



HAL
open science

Blind deconvolution for spike inference from fluorescence recordings

Jérôme Tubiana, Sébastien Wolf, Thomas Panier, Georges Debregeas

► **To cite this version:**

Jérôme Tubiana, Sébastien Wolf, Thomas Panier, Georges Debregeas. Blind deconvolution for spike inference from fluorescence recordings. *Journal of Neuroscience Methods*, 2020, 342, pp.108763. 10.1016/j.jneumeth.2020.108763 . hal-03064793

HAL Id: hal-03064793

<https://hal.science/hal-03064793v1>

Submitted on 15 Jul 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Blind deconvolution for spike inference from fluorescence recordings

Jérôme Tubiana,¹ Sébastien Wolf,² Thomas Panier,³ and Georges Debregeas³

¹*Blavatnik School of Computer Science, Tel Aviv University, Israel*

²*Université Paris Sciences et Lettres, France, Laboratoire de Physique de l'Ecole Normale Supérieure, CNRS UMR 8023 & PSL Research, France, Institut de Biologie de l'Ecole Normale Supérieure, CNRS, INSERM, UMR 8197 & PSL Research, France*

³*Sorbonne Université, CNRS, Institut de Biologie Paris-Seine (IBPS), Laboratoire Jean Perrin (LJP), France*

(Dated: May 5, 2020)

The parallel developments of genetically-encoded calcium indicators and fast fluorescence imaging techniques allows one to simultaneously record neural activity of extended neuronal populations *in vivo*. To fully harness the potential of functional imaging, one needs to infer the sequence of action potentials from fluorescence traces. Here we build on recently proposed computational approaches to develop a blind sparse deconvolution (BSD) algorithm based on a generative model for inferring spike trains from fluorescence traces. BSD features, (1) automatic (fully unsupervised) estimation of the hyperparameters, such as spike amplitude, noise level and rise and decay time constants, (2) a novel analytical estimate of the sparsity prior, which yields enhanced robustness and computational speed with respect to existing methods, (3) automatic thresholding for binarizing spikes that maximizes the precision-recall performance, (4) super-resolution capabilities increasing the temporal resolution beyond the fluorescence signal acquisition rate. BSD also uniquely provides theoretically-grounded estimates of the expected performance of the spike reconstruction in terms of precision-recall and temporal accuracy for each recording. The performance of the algorithm is established using synthetic data and through the SpikeFinder challenge, a community-based initiative for spike-rate inference benchmarking based on a collection of joint electrophysiological and fluorescence recordings. Our method outperforms classical sparse deconvolution algorithms in terms of robustness, speed and/or accuracy and performs competitively in the SpikeFinder challenge. This algorithm is modular, easy-to-use and made freely available. Its novel features can thus be incorporated in a straightforward way into existing calcium imaging packages.

Keywords

Calcium Imaging, Fluorescence Microscopy, Inference, Generative Models

Highlights

- Functional calcium imaging allows one to monitor large neuronal networks yet the fluorescence signal is only an indirect measure of the neural activity.
- Here we introduce a Blind Sparse Deconvolution (BSD) algorithm for inferring spike trains from fluorescence recordings.
- BSD features fully-unsupervised estimation of metaparameters, and temporal super-resolution.
- It provides theoretical bounds on the expected precision-recall performance and temporal accuracy.
- BSD is shown to outperform standard sparse deconvolution algorithm in terms of speed and/or accuracy.
- It is modular, easy-to-use and made freely available to the community on github servers.

Introduction

In the last two decades, functional calcium imaging has emerged as a popular method for recording brain activity *in vivo*. This technique relies on calcium sensors, either synthetic or genetically expressed, that are designed to optically report the transient rise in intra-cellular calcium concentration that accompany spiking events. Compared to standard electrophysiology methods, calcium imaging is non-invasive, allows monitoring extended neuronal networks (up to a few tens of thousands of units) and can be combined with genetic methods in order to target specific neuronal populations. However, calcium imaging only provides a proxy measure of the neuronal activity. The kinetics of the calcium/reporter complexation being relatively slow, a spike-evoked fluorescence transient lasts much longer (0.1-1s)

than the action potential itself (<5ms). As the fluorescence signal is generally noisy and/or weakly sampled, its interpretation heavily relies on deconvolution methods to infer approximated spike trains. With the rapid increase in data-throughput offered by current fast imaging techniques [1–4], these methods need to be fast and unsupervised, as any manual check of the produced inference signals would be prohibitively tedious. Due to the high noise level, simple inference methods such as naive linear deconvolution and Wiener filtering prove inadequate. In the last decade, numerous alternative deconvolution algorithms have thus been proposed [1, 5–18]; among them, a powerful family of algorithms is based on non-negative sparse deconvolution [8, 16, 17]. In short, it consists in solving the inverse problem using the *a priori* knowledge that the spikes are sparse and non-negative. This framework, introduced by Vogelstein *et al.* in [8], was shown to efficiently recover spike trains from fluorescence signals. However, its performance is strongly dependent on the algorithm’s hyperparameters, namely the sparsity prior that controls the mean spike rate and the time constants characterizing the calcium reporter dynamics.

Despite extensive efforts for automatically adjusting these parameters [16, 17], progress are still needed to achieve the adequate robustness of the inference [19]. Another drawback of current algorithms is that the interpretation of the output can be challenging due to a paucity of theoretical understanding of their performance. No information is provided regarding the expected error rate of the inference or the time-precision of the output signal, *i.e.* the probability that a given spike be inferred in advance or delayed with respect to the true spike. Such information would be highly valuable, not only for downstream analysis but also for prior experimental design. In functional imaging, a trade-off has to be made between the sampling frequency, the signal-to-noise ratio and the imaged volume (that in turn sets the number of recorded neurons). Such choice of experimental parameters could be rationalized if one could foresee the achievable performance in terms of spike detection and timing precision for any given configuration. To address these various requirements, the inference algorithm should thus be accurate enough to make the most of the data, while being simple enough to be interpretable. It should provide a robust and unbiased extraction of the experimental parameters, such as signal-to-noise ratio or kernel shape, from the raw fluorescence datasets. It should finally offer a theoretically-grounded estimate of the inference performance for any given value of these parameters.

In the present study, we build on a recently proposed non-negative sparse inference method to develop the so-called *Blind Sparse Deconvolution (BSD)* algorithm. This novel implementation features automatic estimation of the hyperparameters, enhanced speed, similar-to-better reconstruction performances and super-resolution capabilities. We additionally provide thresholding guidelines and theoretical bounds on its performance, in terms of inference efficiency and temporal accuracy as a function of the experimental parameters. These various features are benchmarked on both synthetic and real data, covering a large spectrum of experimental contexts.

I. MATERIALS AND METHODS

A. Generative model

Standard inference methods are based on a generative model, which describes the relationship between a spike train and the resulting fluorescence time trace. It reads:

$$F_i \equiv F(t_i) = a \int K(\tau)N(t_i - \tau)d\tau + b + \epsilon_i \quad (1)$$

where $t_i = i\Delta t$ is the time of measurement, $N(t) = \sum_j \delta_{t,t_j}$ denotes the spike train, b is the baseline fluorescence (spikeless signal), and ϵ_i is a discrete gaussian white noise: $\langle \epsilon_i \rangle = 0$, $\langle \epsilon_i \epsilon_j \rangle = \sigma^2 \delta_{i,j}$. The convolution kernel $K(t)$, which reflects the complexation kinetics, is of the form:

$$K(t) \propto (e^{-\frac{t}{\tau_d}} - e^{-\frac{t}{\tau_r}}) \mathbb{1}_{t>=0} \quad (2)$$

where the rise and decay time constants τ_r, τ_d – typically in the range of 10-100ms and 50-1000ms, respectively – mostly depend on the calcium indicator but can also vary with the targeted neuron. In the following, we normalize K such that $\max_t K(t) = 1$, hence each spike produces a transient of maximum height a . The signal-to-noise ratio (SNR) is thus defined as $\text{SNR} = \frac{a}{\sigma}$. The noise stems from fluctuations of intra-cellular chemical concentration, light source and detector noise, incorrect baseline estimation, and other modeling errors. Typical SNR values range from below 1 to 10.

This description corresponds to configuration in which reporter sensitivity and acquisition rates allow the detection of individual spikes. In many real situations, this is not the case and one can only detect bursts of spikes over timescales shorter than the sampling window. In this case, the signal N to be inferred is not binary anymore but

remains sparse, and it reflects the instantaneous spike rate. The amplitude a is a characteristic scale of the (non-zero) signal e.g. such that $\frac{\text{Var}(N)}{\langle N \rangle} = a$.

B. Non-Negative Sparse Deconvolution

We recall first the non-negative sparse deconvolution approach for inferring $N(t)$. We rewrite Eqn. 1 as:

$$\begin{aligned} F_i &= a \sum_l K(t_i - t_l) + b + \sigma \epsilon_i \\ F_i &= a \sum_{j=1}^T K[\Delta t(i - j + 1)] N_j + b + \sigma \epsilon_i \\ \Leftrightarrow \mathbf{F} &= a\mathcal{K}\mathbf{N} + \mathbf{b} + \sigma \epsilon \end{aligned} \quad (3)$$

where $i \in [1, T]$ is the time frame index, $\Delta t \equiv \frac{1}{f}$ is the sampling interval, \mathcal{K} is the convolution matrix $\mathcal{K}_{ij} = K[\Delta t(i - j + 1)]$ and $N_j = \int_{t'=(j-1)\Delta t}^{j\Delta t} N(t') \in \mathbb{N}$ is the number of spikes in the time interval $[(j-1)\Delta t, j\Delta t]$ [42]. Note that the first and second lines are not equivalent. The second expression implicitly assumes that:

- The boundary condition $N(t) = 0, \forall t < 0$ holds, which is generally true in recordings that start during inactive periods. This simplification can be easily relaxed for inference. [43]
- One can approximate $K(t_i - t_l) = K(i\Delta t - t_l)$ as $K[i\Delta t - (j_l - 1)\Delta t]$ where $j_l - 1 = \lfloor \frac{t_l}{\Delta t} \rfloor$. This discretization error is negligible when Δt is small [44], yet it ensures that the matrix \mathcal{K} is translation invariant, *i.e.* $\mathcal{K}_{ij} = \phi(i - j)$.

