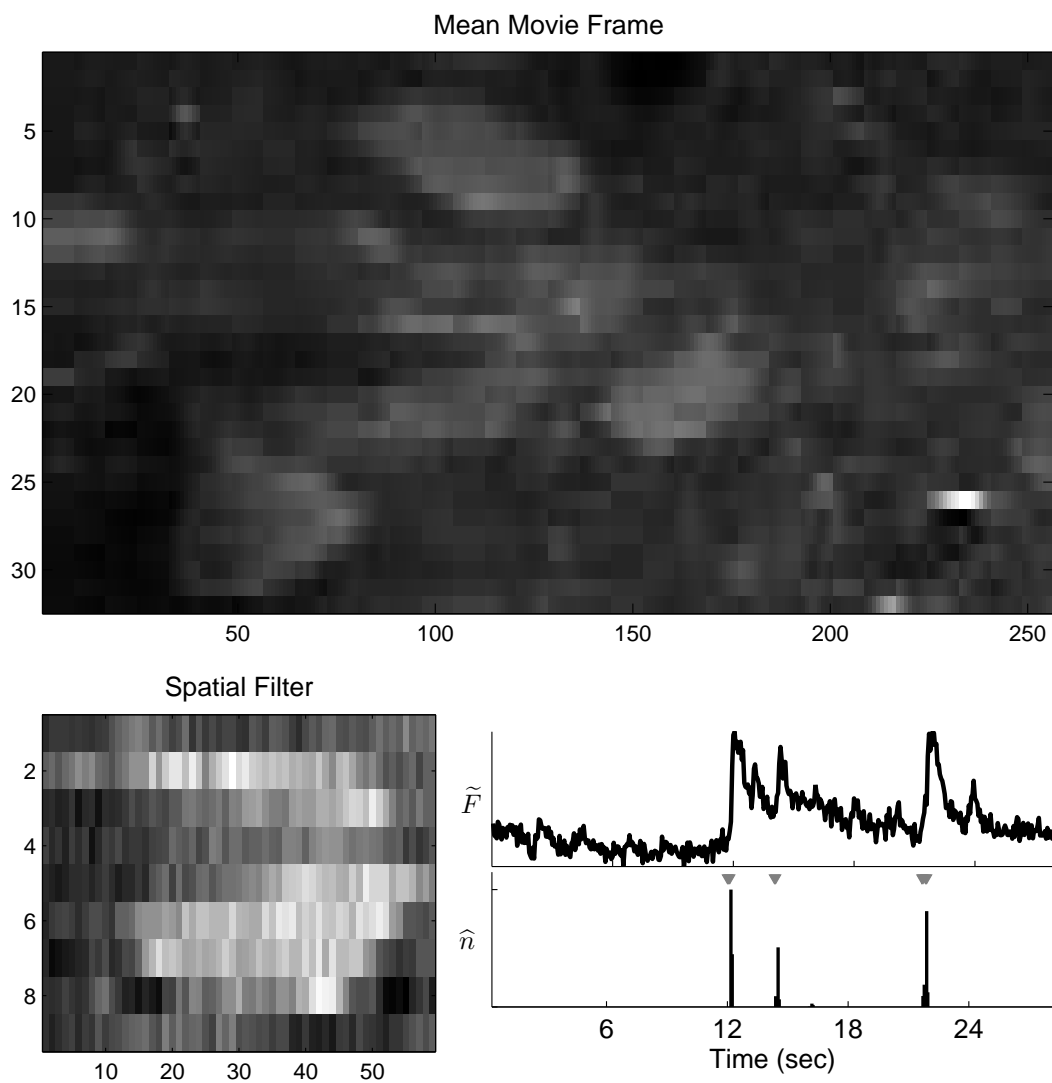


Figure 12: Given only a fluorescence movie, recorded in vivo, we can learn the parameters necessary to correctly infer the spike trains. Left: mean frame. Left: projection of movie onto mean frame. Left: the FANSI filter's inference.