From Eqn. 3, a naive estimate for N can be written as:

$$\begin{aligned} \hat{\mathbf{N}} &= \frac{1}{a} \mathcal{K}^{-1}(\mathbf{F} - \mathbf{b}) \\ \Leftrightarrow \hat{\mathbf{N}} &= \arg \min_{\mathbf{N}} \left\{ \frac{1}{2} \sum_{i=1}^T [F_i - a(\mathcal{K}\mathbf{N})_i - b]^2 \right\} \end{aligned} \quad (4)$$

In practice however, this approach fails to recover any spike at typical noise level $SNR = 2.5$ as shown in Section II-A and illustrated in Figure 1. To understand this failure, one may reason in the continuous framework for which Eqn. 4 reads $\hat{N}(t) \propto \int K^{-1}(\tau) [F(t - \tau) - b] d\tau$. Here, the inverse convolution kernel K^{-1} is proportional to $\delta''(t) - \left[\frac{1}{\tau_r} + \frac{1}{\tau_d} \right] \delta'(t) + \frac{1}{\tau_r \tau_d} \delta(t)$ thus the naive deconvolution reads:

$$\hat{N} \propto \partial_t^2 F(t) + \left[\frac{1}{\tau_r} + \frac{1}{\tau_d} \right] \partial_t F(t) + \frac{1}{\tau_r \tau_d} F(t) \quad (5)$$

A naive estimator of the signal thus involves computing the derivatives of the original signal, and is therefore extremely sensitive to high frequency noise. An intuitive solution to mitigate this issue consists in filtering out the high frequency component before carrying out the deconvolution, as is the basis of the Wiener deconvolution method. Vogelstein et al. showed that it also performs poorly because such filtering smoothes out the fast rise of the spike-evoked fluorescence transients. In contrast, non-negative sparse deconvolution estimators achieve both filtering of the noise while preserving the high-frequency signal. They are given by the outcome of the following optimization problem:

$$\hat{\mathbf{N}} = \arg \min_{\mathbf{N} \geq 0} \left\{ \frac{1}{2} \sum_{i=1}^T [F_i - a(\mathcal{K}\mathbf{N})_i - b]^2 + \lambda N_i \right\} \quad (6)$$

or equivalently :

$$\hat{N} = \frac{1}{a} \arg \min_{N' \geq 0} \left\{ \frac{1}{2} \sum_{i=1}^T (F_i - (\mathcal{K}N')_i - b)^2 + \lambda' N'_i \right\} \quad (7)$$

where $\lambda, \lambda' = \frac{\lambda}{a}$ are L_1 penalty coefficients that control the sparsity of the optimum (the higher λ , the sparser the optimum). When $\lambda = 0$ and the $N \geq 0$ constraint is relaxed, the optimal value \hat{N} is exactly given by Eqn. 4. As shown in the next Section, the choice of λ is crucial for efficient denoising and proper spike inference. Notice that the optimization problem is convex and can be solved efficiently in $\mathcal{O}(T)$ for double-exponential kernels using the interior-point method, see [8]. This unusual linear scaling for a matrix inversion-like operation owes to the fact that \mathcal{K}^{-1} is tridiagonal for double-exponential kernels: $\mathcal{K}_{ij}^{-1} \propto \delta_{ij} - \gamma_1 \delta_{i,j+1} + \gamma_2 \delta_{i,j+2}$, with $\gamma_1 = \exp\left(-\frac{\Delta t}{\tau_r}\right) + \exp\left(-\frac{\Delta t}{\tau_d}\right)$, $\gamma_2 = \exp\left(-\frac{\Delta t}{\tau_r} - \frac{\Delta t}{\tau_d}\right)$. In [17], the authors apply the Pool-Adjacent Violator Algorithm originally developed for isotonic regression problems to solve this optimization in an even faster but approximate way.

C. Determination of the sparsity prior λ and signal threshold value

The choice of the regularization parameter λ is crucial. If it is too large, the inferred spike train is $\hat{N} = 0$ and all spikes are missed, whereas if it is too small, noise-induced transients are interpreted as spikes yielding large false positive rates. Intuitively, we expect the optimal choice to depend on the parameters of the generative model (noise level, spike amplitude, etc.) Here we review the expressions of λ previously used and we then introduce our method. We adopt the convention from Eqn. 7 and drop the primes. We assume for now that all generative model parameters are known.

1. Existing methods: fast-oopsi and constrained-oopsi

In [8], the authors derive the non-negative sparse deconvolution from an approximate Maximum A Posteriori principle. They assume that the spike count N_i at time step i follows a Poisson distribution of mean $\nu \Delta t$, where ν is the firing rate. After approximating the Poisson prior with an exponential distribution, they compute the negative log-likelihood $-\log P(F, N)$, which they find to be proportional to (7) with a sparsity prior λ given by:

$$\lambda_{oopsi} = \frac{\sigma^2}{a\nu\Delta t} \quad (8)$$

This approach thus provides an analytical expression for λ . However, using an exponential approximation instead of Poisson can largely overestimate the threshold required (it is illustrated for a simple example in Annex I), and result in improper behavior in several realistic experimental conditions as shown in Section II-A and illustrated in Figure 1. To address this issue, a non-analytical method called constrained-oopsi (referred to as con-oopsi in the following) was recently introduced in [16]. The authors propose the following constrained deconvolution:

$$\begin{aligned} \hat{N} &= \arg \min_{N \geq 0} \sum_i N_i \\ \text{subject to } &\sum_i [F_i - (\mathcal{K}N)_i]^2 \leq \sigma^2 T \end{aligned} \quad (9)$$

Where T is the number of observations. The problem can be rewritten using the Karush-Kuhn Tucker conditions by introducing the Lagrangian $\mathcal{L} = \sum_i N_i + \rho \sum_i [F_i - (\mathcal{K}N)_i]^2$ where ρ is the Lagrange multiplier associated with the constraint. There exists ρ such that the critical point N^* of \mathcal{L} is the solution of the constrained optimization problem. Clearly, N^* satisfies the constraint only if ρ is non-negative; in this case \mathcal{L} is convex and the critical point is a minimum of \mathcal{L} . Overall, the optimization problem can be rewritten as:

$$\hat{N}(\rho) = \arg \min_{N \geq 0} \left\{ \sum_i N_i + \rho \sum_i [F_i - (\mathcal{K}N)_i]^2 \right\} \quad (10)$$

Identifying $\lambda = \frac{1}{2\rho}$, the constrained deconvolution is equivalent to a sparse deconvolution with an adaptive sparsity prior λ . Since $\sum_i \hat{N}_i(\rho)$ is a decreasing function of ρ , the expression for λ reads:

$$\lambda_{con-oopsi} = \max\{\lambda \in \mathbb{R}^+, \sum_i [F_i - (\mathcal{K}\hat{N})_i]^2 \leq \sigma^2 T\} \quad (11)$$

In practice, $\lambda_{con-oopsi}$ is found by alternatively solving Eqn. 7 and updating λ , decreasing it if the reconstruction error is too large, and increasing it otherwise. This non-analytical approach performs better than fast-oopsi (see Section II-A and Figure 1).

2. Blind Sparse Deconvolution

We propose a different analytical expression for λ , inspired by [20]. It is deduced from the analysis of the optimization problem for two simple configurations, in which there is either zero or one spike in the original signal. We show that this solution combines the computational speed of fast-oopsi and the robustness of con-oopsi.

3. Spikeless Signal

In the following, we use matrix notations and rewrite the cost function as :

$$\mathcal{L}(\mathbf{N}) = \frac{1}{2} [\mathbf{F} - \mathcal{K}\mathbf{N}]^T [\mathbf{F} - \mathcal{K}\mathbf{N}] + \lambda \mathbf{1}^T \mathbf{N} \quad (12)$$

The gradient writes:

$$\begin{aligned} -\nabla_{\mathbf{N}} \mathcal{L} &= \mathcal{K}^T (\mathbf{F} - \mathcal{K}\mathbf{N}) - \lambda \mathbf{1} \\ -\nabla_{\mathbf{N}} \mathcal{L} &= -(\mathcal{K}^T \mathcal{K}) \mathbf{N} + \mathcal{K}^T \mathbf{F} - \lambda \mathbf{1} \end{aligned} \quad (13)$$

Let's first assume that the signal is spikeless, such that $F_i = \sigma \epsilon_i$, where ϵ_i is a gaussian white noise. Since $(\mathcal{K}^T \mathcal{K})N > 0$, we have:

$$\begin{aligned} -\frac{\partial \mathcal{L}}{\partial N_i} &< \sigma (\mathcal{K}^T \epsilon)_i - \lambda \sim \mathcal{N}(-\lambda, \sigma^2 \|K\|^2) \\ &\Rightarrow P \left[-\frac{\partial \mathcal{L}}{\partial N_i} > 0 \right] < \Phi \left[\frac{\lambda}{\sigma \|K\|} \right] \end{aligned} \quad (14)$$

where $\Phi(x) = \int_x^{+\infty} \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}}$, and $\|K\| \equiv \sqrt{\sum_{i=-\infty}^{\infty} K^2 [i\Delta t]}$

Therefore, if $\lambda = \lambda_1 \equiv z_1 \sigma \|K\|$ with z_1 large enough, the gradients are almost always negative, and the global optimum of \mathcal{L} is $\hat{\mathbf{N}} = \mathbf{0}$. Hence for instance, setting $z_1 = 2.326$, yields a probability of false positive event per time bin $P_{FP} < 0.01$.

4. Single spike signal

We now examine a configuration in which a single spike is present in the data :

$$\begin{aligned} N_i^0 &= \delta_{i,i_0} \\ F_i &= aK [\Delta t(i - i_0 + 1)] + \sigma \epsilon_i \end{aligned} \quad (15)$$

The gradient writes:

$$-\nabla_{\mathbf{N}} \mathcal{L} = -(\mathcal{K}^T \mathcal{K})(\mathbf{N} - a\mathbf{N}^0) + \sigma K^T \epsilon - \lambda \mathbf{1} \quad (16)$$

We look for an optimum of the form $\hat{N}_i = an\delta_{i,i_0}$. The optimization with respect to n gives:

$$\begin{aligned} n &= \max \left\{ 1 - \frac{\lambda}{a \sum_{i=1}^T K^2 [(i - i_0 + 1)\Delta t]} + \frac{\sigma \sum_{i=1}^T K [(i - i_0 + 1)\Delta t] \epsilon_i}{a \sum_{i=1}^T K^2 [(i - i_0 + 1)\Delta t]}, 0 \right\} \\ n &\sim \max \left[\mathcal{N} \left(1 - \frac{\lambda}{a \|K\|^2}, \frac{\sigma^2}{a^2 \|K\|^2} \right), 0 \right] \end{aligned} \quad (17)$$

where the last line assumes that $\sum_{i=1}^T K^2 [(i - i_0)\Delta t] \approx \sum_{i=-\infty}^{\infty} K^2 [i\Delta t] \equiv \|K\|^2$, which is true provided that i_0 is far from the boundaries. Thus, if the spike position is known in advance, the inferred spike is a thresholded gaussian variable.

Importantly, the effective noise level $\sigma' = \frac{\sigma}{a\|K\|}$ that appears in this expression is smaller than $\frac{\sigma}{a}$ by a factor $\frac{1}{\|K\|}$. This has an important consequence: since $\max(K) = 1$, the norm $\|K\| = \sqrt{\sum_i K(i\Delta t)^2}$ of the discretized kernel is proportional to \sqrt{M} where M is the typical number of time frames over which K is non-zero. Thus, $\sigma' \sim \frac{\sigma}{a\sqrt{M}}$ as if the noise had been averaged over the duration of the transient. This suggests that even signals with low SNR can be efficiently inferred provided that the spike-induced fluorescence transient is sufficiently well sampled.

Eqn. 17 shows that when λ is too large, the probability that a given spike is undetected reads:

$$P_{FN} = P[n = 0] = \Phi \left[\frac{\|K\| \left(a - \frac{\lambda}{\|K\|^2} \right)}{\sigma} \right] \quad (18)$$

Therefore, setting $\lambda \leq \lambda_2 \equiv \|K\|^2 a - z_2 \sigma \|K\|$, with, say, $z_2 = 2.326$, guarantees a low false negative rate (FNR) as the probability that a spike is detected is then larger than 0.99.

The sparsity prior λ_{BSD} is chosen to minimize both the FPR and FNR. For $a = 1$, $\sigma = 0.1$, $\tau_r = 0.1$, $\tau_d = 0.5$, $f = 10Hz$, $\lambda_2 = 4.1379$ is much higher than λ_1 . In this case, setting $\lambda_{BSD} = \lambda_1$ is the best solution, as smaller values of λ_{BSD} lead to less signal deformation. In contrast, for configurations such that $\lambda_1 > \lambda_2$, *i.e.* when $\sigma > \sigma^{max} = \frac{a\|K\|}{z_1 + z_2}$, it is impossible to satisfy both constraints (low FPR and low FNR); in this case we use the crossover value $\lambda = \frac{z_1 a \|K\|}{z_1 + z_2} = \lambda_1(\sigma^{max})$.

5. Sparsity prior for BSD

To summarize, in our Blind Sparse Deconvolution (BSD) algorithm, the sparsity prior is set analytically as:

$$\lambda_{BSD} = z_1 \|K\| \min \left(\sigma, \frac{a\|K\|}{z_1 + z_2} \right) \quad (19)$$

where $\|K\| = \sqrt{\sum_{i=-\infty}^{+\infty} K(i\Delta t)^2}$ is the L_2 norm of the discretized convolution kernel K , and z_1, z_2 are two numbers ~ 2 that control the precision and recall, respectively.

We expect the sparse deconvolution to perform more consistently with λ_{BSD} than with λ_{oopsi} . Indeed, if $a = 1$, $\sigma = 0.1$, $\tau_r = 0.1$, $\tau_d = 0.5$, $\Delta t = 0.1s$, $\nu = 1Hz$, we find $\lambda_{oopsi} = 0.1$ and $\lambda_1 = 0.51$; using λ_{oopsi} therefore results in multiple noise-induced false spikes. Conversely, for $\sigma = 0.25$ and $\nu = 0.1Hz$, we have $\lambda_1 = 1.25$, $\lambda_2 = 3.39$, $\lambda_{oopsi} = 6.25$; in this case, λ_{oopsi} is too large and most, if not all the spikes are missed.

Robust performances are also expected using $\lambda_{con-oopsi}$, although a slightly larger FPR is expected compared to λ_{BSD} . Indeed, in the absence of spikes, $N = 0$ satisfies the reconstruction constraint in the large T limit and is correctly found by the algorithm. In the presence of spikes, we expect $\lambda_{con-oopsi}$ to be slightly lower than λ_{BSD} , resulting in small overfitting of the noise. Indeed, as soon as $\lambda > 0$, the spikes are on average underestimated, see Eqn. 17; therefore, any good choice of λ that perfectly filters the noise also underestimates the reconstructed trace KN , yielding a reconstruction error $\sum_{i=1}^T (F_i - anK[(i - i_0)\Delta t])^2 > \sigma^2 T$ and is therefore not a valid solution for the constrained optimization. Instead, λ is further decreased until false (noise-induced) spikes appear and reduce the reconstruction error below $\sigma^2 T$.

6. Thresholding the BSD inferred signal

Some applications, such as network connectivity inference, may require to threshold the signal in order to get a binary spike train. Unlike previous methods, BSD provides rationale for choosing a threshold. Indeed, the computations performed in Sections I-C3-4 show that the (unnormalized) inferred spikes in the absence (*resp.* presence) of spikes are thresholded gaussian variables, with means $-\frac{\lambda}{\|K\|^2}$ and $a - \frac{\lambda}{\|K\|^2}$, respectively, and identical variance $\frac{\sigma^2}{\|K\|^2}$. Picking a threshold that separates the two distributions yields:

$$\theta = \min \left[z_3 \frac{\sigma}{\|K\|}, u \left(a - \frac{\lambda_{BSD}}{\|K\|^2} \right) \right] \quad (20)$$

where z_3 is a quantile of the normal distribution, and u a number between 0 and 1, say, 0.5. When θ equals the left term, the vast majority of the noise is efficiently filtered out such that any non-zero value in the output signal can be safely assigned to a spike; the right-hand term in turn prevents the threshold from becoming larger than the signal itself.

D. Theoretical limits on the precision-recall and temporal resolution

We now present a similar analysis to derive theoretical estimates of the true and false positive rates, and of the temporal accuracy of the predicted spikes trains as function of the generative model parameters. The corresponding scripts are also implemented in the BSD package, and can be directly applied for a given fluorescence recording once the generative model parameters have been inferred.

1. Precision-Recall for isolated spikes

The theoretical false positive and negative rates (FPR, FNR) are first computed within the sparse deconvolution framework with λ_{BSD} . For the false positive rate, the computation was performed in Section I-C3: we obtain a probability of false positive rate per time bin:

$$P_{FP} = \Phi \left[z_1 \min \left(1, \frac{a\|K\|}{\sigma(z_1 + z_2)} \right) \right] \quad (21)$$

For the false negative rate, we follow a similar reasoning as in Section I-C4: we consider a signal of the form $F_i = an_0K[\Delta ti - t_0] + \sigma\epsilon_i$ with $t_0 = (i_0 - 1)\Delta t + \delta t_0$ and $0 \leq \delta t_0 < \Delta t$. Note that we have now relaxed the previously made approximation $K[\Delta ti - t_0] \approx K[\Delta t(i - i_0 + 1)]$ in order to probe the effect of intermittent sampling. We obtain a lower bound [45] for the probability of false negative per spike:

$$\begin{aligned} \hat{N} &= \arg \max_{N \geq 0} \mathcal{L}(N) \approx an\delta_{i, i_0} \\ \implies n &\sim \max \left[\mathcal{N} \left(n_0 \cos \theta(-\delta t_0) - \min \left(z_1 \tilde{\sigma}, \frac{z_1}{z_1 + z_2} \right), \tilde{\sigma}^2 \right), 0 \right] \\ \implies P_{FN}(\delta_0) &= \Phi \left[\frac{n_0 \cos \theta(-\delta t_0) - \min \left(z_1 \tilde{\sigma}, \frac{z_1}{z_1 + z_2} \right)}{\tilde{\sigma}} \right] \end{aligned} \quad (22)$$

where:

$$\begin{aligned} \cos \theta(\delta t) &= \frac{\sum_{l=-\infty}^{\infty} K[\Delta tl] K[\Delta tl + \delta t]}{\sum_{l=-\infty}^{\infty} K[\Delta tl]^2} \\ \tilde{\sigma} &= \frac{\sigma}{a\|K\|} \end{aligned} \quad (23)$$

Note that the probability depends on δt_0 ; for instance if $\tau_r = 0$ and $\delta t_0 \ll \Delta t$, spikes emitted right after a measurement yield low-amplitude fluorescent transients and are thus likely to be missed. Overall, the probability of false negative is given by:

$$P_{FN} = \frac{1}{\Delta t} \int_{\delta t_0=0}^{\Delta t} P_{FN}(\delta t_0) d\delta t_0 \quad (24)$$

2. Temporal Resolution for isolated spikes

Intuitively, the temporal resolution depends on three factors: the signal-to-noise ratio, the sampling rate, the shape of the fluorescence kernel. It is characterized by the point-spread function (PSF) of the inferred spikes with respect to the true spikes, namely the conditional average given a fluorescence signal with a single-spike at $t = 0$:

$$PSF_\tau = \mathbb{E} \left[\hat{N}_{t_0+\tau} | N_t^0 = \delta_{t,t_0} \right] \quad (25)$$

Where $\hat{\mathbf{N}} = \arg \min_{\mathbf{N}} \mathcal{L}(\mathbf{N})$ is the spike train inferred from the fluorescence signal see Eqn. 12 and the expectation is taken over the Gaussian noise realizations, see Eqn. 1. The distribution of $\hat{\mathbf{N}}$ is not tractable and, therefore, the PSF cannot be derived by exact analytical computation. Instead, we use two heuristics to obtain analytical insights into the width of the PSF and an efficient numerical approximation of the PSF.

For the analytical computation, we focus on the distribution of the initial negative gradients $-\frac{\partial \mathcal{L}}{\partial N_i} |_{\mathbf{N}=0}$ rather than the one of $\hat{\mathbf{N}}$. Consider indeed the gradient descent optimization dynamics: because of the L_1 penalty, large components N_i tend to grow faster and to screen neighboring small components, yielding sparse solutions with only few non-zero components. It is therefore likely that the largest components of \mathbf{N} after one gradient descent step (after which $N_i \propto -\frac{\partial \mathcal{L}}{\partial N_i} |_{\mathbf{N}=0}$) remains the largest at the end of the optimization. Hence if the initial negative gradient is larger at position $i_0 + \delta$ than at position i_0 , we expect the inferred spike $\hat{\mathbf{N}}$ to be similarly delayed with respect to the true spike position. The probability of such an error can be computed as:

$$\begin{aligned} -\frac{\partial \mathcal{L}}{\partial N_{i_0}} |_{N=0} &= a \|K\|^2 + \sigma \sum_i K[\Delta t(i - i_0 + 1)] \epsilon_i - \lambda \\ -\frac{\partial \mathcal{L}}{\partial N_{i_0+\delta}} |_{N=0} &= a \|K\|^2 \cos \theta(\delta \Delta t) + \sigma \sum_i K[\Delta t(i - i_0 - \delta + 1)] \epsilon_i - \lambda \\ \Rightarrow \Delta \left[-\frac{\partial \mathcal{L}}{\partial N} \right] &\sim -2a \|K\|^2 \sin^2 \frac{\theta(\delta \Delta t)}{2} + 2\sigma \|K\| \sin \frac{\theta(\delta \Delta t)}{2} \mathcal{N}(0, 1) \end{aligned} \quad (26)$$

where the angle $\theta(\delta t)$ is defined in Eqn. 23.

Thus, the initial gradient at the offset time $i_0 + \delta$ is higher than its value at the spike time i_0 with probability $\Phi \left[\frac{a \|K\| \sin \frac{\theta(\delta \Delta t)}{2}}{\sigma} \right]$, typically resulting in a time-shifted inferred spike. This results in a typical timing error δt on the spike position of the order of:

$$\delta t \text{ s.t. } \sin \frac{\theta(\delta t)}{2} = \frac{\sigma}{a \|K\|} \quad (27)$$

This timing error is a non-trivial function of the kernel K and the noise level. The higher the effective noise level $\frac{\sigma}{a \|K\|}$, the higher δt . The second factor is small for rapidly growing $\theta(\delta t)$, *i.e.* when the overlap between the kernel $K(t)$ and its lagged version $K(t + \delta t)$ is a fast decaying function of δt . Hence, the 'sharper' the kernel, the lower δt .

For the numerical approximation of the PSF, we restrict the optimization search space to solutions of the form $\hat{N}_i(\tau, n) = a n \delta_{i, i_0+\tau}$. Using this simplification, the optimization over n can be carried out analytically (similarly as Eqn. 22) for each τ , and PSF_τ is given by the probability that the optimal solution is located at τ . After rearrangement, we obtain:

$$PSF_\tau \approx P(X_\tau \geq 0 \ \& \ X_\tau \geq X_{\tau'} \ \forall \tau') \quad (28)$$

Where $X_\tau, \tau \in \mathbb{Z}$ is a Gaussian vector with mean $\mathbb{E}[X_\tau] = \cos(\theta(\tau \Delta t)) - \min \left[\frac{z_1 \sigma}{a \|K\|}, \frac{z_1}{z_1 + z_2} \right]$ and covariance $\text{Cov}(X_\tau, X_{\tau'}) = \frac{\sigma}{a} \cos \theta(\Delta t |\tau - \tau'|)$.

Eqn. 28 is efficiently evaluated by Monte Carlo, with values of $|\tau| \leq 20$ and is valid for the discrete generative model; the computation can be generalized for the continuous generative model, see Annex D. The final formula is given by:

$$PSF(\delta t) = \frac{1}{\Delta t^2} \int_{\delta t_0=0}^{\Delta t} \sum_{\tau \in \mathbb{Z}} PSF_{\tau}(\delta t_0) \mathbb{1}_{\tau \Delta t - \delta t_0 < \delta t} \quad (29)$$

Where $PSF_{\tau}(\delta t_0)$ is the discrete PSF of Eqn. 28, but with a term $\cos(\theta(\tau \Delta t - \delta t_0))$ instead of $\cos(\theta(\tau \Delta t))$ in the mean value of X_{τ} .