4 Discussion

Summary

Extensions

Thresholding

4.1 Incorporating a nonlinear observation model

$$\mathbf{F}_t = \boldsymbol{\alpha} S(C_t + \beta) + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (34)$$

where $S(x) = \frac{x^{n_d}}{x^{n_d} + k_d}$

note: initialize with linear result, but add a constant wherever constraint is not satisfied

Figure 13: Saturation

4.2 Dynamic prior

Model let $\lambda = (\lambda_1, \dots, \lambda_T)^\top$

$$C_t = \gamma C_{t-1} + n_t, \quad n_t \sim \text{Poisson}(n_t; \lambda_t \Delta) \quad (35)$$

Inference

$$\mathcal{L} = \frac{1}{2\sigma^2} \left\| \mathbf{F} - \alpha(\mathbf{C}^\top + \beta + \mathbf{1}^\top) \right\|^2 + (\mathbf{MC})^\top \lambda \Delta - z \log(\mathbf{MC})^\top \mathbf{1} \quad (36a)$$

$$\mathbf{g} = -\frac{\alpha}{\sigma^2} (\mathbf{F} - \alpha(\hat{\mathbf{C}}_z^\top + \beta)) + \mathbf{M}^\top \lambda \Delta - z \mathbf{M}^\top (\mathbf{M} \hat{\mathbf{C}}_z)^{-1} \quad (36b)$$

Acknowledgments 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. BOW was supported by NDS grant F30 NS051964. The authors would like to thank A. Packer for helpful discussions.

References

- [1] R. Yuste and A. Konnerth. *Imaging in Neuroscience and Development, A Laboratory Manual*, 2006.
- [2] 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 Ca^{2+} imaging of local neuronal circuits. *Neuron*, 53(6):789–803, Mar 2007.
- [3] Werner Göbel and Fritjof Helmchen. In vivo calcium imaging of neural network function. *Physiology (Bethesda)*, 22:358–365, Dec 2007.
- [4] Liqun Luo, Edward M Callaway, and Karel Svoboda. Genetic dissection of neural circuits. *Neuron*, 57(5):634–660, Mar 2008.
- [5] 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.
- [6] Marco Mank and Oliver Griesbeck. Genetically encoded calcium indicators. *Chem Rev*, 108(5):1550–1564, May 2008.
- [7] Damian J Wallace, Stephan Meyer zum Alten Borgloh, Simone Astori, Ying Yang, Melanie Bausen, Sebastian Kglar, Amy E Palmer, Roger Y Tsien, Rolf Sprengel, Jason N D Kerr, Winfried Denk, and Mazahir T Hasan. Single-spike detection in vitro and in vivo with a genetic Ca^{2+} sensor. *Nat Methods*, 5(9):797–804, Sep 2008.
- [8] Kenichi Ohki, Sooyoung Chung, Prakash Kara, Mark Hübner, Tobias Bonhoeffer, and R. Clay Reid. Highly ordered arrangement of single neurons in orientation pinwheels. *Nature*, 442(7105):925–928, Aug 2006.
- [9] J.N.D. Kerr, C.P.J. de Kock, D.S. Greenberg, R.M. Bruno, B. Sakmann, and F. Helmchen. Spatial organization of neuronal population responses in layer 2/3 of rat barrel cortex. *Journal of Neuroscience*, 27(48):13316, 2007.
- [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*, 11(7):749 – 751, Jun 2008.
- [11] Terrence F Holekamp, Diwakar Turaga, and Timothy E Holy. Fast three-dimensional fluorescence imaging of activity in neural populations by objective-coupled planar illumination microscopy. *Neuron*, 57(5):661–672, Mar 2008.
- [12] JT Vogelstein, BO Watson, Packer AM, R Yuste, Jedynak B, and L Paninski. Spike inference from calcium imaging using sequential monte carlo methods. *Biophysical Journal*, 2009.

- [13] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, Oct 1999.
- [14] D.D. Lee and H.S. Seung. Algorithms for non-negative matrix factorization. *Advances in Neural Information Processing Systems*, 13:556–562, 2001.
- [15] 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.
- [16] N. Wiener. *Extrapolation, interpolation, and smoothing of stationary time series: with engineering applications*. MIT Press, Cambridge, Mass., 1949.
- [17] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2001.
- [18] S. Boyd and L. Vandenberghe. *Convex Optimization*. Oxford University Press, 2004.
- [19] Lawrence R Rabiner. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, 72(2):257–286, February 1989.
- [20] Emre Yaksi, Benjamin Judkewitz, and Rainer W Friedrich. Topological reorganization of odor representations in the olfactory bulb. *PLoS Biol*, 5(7):e178, Jul 2007.

A Wiener Filter

Sections 2.3 outline one approach to solving Eq. (6), by approximating the Poisson distribution with an exponential distribution, and imposing a non-negative constraint on the inferred $\hat{\mathbf{n}}$. Perhaps a more straightforward approach would be to approximate the Poisson distribution with a Gaussian distribution. In fact, as rate increases above about 10 spikes/sec, a Poisson distribution with rate $\lambda\Delta$ is well approximated by a Gaussian with mean and variance $\lambda\Delta$. Given such an approximation, instead of Eq. (9), we would obtain:

$$\hat{\mathbf{n}}_w \approx \underset{\mathbf{n}_t \in \mathbb{R}, \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 + \frac{1}{2\lambda\Delta} (n_t - \lambda\Delta)^2 \right) \quad (37)$$

As above, we can rewrite Eq. (37) in matrix notation in terms of \mathbf{C} :

$$\hat{\mathbf{C}}_w = \underset{\mathbf{C}_t \in \mathbb{R}, \forall t}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{F} - \mathbf{C}\|^2 + \frac{1}{2\lambda\Delta} \|\mathbf{MC} - \lambda\Delta\mathbf{1}\|^2 \quad (38)$$

which is quadratic in \mathbf{C} , and may therefore be solved analytically using quadratic programming, $\hat{\mathbf{C}}_w = \hat{\mathbf{C}}_0 + \mathbf{d}_w$, where $\hat{\mathbf{C}}_0$ is the initial guess and $\mathbf{d}_w = \mathbf{H}_w \backslash \mathbf{g}_w$, where

$$\mathbf{g}_w = \frac{1}{\sigma^2} (\mathbf{C}'_0 - \mathbf{F}) + \frac{1}{\lambda\Delta} ((\mathbf{M}\hat{\mathbf{C}}_0)' \mathbf{M} - \lambda\Delta \mathbf{M}' \mathbf{1}) \quad (39)$$

$$\mathbf{H}_w = \frac{1}{\sigma^2} \mathbf{I} + \frac{1}{\lambda\Delta} \mathbf{M}' \mathbf{M} \quad (40)$$

Note that this solution is the optimal linear solution, under the assumption that spikes follow a Gaussian distribution, and is often referred to as the Wiener filter, regression with a smoothing prior, or ridge regression. To estimate the parameters for the Wiener filter, we take the same approach as above:

$$\hat{\boldsymbol{\theta}}_w \approx \underset{\boldsymbol{\theta}_w}{\operatorname{argmax}} P[\mathbf{F} | \hat{\mathbf{n}}_w, \boldsymbol{\theta}_w] P[\hat{\mathbf{n}}_w | \boldsymbol{\theta}_w] \quad (41a)$$

$$= \underset{\boldsymbol{\theta}_w}{\operatorname{argmax}} -\frac{T}{2} \log(4\pi^2 \sigma^2 \lambda\Delta) - \frac{1}{2\sigma^2} \|\mathbf{Y}_w + \boldsymbol{\eta}_w \mathbf{X}_w\|^2 - \frac{1}{2\lambda\Delta} \|\hat{\mathbf{n}}_w - \lambda\Delta \mathbf{1}\|^2 \quad (41b)$$

where \mathbf{Y}_w , $\boldsymbol{\eta}_w$, and \mathbf{X}_w are defined as their subscriptless counterparts in Eq. (23).