E. Hyperparameters learning

All sparse deconvolution methods rely on the knowledge of the generative model's parameters. However, owing to the variability in the calcium reporters intracellular concentration and other biochemical cellular processes, these parameters may significantly vary from experiment to experiment, and for different neuronal types. In the fast-oopsi implementation, the authors proposed to infer the parameters (a, b, σ, ν) in an iterative way: an initial guess is made, deconvolution is performed, parameters values are then updated based on the deconvolution result, whereas for con-oopsi, the authors propose to estimate σ, K only once. We follow the same iteration-based approach as fast-oopsi, but the parameters are inferred and refined differently; we also add a method to infer and refine the kernel K .

1. Initial estimation of the parameters

We are given a time series of the form $\mathbf{F} = a\mathbf{K}\mathbf{N} + \sigma\epsilon + b$, with unknown a, b, σ, K . In the following, we assume that the baseline is constant or equivalently that the variable baseline has been previously estimated and subtracted from the signal. From the knowledge that N is non-negative and sparse, we deduce that:

- The baseline b is essentially the most often observed value of F ; the data histogram is computed, and b is estimated as the center of the interval with highest frequency. Using the median of S also provides a good estimator.
- All activity below the baseline originates from the noise, hence $F' = F[F < b] - b$ follows a half-Gaussian distribution $\min[\mathcal{N}(0, \sigma^2), 0]$; it is fitted to deduce σ .

We estimate the convolution kernel K through the signal auto-correlation matrix. Indeed, observe that :

$$A_F(l) \equiv \langle F_i F_{i+l} \rangle - \langle F \rangle^2 = a^2 \sum_{j,k} [\langle N_j N_{j+k} \rangle - \langle N \rangle^2] \sum_i K[(i-j-1)\Delta t] K[(i+l-k-j-1)\Delta t] + \sigma^2 \delta_{l,0} \quad (30)$$

Under the assumption that the spiking events N_i are independent, identically distributed Poisson variables, we have $\langle N_j N_{j+k} \rangle - \langle N \rangle^2 = a^2 \nu \Delta t \delta_{k,0}$ and Eqn. 30 can be simplified as:

$$\begin{aligned} A_F(l) &= a^2 \nu \Delta t \sum_{j=-\infty}^{\infty} K[j\Delta t] K[(l+j)\Delta t] + \sigma^2 \delta_{l,0} \\ \iff (A_F(l) - \sigma^2 \delta_{l,0}) &\propto \sum_{j=-\infty}^{\infty} K[j\Delta t] K[(l+j)\Delta t] \end{aligned} \quad (31)$$

The auto-correlation matrix can be estimated from the data as $\hat{A}_F(l) = \frac{1}{T} \sum_i F_i F_{i+l} - \left[\frac{\sum_i F_i}{T} \right]^2$. Together with the previous estimate of σ , the left-hand side of the equation can thus be estimated. In practice, \hat{A}_F obtained is not necessarily positive definite, because the estimate of σ can be incorrect - this can lead to very bad estimates of τ_r, τ_d . To mitigate this issue, we subtract $\min(\sigma^2, \lambda_{min})$ instead of σ^2 , where λ_{min} is the smallest eigenvalue of the Toeplitz autocorrelation matrix. The right-hand side is the overlap between the kernel K and its delayed version $K'(t) = K(t + l\Delta t)$. We can normalize both terms to 1 for $l = 0$, and use a least square fit to estimate K .

Lastly, the spike amplitude a and frequency ν can be deduced from the following equations, that hold under the model assumption:

$$\begin{aligned} \langle F \rangle &= a\nu\Delta t \sum_i K[i\Delta t] \\ \langle F^2 \rangle - \langle F \rangle^2 &= a^2\nu\Delta t \sum_i K[i\Delta t]^2 \end{aligned} \quad (32)$$

Although they yield very good results for synthetic datasets, these estimators can fail in several frequently encountered situations in practice:

- When the neural activity is not sparse, we do not expect b to be the most frequent fluorescence value. An error in the estimation of b can result in a misestimation of σ as well.
- When the neuron displays bursting activity (*i.e.* several spikes in short time intervals), the hypothesis that the N_i are independent usually fails. This may result in overestimating τ_r and/or τ_d .
- In the same situation, Eqn. 32 is incorrect and a can be overestimated.
- When the noise exhibits temporal correlation (streaking artefacts in light sheet imaging, small sample drifts, fluctuations in laser intensity, etc.), the white-noise hypothesis does not hold, which may result in a misestimation of τ_r and τ_d .

When the estimated time constants τ_r and τ_d differ from their true values, τ_r^0 and τ_d^0 , systematic estimation errors arise. Suppose for instance that $\tau_r < \tau_r^0$ and $\tau_d = \tau_d^0$. Then a spike-induced fluorescence transient tends to exhibit a faster initial rise than expected. Hence, from a Bayesian perspective, such a transient is likely to be interpreted as two small consecutive spikes. Hence, inferred spikes will tend to be duplicated. In general, the nature of the error depends on the kernel mismatch; some simulation results are presented in Annex B. These results highlight the need to refine the kernel parameters estimators.

2. Iterative parameter estimation: adaptive blind deconvolution

At fixed hyperparameters, the cost function to minimize is the sum of a reconstruction error and a sparse penalty on N :

$$\mathcal{L}(\mathbf{N}, K, b, \sigma, a) = \frac{1}{2} \|\mathbf{F} - \mathcal{K}\mathbf{N} - \mathbf{b}\|^2 + \lambda_{BSD}(\|K\|, \sigma, a) \|\mathbf{N}\|_1 \quad (33)$$

We can further refine the model hyperparameters by jointly minimizing this cost function with respect to both the spike train and the hyperparameters. In other words, we look for the convolution kernel that achieves the best trade-off between inferred spikes sparsity and reconstruction error. Such optimization yields non-trivial kernels, since a very sharp kernel ($\tau_r, \tau_d \rightarrow 0$) would give a perfect reconstruction but dense spikes, whereas a wide kernel would give very sparse spikes but poor reconstruction. Furthermore, it can be shown that such sparsity-reconstruction trade-off maximization, which is also featured in Sparse Dictionary Learning [21] or Blind Deconvolution in image deblurring [22] is equivalent to a Maximum A Posteriori optimization of the likelihood, assuming a gaussian noise, a sparse prior for spikes and a flat prior for the hyperparameters [22]. In practice, we optimize over K and b iteratively through the following coordinate descent algorithm:

$$\begin{aligned} \hat{\mathbf{N}}^{(t)} &\leftarrow \arg \min_{\mathbf{N}} \mathcal{L}(\mathbf{N}, \mathbf{K}^{(t-1)}, \mathbf{b}^{(t-1)}, \sigma, \mathbf{a}) \\ (\hat{K}^{(t)}, \hat{\mathbf{b}}^{(t)}) &\leftarrow \arg \min_{K, b} \mathcal{L}(\hat{\mathbf{N}}^{(t)}, K, \mathbf{b}, \sigma, a) \end{aligned} \quad (34)$$

The first optimization step was discussed in Section I-B. The second optimization is a parametric temporal regression problem; it can be solved efficiently in $\mathcal{O}\left(\frac{\tau_r + \tau_d}{\Delta t}\right)$ by introducing the cross-correlation $X_\tau = \frac{1}{T} \sum_i F_i \hat{N}_{i-\tau}$ and auto-correlation $A_\tau = \frac{1}{T} \sum_i \hat{N}_i \hat{N}_{i-\tau}$ functions up to some cut-off $\tau_m \sim \frac{\tau_r + \tau_d}{\Delta t}$ (see details in Annex C). The purpose of this step is that if $N^{(t)} = N_0$, then the optimum is exactly K_0 if σ is small or T is large. More generally (N_0, K_0) is a fixed point of the optimization dynamic in the low noise limit, and intuitively, we expect that at finite noise, another

fixed point close to (N_0, K_0) exists and can be reached. We show in Annex C that K^0 is the global optimum in the case of isolated spikes and low noise level. The optimization will not necessarily converge to such solution because the function $\mathcal{L}(N, K)$ is not convex, and only local minima are found. In practice, the optimum is usually very close to the original convolution kernel, and is reached if the initial estimate is good enough. The convergence can be improved by thresholding the spikes before updating the kernel, as it prevents false spikes from contributing to the cross-correlation. The iterative process is no longer an optimization but it still converges.

The noise σ and spike amplitude a can be refined as well, using

$$\hat{a} = \frac{\sum_t \hat{N}'_t}{\sum_t 1_{\hat{N}'_t > 0}} + \frac{\lambda}{\|K\|^2} \quad (35)$$

Where the last term corrects the bias due to the sparse prior (see Eqn. 17)

$$\hat{\sigma} = \sqrt{\frac{1}{T} \sum_t [(F_i - (\mathcal{K}N')_i - b)^2]} \quad (36)$$

F. Super-resolution

Most fluorescent microscopy techniques –two-photon, confocal or light-sheet – involves the sequential scanning of a laser beam at different locations within the sample. Hence, for a given dwell time of the laser at each neuron position, there is a trade-off between the sampling rate and the total number of sampled neurons. In other experimental fields, resolution limitations due to recording constraints have been significantly circumvented through signal processing algorithms. For instance, super-resolution microscopy achieves imaging at higher resolution than the diffraction limit [23–27], and compressed sensing applied to MRI allows to drastically reduce the number of measurements required to reach a given resolution [28]. These algorithms rely on the hypothesis that the original signal is sparse in a certain basis; it is therefore tempting to apply them to our problem, given that neural spikes are sparse in the canonical basis. This possibility had been discussed in the context of bayesian inference [7]. Temporal resolution was shown to be slightly improved in very specific settings, *i.e.* when using prior knowledge of inputs (stimulus) and spiking history dependence of the neuronal activity. In this section, we extend the Blind Sparse Deconvolution framework to super-resolution, *i.e.* we develop a method to infer spiking events timing with a temporal resolution beyond the sampling rate.

1. Qualitative analysis

We start off with a qualitative analysis and consider the fluorescence signal produced by an isolated single spike of amplitude n :

$$F_i = anK(\Delta ti - t_0) + \sigma\epsilon_i + b \quad (37)$$

Denoting $j = \lfloor \frac{t_0}{\Delta t} \rfloor$, $\delta t = t_0 - j\Delta t \in [0, \Delta t]$, $\lambda_d = e^{-\frac{\delta t}{\tau_d}}$, $\lambda_r = e^{-\frac{\delta t}{\tau_r}}$, and assuming for simplicity that $b = 0$ and that K is unnormalized, we write, for $i > j$:

$$\begin{aligned} F_i &= an \left[\lambda_d^{i-j-\frac{\delta t}{\Delta t}} - \lambda_r^{i-j-\frac{\delta t}{\Delta t}} \right] + \sigma\epsilon_i \\ \iff F_i &= an\lambda_d^{-\frac{\delta t}{\Delta t}} \lambda_d^{i-j} - an\lambda_r^{-\frac{\delta t}{\Delta t}} \lambda_r^{i-j} + \sigma\epsilon_i \end{aligned} \quad (38)$$

Thus, the observed fluorescence is a double exponential with non-equal coefficients of the form $f(i) = A\lambda_d^i + B\lambda_r^i$. Fitting the coefficients with a least-square method yields estimates of $an\lambda_d^{-\frac{\delta t}{\Delta t}}$ and $an\lambda_r^{-\frac{\delta t}{\Delta t}}$, which can be converted to estimates of δt and n . Thus, it is possible in principle to find the exact spike position in the noiseless case, if we know a priori that the signal contains a single spike. Notice that this is possible only if $\lambda_r > 0$, *i.e.* $\tau_r > 0$; if $\tau_r = 0$, the observed fluorescence is a single exponential of amplitude $an\lambda_d^{-\frac{\delta t}{\Delta t}}$, and we cannot recover both n and τ without ambiguity [46]. In the case of a noisy signal, we expect that super-resolution can be achieved only if $\frac{\tau_r}{\Delta t}$ is large enough

with respect to some function of σ . Notice also that if multiple spikes occur within the same time bin, the observed fluorescence transient is still a double exponential with non-equal coefficients, and it cannot be distinguished from the one produced by a single large spike at some average position. More generally, resolving two spikes in the same time bin would require the use of more complex convolution kernels.

2. Generative model

With these limitations in mind, we now extend the deconvolution framework to implement super-resolution. The fluorescence signal is constructed using a discrete generative model at a fine-grained time scale $\frac{\Delta t}{s}$, where s is a non-zero integer, which is then down-sampled by the same factor s . This yields the following generative model:

$$\begin{aligned} F_k^s &= a \sum_{j=1}^{sT} K \left[\frac{\Delta t}{s}(k-j+1) \right] N_j^s + b + \sigma \epsilon_k^s \\ F_i &\equiv F_{is}^s = a \sum_{j=1}^{sT} K \left[i\Delta t - (j-1)\frac{\Delta t}{s} \right] N_j^s + b + \sigma \epsilon_i \\ \Leftrightarrow \mathbf{F} &= a\mathcal{K}\mathbf{N}^s + \mathbf{b} + \sigma\epsilon \end{aligned} \quad (39)$$

where F_i is the fluorescence measurement at $t_i = i\Delta t$ and $N_j^s = \int_{(j-1)\frac{\Delta t}{s}}^{j\frac{\Delta t}{s}} N(t)dt$ is the number of spikes emitted in the time interval $[(j-1)\frac{\Delta t}{s}, j\frac{\Delta t}{s}]$. [47] The convolution matrix \mathcal{K} is now rectangular, of size $T \times sT$. It is not translation invariant anymore with respect to the spikes index j as the norm of the transient, $\|K_j\| = \sqrt{\sum_i \mathcal{K}_{ij}^2}$ now depends on j . Indeed, writing $j = (p-1)s + r$, we have:

$$\|K_j\| = \sqrt{\sum_{i=1}^T K \left[i\Delta t - (j-1)\frac{\Delta t}{s} \right]^2} = \sqrt{\sum_{i=1}^T K \left[(i-p)\Delta t + \frac{s+1-r}{s}\Delta t \right]^2} \approx \sqrt{\sum_{k=-\infty}^{\infty} K \left[k\Delta t + \frac{s+1-r}{s}\Delta t \right]^2} = f(r) \quad (40)$$

Typically, spikes occurring right after a fluorescence measurement (small r) have smaller $\|K_j\|$ than spikes occurring right before a measurement (large r).

3. Sparse Deconvolution

A sparse deconvolution algorithm is applied to estimate the spikes N^s :

$$\hat{\mathbf{N}}^s = \arg \min_{N^s \geq 0} \left\{ \frac{1}{2} \sum_{i=1}^T [F_i - a(\mathcal{K}N)_i - b]^2 + \sum_{j=1}^{Ts} \lambda^j N_j^s \right\} \quad (41)$$

Notice that, although \mathcal{K} is not invertible anymore, the optimum is still well-defined because of the sparsity penalty and non-negativity constraint. Compared to Eqn. 6, the main difference is that λ is not uniform anymore: $\lambda_j \propto \|K_j\|$. This property has an important consequence, as can be seen by considering the limit case $\tau_r = 0$, $\sigma \ll a$. As discussed previously, a transient observed for $i \geq i_0$ can be interpreted either as a small spike right before the i_0 measurement, or a 'large' one right after the $i_0 - 1$ measurement. Thus, using a constant λ would systematically select the small spike interpretation, *i.e.* the inferred spike train would be systematically delayed with respect to the original spike train. This behavior is not desirable, and we would rather have both solutions to be degenerate global optima. This can be achieved by setting λ to a smaller value right after the $i_0 - 1$ measurement. We show in Annex E that both efficient noise filtering and unbiased estimation of spike timing for isolated spikes can be obtained with the following expression for λ_{BSD}^j :

$$\lambda_{BSD}^j = z_1 \|K_j\| \min \left(\sigma, \frac{a \sum_{r=1}^s \|K_r\|}{z_1 + z_2} \right) \quad (42)$$

In practice, the optimization can also be performed efficiently using the interior-point method [8]. Adding a small L_2 penalty $\sum_j \mu_j N_j^2$, with $\mu_j \propto \|K_j\|^2$ often provides better conditioning of the hessian, and faster convergence. It also ensures the unicity of the solution, in particular when $\tau_r = 0$. The kernel inference can also be adapted efficiently to support super-resolution, see Annex F.

G. Methods for Performance Evaluation

The algorithm is evaluated based on three criteria: whether the spikes are correctly detected or not; if they are detected, whether their timing is accurate; and whether the inferred generative model parameters match the ground truth values.

1. Spike detection

For the spike detection assessment, we adopt the SpikeFinder main metric [13], i.e. the Pearson correlation between the ground-truth spike train discretized at a frequency f_{eval} and the inferred spike train resampled at the same frequency. Resampling is performed as follows: (i) if $f_{\text{eval}} = f/s$ for integer s , s consecutive frames are summed; (ii) if $f_{\text{eval}} = f \times s$ for integer s , we write $N'_j = N_{\lceil j/s \rceil}/s$. (iii) otherwise, (e.g. for SpikeFinder), the signal is first resampled to the closest multiplier or divider of f_{eval} by linear interpolation, then resampled using (i) or (ii). When $f_{\text{eval}} = f$, the metric penalizes equally an undetected spike and a spike that is detected either in advance or delayed with respect to the true spike; when $f_{\text{eval}} < f$, the metric is tolerant to timing errors of order $1/f_{\text{eval}}$.

2. Spike timing

The accuracy of the timing is assessed by measuring the point spread function (PSF), *i.e.* the average (over noise) of the inferred spike train in presence of an isolated single spike:

$$PSF_\tau = \mathbb{E} \left[\hat{N}_{t_0+\tau} | N_t^0 = \delta_{t,t_0} \right] \quad (43)$$

Where $\tau \in \mathbb{Z}$ is the lag. Ideal point spread functions are centered around 0 and decay fast to ~ 0 with $|\tau|$. [48]. Then, the timing error can be characterized by the mean and standard deviation of a Gaussian curve fitted the point spread function. [49]

A naive estimator of PSF would be the cross-correlation between the ground-truth spikes and the inferred spikes $X_\tau = \frac{1}{T} \sum_{i=1}^T N_i^0 \hat{N}_{i+\tau}$. Such estimator is correct when the spike train are Poisson-distributed, but is biased when the spikes are temporally correlated: in the best case scenario where the spike is perfectly recovered $\hat{N}_i = N_i^0 \forall i$, we obtain the auto-correlation of the spike trains $X_\tau = \frac{1}{T} \sum_{i=1}^T N_i^0 N_{i+\tau}^0 \equiv A_\tau$. To overcome this issue, we estimate the PSF as the kernel of a linear temporal regression model:

$$\hat{N}_t = \sum_{\tau=-m}^m PSF_\tau N^0(t-\tau) + \epsilon \quad (44)$$

The least square estimator for the PSF is then:

$$\begin{aligned} \mathbf{PSF} &= \mathcal{A}^{-1} \mathbf{X} \\ \mathcal{A}_{ij} &= A(i-j) \end{aligned} \quad (45)$$

where PSF_τ , X_τ are indexed formally as vectors \mathbf{R} , \mathbf{X} .

This estimator is defined for N^0 , \hat{N} with identical sampling rate; if \hat{N} has lower sampling rate, it is first upsampled using the above procedure.

3. Generative model parameters

Finally, we assess the accuracy of the blind generative model parameters inference. For simulated train spikes, ground truth generative model parameters exist and they can be directly compared with the inferred ones. For real fluorescence data with joint electrophysiological recordings, we derive 'ground truth' generative parameters using the knowledge of the spike positions. The kernel parameters and the baseline are obtained by minimizing the following mean square error:

$$MSE = \sum_i \left(F_i - \sum_j K[(i-j+1)\Delta t] a_j N_j^0 - b \right)^2$$

$$(\tau_r^{GT}, \tau_d^{GT}, b^{GT}) = \arg \min_{\tau_r, \tau_d, b, \mathbf{a} \geq 0} MSE(\tau_r, \tau_d, b, \mathbf{a}) \quad (46)$$

Where the ground-truth spikes are discretized at the fluorescence sampling frequency. Note that we relax the hypothesis that all transients have the same amplitude a , and optimize over all the amplitudes $a_j \geq 0$. This is particularly important for spike bursts, where strong non-linear effects are observed. In practice, the optimization is carried out by using the exact same algorithm as for the fully blind setting, but with a position-dependent sparsity prior λ_i that takes the value 0 for positions where a ground truth spike is present and a large value (e.g. 20) elsewhere. Once inference is performed, the transient amplitude a is computed as the average transient amplitude $(\sum_j a_j) / (\sum_i N_i^0)$.

II. RESULTS

The Blind Sparse Deconvolution (BSD) method, whose algorithmic details were presented in the preceding section, allows for unsupervised spike inference, *i.e.* both the algorithm hyperparameters - sparsity prior λ and generative model parameters - kernel time constants τ_r, τ_d , transient amplitude a , noise level σ - are automatically evaluated. BSD can infer spike trains at or beyond the fluorescence sampling rate. Importantly, the expected performances of BSD can be predicted as well and this possibility is integrated in the publicly available program.

Section II A is dedicated to simulated data; we demonstrate that the choice of sparsity prior outperforms other methods in terms of spike and/or computational speed and that the kernel parameters can be accurately recovered using our iterative approach. We also demonstrate super-resolution capabilities for a wide range of experimental conditions. In Section II C, we apply BSD to the SpikeFinder contest, a collection of joint electrophysiological and fluorescence recordings. We show that i) our choice of sparsity prior outperforms others ii) the kernel inference is accurate, allows robust performance even in the absence of training data and can improve spike detection performance when interneuron variability is important. iii) the predicted temporal errors are consistent with empirical errors and that integrating it into the prediction can improve spike detection. iv) super-resolution significantly increases the temporal accuracy for some datasets. Overall, our best submission is competitive across all datasets with state-of-the-art Machine Learning algorithms, while not requiring any training data or hyperparameter fine-tuning. In Section II D, we show that BSD scales well to large-scale zebrafish recordings. Finally, section II B is dedicated to experimental design: we use BSD to predict the expected accuracy for various standard calcium reporters and imaging parameters.

A. Simulated data

1. Spike detection accuracy

In BSD, the sparsity prior λ_{BSD} is computed analytically and allows one to simultaneously minimize, in a tractable way, both the false-positive rate (FPR) and false-negative rate (FNR) (see Methods). In contrast, the expression λ_{oopsi} used in the fast-oopsi algorithm [8] offers no guarantee that either is small in all situations (see Methods). This issue has motivated the recent development of the constrained-oopsi algorithm [16] where the sparsity prior is determined iteratively. Figure 1 illustrates the strong impact of the chosen value of the sparsity prior on the inference performance, as it compares the results of the four inference algorithms (BSD, oopsi, con-oopsi, and non-negative, *i.e.* with $\lambda = 0$) for a signal with $\sigma = 0.4$, $f = 10Hz$. For fast-oopsi, the sparsity prior is too large, and no spikes are inferred, whereas for both con-oopsi and BSD, the signal is correctly recovered. Here, BSD infers slightly less false

spikes than con-oopsi in this particular configuration. As a baseline, we also show the non-negative deconvolution without any sparsity prior, which is significantly better than naive deconvolution but shows significantly more false positives than con-oopsi and BSD. One may notice that the BSD’s reconstructed signal is systematically lower than the original signal. This difference directly derives from Eqn. 17, which indicates that a spike of amplitude a is reconstructed with an amplitude $a - \frac{\lambda}{\|K\|^2}$. However, this systematic bias in the reconstructed signal does not impact the quality of the inferred spike signal.

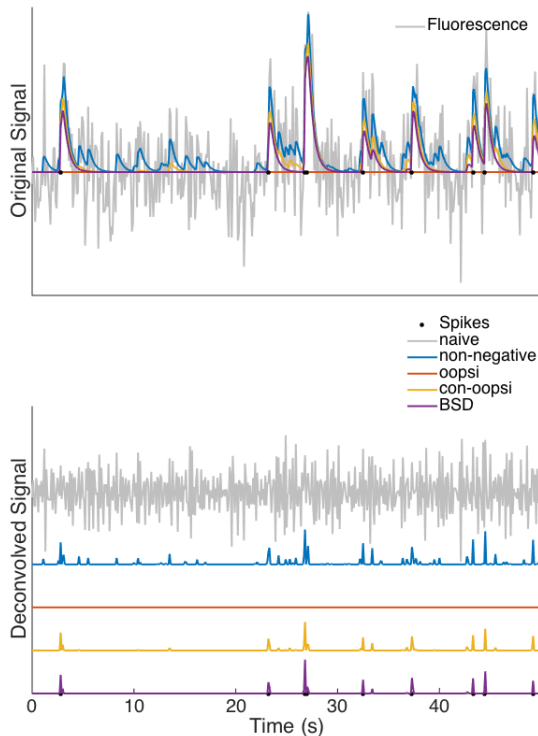


FIG. 1: **Example of fluorescence signal and inference results for various deconvolution frameworks.** The fluorescence signal was generated using parameters $f = 10Hz$, $\sigma = 0.4$, $a = 1$, $\tau_r = 0.1$, $\tau_d = 0.5$. The naive inference corresponds to the result of Eqn. 5.

The performances of the four algorithms are now compared on a systematic benchmark. A random spike train is drawn from a Poisson distribution of mean firing rate $\nu = 0.1Hz$ over a duration $t = 10000s$. The signal is generated according to the discrete model Eqn. 3 with a fixed transient amplitude $a = 1$ and variable sampling frequency f and noise level σ . We use a double exponential kernel K with $\tau_r = 0.1$, $\tau_d = 0.5$ (see Eqn. 2), akin to a GCaMP6 reporter. Spike trains $\hat{\mathbf{N}}$ are inferred *with knowledge of the generative model parameters* and compared with the original spike train \mathbf{N} using the metric defined in Section I G 1. We show in Fig. 2 the correlation as function of the signal-to-noise ratio for various sampling frequencies, and at evaluation frequencies $f_{eval} = f$ (top row) and $f_{eval} = 10Hz$ (bottom row).

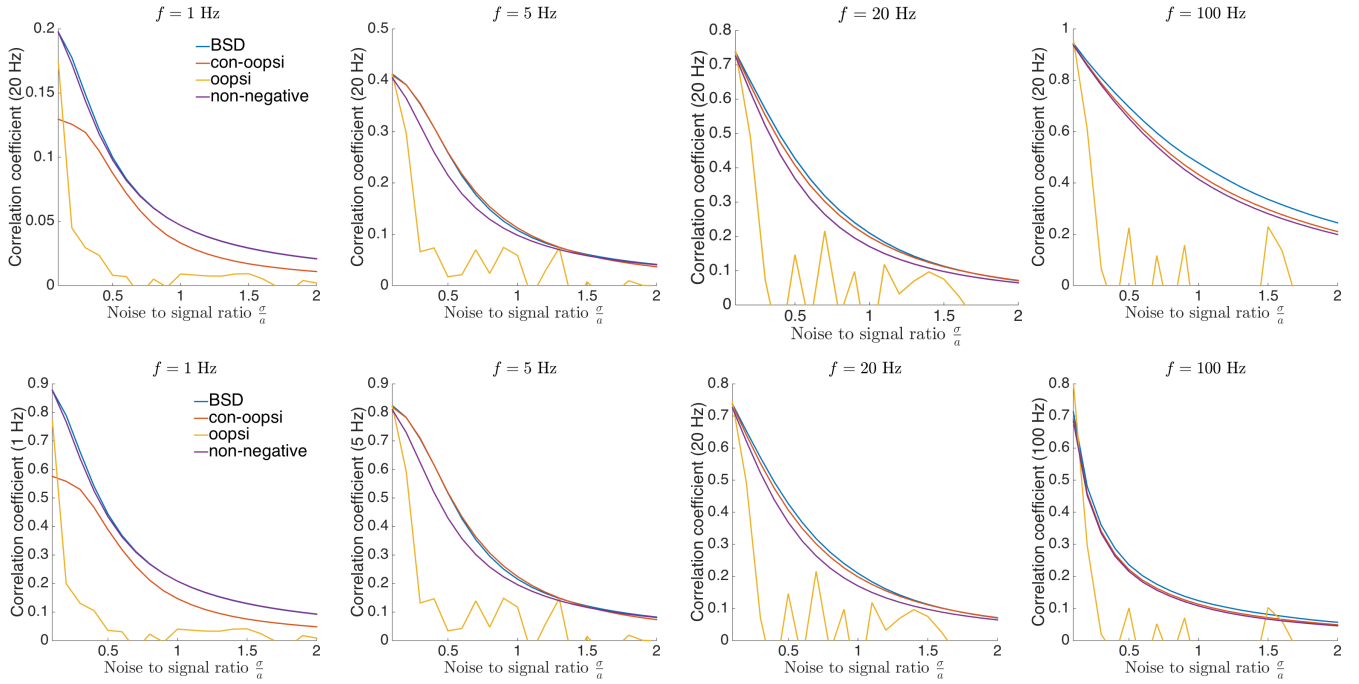


FIG. 2: Comparison of the reconstruction performances on synthetic data between BSD, con-oopsi, oopsi, and non-negative deconvolution. For each algorithm, the correlation between the inferred and ground truth spike train is shown as a function of the noise-to-signal ratio. The different plots correspond to various sampling frequencies f and evaluation frequency (indicated on the y axis legend).

We observe that the fast-oopsi algorithm sometimes performs very well ($f = 100\text{Hz}$, high SNR) sometimes equivalently as con-oopsi and BSD ($f = 100\text{Hz}$, lower SNR), but often very poorly ($f = 20\text{Hz}$, $\sigma > 0.2$); Such unreliability may be highly detrimental in actual experiments. BSD and con-oopsi yield comparable results, with BSD slightly outperforming con-oopsi in most configurations and particularly at low sampling rate ($f = 1\text{Hz}$). The non-negative deconvolution without sparsity prior baseline is more stable than oopsi (as was previously reported in [29]) and performs equivalently as BSD and con-oopsi at high SNR, but, as expected, has a significantly higher false positive rate than BSD and con-oopsi at low SNR. Similar results are found when increasing the firing rate (see same simulations for $\nu = 1\text{Hz}$, Supplementary Figure 11); then differences tend to vanish at large firing rates (same experiment for $\nu = 5\text{Hz}$, Supplementary Figures 11).

We also compare the computational cost of the various algorithms. BSD, non-negative deconvolution and fast-oopsi all share the same core algorithm and therefore have similar computational cost. In contrast, the con-oopsi implementation is slower because the sparse deconvolution has to be performed many times with different values of λ until convergence is reached. In practice, for experiments performed on a MacBook Air 2013, with 1.3 GHz Intel Core i5, we find a 3 to 25-fold increase in computation speed, depending on the array size. Our experiments shows that the number of iterations can be surprisingly large in practice. In particular, if the noise level is underestimated by con-oopsi, the error constraint is tighter and adding the positivity constraint may lead to no solutions at all - yielding many iterations in vain and increased computational time, see Table I. This reflects in the fact that the computing time is largely dependent on whether or not the noise level is provided.

Notice that the exact gain in speed depends on which version of con-oopsi is used (here, Matlab implementation, con-oopsi version of Dec. 2015, with cvx). Although we did not test the PAVA optimizer [17], we expect a gain of the same order of magnitude between constrained-PAVA and BSD-like PAVA. Such a difference in computation load may prove highly beneficial for real-time inference in high data-throughput recordings, as illustrated in Section II-D.

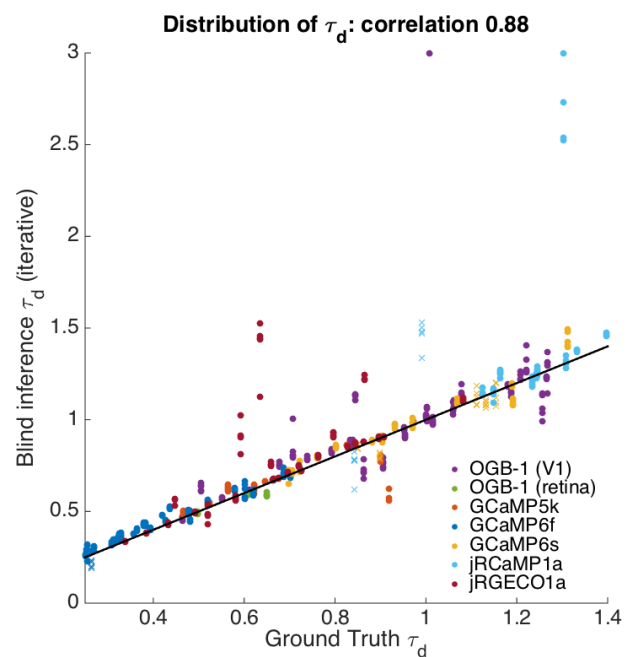
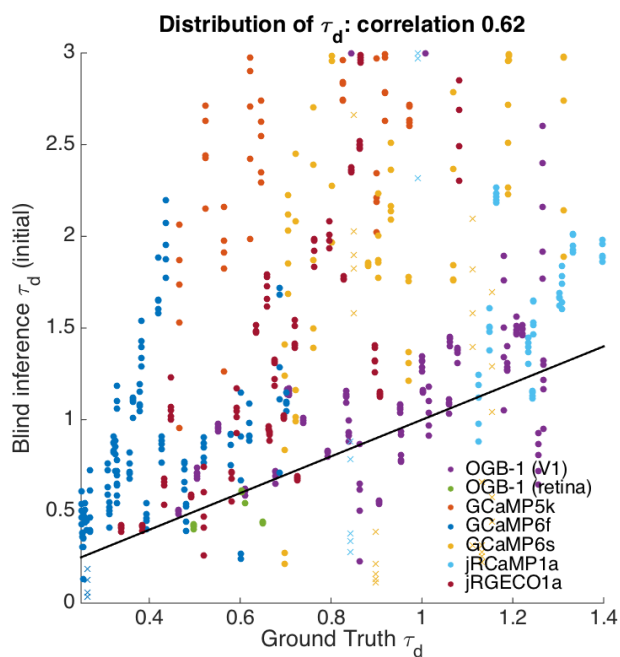
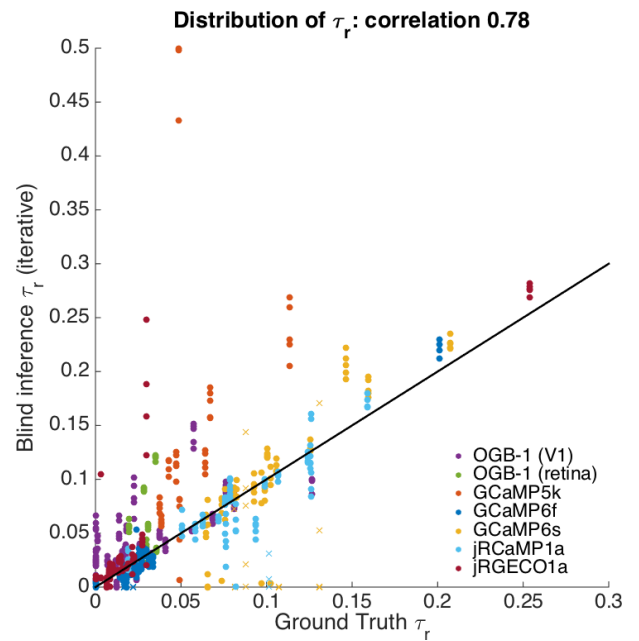
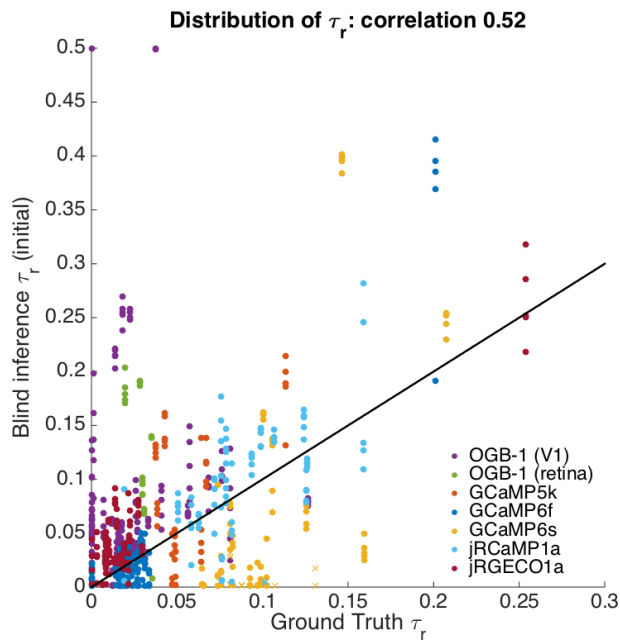
N_{frames}	BSD/ fast-oopsi (s)	con-oopsi (s)	con-oopsi σ user-provided
10^4	0.7	2.4	2.1
$5 \cdot 10^4$	2.6	8.2	8.3
$2 \cdot 10^5$	10	84	45
$5 \cdot 10^5$	27	529	128
10^6	49	1594	290

TABLE I: Comparison of BSD, fast-oopsi and con-oopsi computational speed. For con-oopsi, under-estimation of the noise level σ , even for synthetic data, can lead to a large increase in computational time

2. Kernel inference accuracy

When deconvolution is performed with an incorrect kernel, systematic biases arise in the inferred spike trains and the accuracy decreases both in terms of spike timings and spike detections: for instance, spikes may be split into two time frames if the rise time is too short, or some spikes may be missed if the decay time is too long, see Annex B and Figure 10 for a detailed study. We now relax the assumption that the generative model parameters are known and attempt to retrieve them from the raw fluorescence recordings. For long recordings of Poisson-distributed spike trains with sufficiently high signal-to-noise ratio and purely white Gaussian noise, the task is relatively easy. Indeed, the autocorrelation function follows exactly Eqn. 31, and the initial kernel estimation is excellent. However, most real datasets feature bursts of spikes, temporally correlated noise, artifacts, limited spike counts and low signal-to-noise ratio. Under these conditions, Eqn. 31 becomes incorrect.

We thus build a more realistic collection of synthetic fluorescence traces as follows: for each recording of the SpikeFinder dataset (see next section), we generate five synthetic fluorescence recordings according to the generative model of Eqn. 1 using the actual spike train measured by electrophysiology and the rise and decay time constants and signal-to-noise ratio inferred from the recording (see next section). This synthetic dataset has diverse spike counts, signal-to-noise ratios, sampling rates, kernel parameters and spiking patterns. We jointly infer the spike trains and kernel parameters using either BSD or adaptive BSD (iterative parameter estimation). Results are shown in Figure 3. Panels (a) and (c) show the outcome of the initial kernel estimation, compared with the ground truth kernel values. We find that the decay time is often vastly overestimated. This systematic bias reflects the existence of large temporal correlations in the experimental spike signal whereas the inference model assumes none. The algorithm thus tends to confuse spike bursts with long fluorescence transients, and the kernels are therefore poorly inferred at first. However, after kernel refinement (Panels (b) and (d)), the inferred kernel becomes very similar to its experimental counterpart in most experiments. Upon closer inspection, we can identify the main sources of error for the algorithm. Coarse errors arise when the algorithm starts from largely incorrect initial kernel values and ends up trapped in a wrong local minimum of the cost function. These errors can be easily avoided by providing realistic bounds for the rise and decay time constants. For all the remaining cases, three main factors affect the quality of inference parameter, see Supplementary Figure 13: the spike count (higher is better), the effective signal-to-noise ratio (higher is better) and the spike 'burstiness', i.e. the deviations from Poisson-distributed spikes (lower is better).



Kernel inference:
BSD

Kernel inference:
Adaptive BSD

FIG. 3: Kernel inference benchmark for realistic simulated data For each recording of the SpikeFinder datasets (see Section II C), we generate five simulated fluorescence traces according to Eqn. 3 using the actual spike trains, and signal-to-noise ratios and kernel parameters estimated, see Section II C for derivation of the generative model parameters. Then, BSD and adaptive BSD are applied to infer the spike trains and the kernel parameters. Panels (a),(b): ground-truth rise and decay time constants vs inferred values using BSD. Panels (c),(d): Same, with adaptive BSD. In all panels, crosses denote spike recordings with fewer than 20 spikes.

3. Super-resolution

Figure 4a shows an example of reconstruction of a signal generated at $f_0 = 20\text{Hz}$, and sampled at 5Hz . We observe a good agreement with the original spike train. In particular, it appears that in spite of the sparse sampling, the onsets of the green and dark curves transients are very close to one another.

We test our algorithm on synthetic datasets generated using the model Eqn. 39 at $f_0 = 500\text{Hz}$, with $\tau_r = 0.1$, $\tau_d = 0.5$, spike frequency $\nu = 2\text{Hz}$. The fluorescence signal is down sampled to recording frequencies ranging from $f = 1\text{ Hz}$ to 500Hz . Spike trains are inferred with and without super-resolution. For super-resolution, we use a frequency gain $s = \frac{f_0}{f}$ in order to reconstruct a spike train at the original frequency f_0 . We perform the spike inference for various sampling frequencies and noise levels, and we estimate the point spread function of the inferred spike train (*i.e.* the average response to a single spike, see the method in Section I-G2). Results are depicted in Figure 4 (b) and (c). They demonstrate that super-resolution is perfectly workable at small noise levels, and that a significant resolution gain can be achieved at intermediate noise level typical of actual experimental conditions. For instance, at $f = 10\text{ Hz}$, $\text{SNR} = 5$, the point spread function width is $\sim 2\times$ smaller than without super-resolution. Figure 4 (c) shows that the gain in resolution becomes significant as soon as $f \gtrsim 4\text{Hz}$.

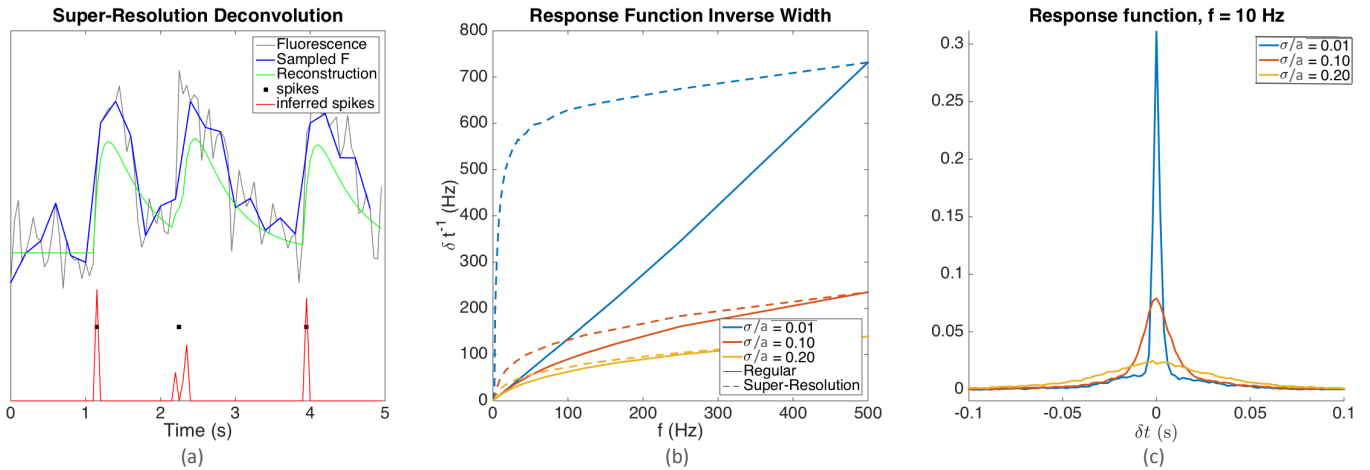


FIG. 4: **Super-resolution for synthetic data** (a) Example of super-resolution inference: a fluorescence signal is generated at $f_0 = 20\text{ Hz}$, sampled at 5 Hz and reconstructed at 20 Hz . Parameters: $\tau_r = 0.1$, $\tau_d = 0.5$, $\sigma = 0.2$, $a = 1$. (b) point spread function at 10Hz for various noises. Notice a smaller width than the sampling interval 0.1s (c) Inverse width δt^{-1} of the point-spread function, as function of the sampling frequency for regular reconstruction (full) and SR reconstruction (dotted).

B. Theoretical limits of calcium reporter accuracy

The fundamental motivations for spike train inference are to denoise the fluorescence signal and to improve the temporal resolution of the neural recording. Theoretically, if the generative model is correct, the convolution kernel is known and the signal is noiseless, then perfect retrieval of the spike train in terms of detection and timing can be achieved. Because of the noise, the accuracy is in practice limited by the rise and decay times: some spikes can be missed, have a wrong timing or be split across two successive time bins (see Annex B and Fig 10). These limitations have been characterized quantitatively in [30] in the context of Bayesian inference, when the noise is Poisson-like, τ_r is negligible and without super-resolution. However, no such analysis has been performed for sparse deconvolution algorithms.

We derived in Section ID theoretical bounds of performance of our deconvolution algorithm; we present hereafter the main findings. We display in Figure 5a the true-positive rate (TPR) for different sampling rates, as a function of the noise level for $\tau_r = 0.1$, $\tau_d = 0.5$. The false-positive rate (FPR) is set at $0.01/\text{frame}$ (*i.e.* $\lambda_{\text{BSD}} = \lambda_1$ and $z_1 = 2.366$, see Methods). An important insight of this graph is that at low sampling rate, the TPR quickly decays with the noise level because spikes emitted shortly after a measurement are often completely missed. Conversely, improving the TPR (with *e.g.* $z_2 = 2.366$) yields a large number of false positives at low sampling rate. Some

point spread functions are displayed in Figure 5b-inset for typical calcium indicators. We used the parameters: (i) GCaMP6s: $\tau_r = 180ms$, $\tau_d = 0.55$, (ii) GCaMP5k: $\tau_r = 58ms$, $\tau_d = 0.52s$, (iii) GCaMP6f: $\tau_r = 25ms$, $\tau_d = 0.38s$, (iv) OGB1-like: $\tau_r = 20ms$, $\tau_d = 80ms$. For the 3 first sets of time constants, the values are deduced from the fluorescence recordings on mice V1 cells reported in [31].

We also display in Figure 5b the width δt of the point spread function (extracted using a gaussian fit) as a function of the noise level, for a sampling frequency $f = 60Hz$. As expected, the temporal resolution of the spikes can be lower than the sampling period if the noise is large, and we observe that reporters with large τ_r yield lower temporal resolution.

We finally examine the impact of the sampling frequency on the temporal resolution. In an experiment with a fixed number of sampled neurons, increasing the sampling rate $f \equiv \Delta t^{-1}$ by a factor s typically comes at the cost of reducing the exposure time τ_e by the same factor s , which in turn increases the noise σ by \sqrt{s} . Therefore, there is no guarantee that increasing the sampling frequency improves the time-resolution. We display in Figure 5c for the same set of calcium indicators, the inverse width δt^{-1} as a function of the sampling rate f , for various signal-to-noise ratios (SNRs) at a reference frequency of $10Hz$. We see that δt^{-1} saturates at a value that depends on the SNR and on the rise and decay constant times. For instance, with GCaMP6s and $SNR_{10Hz} = 5$, increasing the frequency beyond $50Hz$ does not result in improving the temporal resolution.

The observation that the temporal resolution saturates at high sampling rate can be understood by examining Eqn. 27: asymptotically, we have $\|K\| \propto \sqrt{\Delta t}$, and since $\sigma \propto \sqrt{\Delta t}$, the effective noise level $\frac{\sigma}{a\|K\|}$ reaches a well-defined limit - and so does δt .

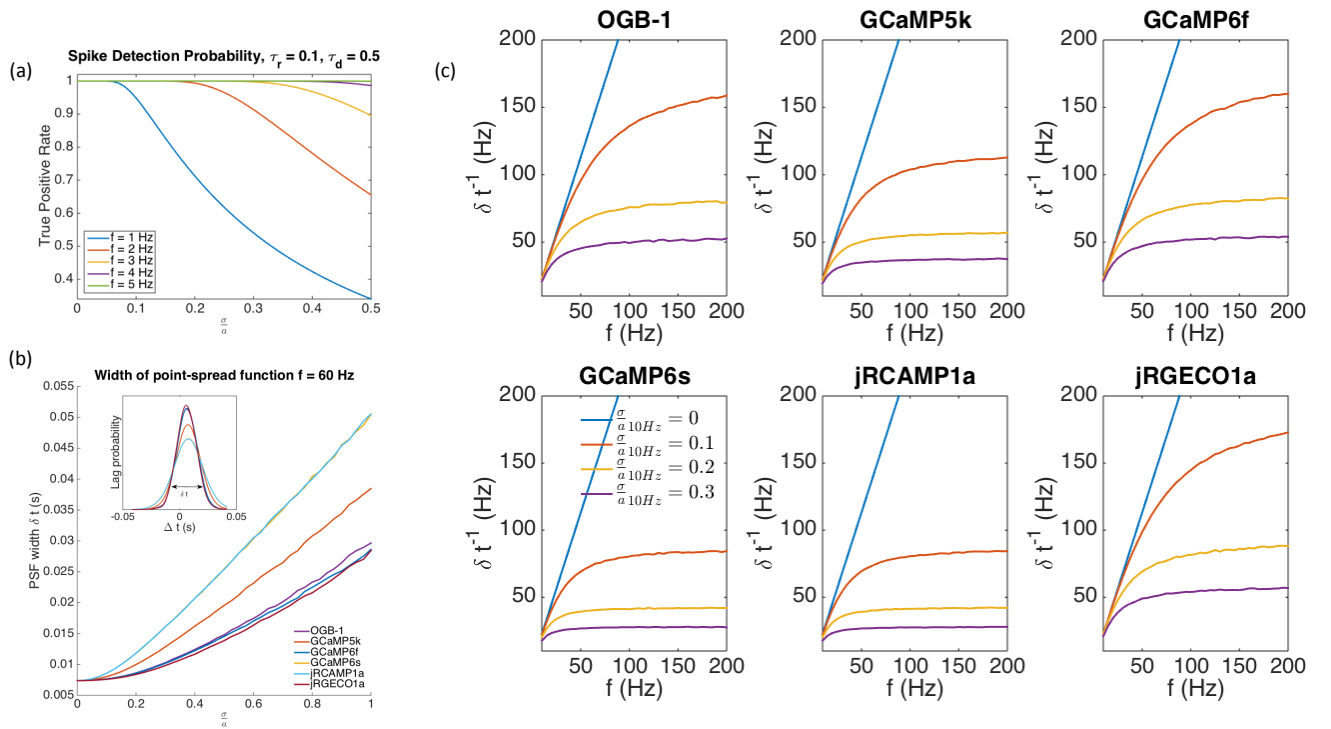


FIG. 5: **Theoretical limits of calcium reporters** (a) True positive rate of BSD as a function of the noise level for different sampling rates, for a fixed FPR of 0.01/frame, (b) Width of the point spread function as function of the noise for various calcium indicators, at fixed frequency $f = 60Hz$, (c) Width of the point spread function as a function of the sampling frequency for various calcium indicators and reference noises.

C. Joint electrophysiological and fluorescence recordings: the SpikeFinder contest

SpikeFinder (<http://spikefinder.codeneuro.org/>) [32] is a public contest for spikes train inference from calcium recordings. It consists of a compilation of 10 datasets of joint electrophysiological and calcium recordings from

mouse V1 or retina, for six different calcium probes (OGB-1, GCaMP5k, GCaMP6s, GCaMP6f, jRCaMP1a, jRGECO1a) at various sampling rate, see Table II. The data was compiled from various sources [14, 31, 33, 34]. All the recordings are upsampled to $f = 100Hz$, and the metric used to assess performance is the correlation between the inferred spike train and the electrophysiological spike train, computed after downsampling at $20Hz$. Dozens of algorithms based on Supervised Machine Learning, generative models and others have been benchmarked on SpikeFinder, and the results are publicly available.

Dataset	1	2	3	4	5	6	7	8	9	10
Calcium Indicator	OGB-1	OGB-1	GCaMP6s	OGB-1	GCaMP6s	GCaMP5k	GCaMP6f	GCaMP6s	jRCaMP1a	jRGECO1a
Brain region	V1	V1	V1	Retina	V1	V1	V1	V1	V1	V1
Number of recordings	11	21	13	6	9	9	37	21	20	27
Spikes per recording	1012	525	1240	1304	872	363	127	96	75	232
Sampling rate (Hz)	40	12	60	8	60	50	60	60	15	30
Rise time (s)	0.02	0.02	0.00*	0.03	0.00*	0.05	0.02	0.10	0.09	0.01
Decay Time (s)	0.95	1.00	1.10	0.55	0.99	0.65	0.33	0.97	1.32	0.68
Noise level $\frac{\sigma}{\mu}$	1.06	0.68	0.83	1.06	0.90	0.41	0.26	0.21	0.39	0.29

TABLE II: SpikeFinder data sets summary. For recording-specific variables, the median value across recordings was provided. The rise and decay time constants and signal-to-noise ratios, as defined by Eqn. 2, are obtained by a parametric linear regression of the true spikes against the fluorescence, see Section IG 3. * The rise time for datasets 3 and 5 are unreliable because of important delays between the fluorescence and electrophysiological recordings

We benchmark several inference algorithms and postprocessing variants on the SpikeFinder dataset using the following pipeline:

- **Preprocessing.** For all datasets, we started from the raw data, before upsampling or detrending. Indeed, when the fluorescence is upsampled, the noise becomes temporally correlated between time frames and the white noise assumption used to compute the sparsity prior λ is violated. Then, each fluorescence recording is normalized (zero mean and unit variance). The variable component of the baseline is further removed by subtracting a moving percentile (quantile $q = 0.15$, variable window size $\in [10s, 60s]$ adjusted by validation). Such baseline slow modulations are unrelated to the calcium signal but reflect experimental artefact such as photobleaching or minute axial motions of the specimen.
- **Initial values of rise and decay time.** We tested four initialization procedures: **Blind**: no information is provided and the automated initial kernel estimation described in Section IE 1 is used; **Literature** parameters are taken from the literature, see values and sources in Supplementary Table VI; **Ground Truth** for each recording, we compute ground truth rise and decay time constants by temporal regression of the fluorescence from the true spikes, see Section IG 3. **Train Set** for each dataset, we use the median values of the ground truth rise and decay time as initial values, see values in Table II. We also use the minimum and maximum values found in the dataset as bounds for adaptive BSD. Note that the ground-truth initialization is not a valid algorithm as it requires prior knowledge of the spikes for each neuron. We thus only use it as a reference for validating the kernel inference and as an upper bound of performance.
- **Inference Algorithm** We tested four sparse deconvolution variants: non-negative deconvolution (*i.e.* with $\lambda = 0$), con-oopsi, BSD and adaptive BSD.
- **Postprocessing.** The inferred signal is upsampled to $100Hz$ by linear interpolation, and we add a temporal offset (adjusted on the training set). it is mostly relevant for datasets 3 and 5 to account for a delay between the fluorescence time and electrophysiological time, in agreement with other submissions [32]. Then, the inferred spike is optionally convolved with the point-spread function computed in Section ID 2, using the parameters inferred from the recording. This can be viewed as a poor man (but computationally efficient) version of posterior averaging: one first computes the most likely spike train given the fluorescence using BSD, and then convolve it with the point-spread function (PSF) to account for the temporal uncertainty, see examples in Fig. 6

Quantitative results are displayed in Table III and selected examples are shown in Figure 6.

Dataset	1	2	3	4	5	6	7	8	9	10	Mean	
Calcium Indicator	OGB-1	OGB-1	GCaMP6s	OGB-1	GCaMP6s	GCaMP5k	GCaMP6f	GCaMP6s	jRCAMP1a	jRGECO1a	Score	
Blind	non-negative	0.378	0.439	0.223	0.416	0.189	0.439	0.649	0.610	0.579	0.760	0.468
	con-oopsi	0.370	0.452	0.253	0.551	0.213	0.444	0.659	0.619	0.586	0.728	0.487
	BSD	0.425	0.439	0.277	0.545	0.279	0.533	0.707	0.595	0.583	0.747	0.513
	adaptive BSD	0.451	0.450	0.294	0.542	0.332	0.621	0.747	0.686	0.625	0.835	0.558
	non-negative * PSF	0.405	0.457	0.333	0.442	0.291	0.464	0.631	0.547	0.526	0.747	0.484
	BSD * PSF	0.450	0.466	0.410	0.544	0.344	0.554	0.681	0.575	0.553	0.758	0.533
	adaptive BSD * PSF	0.487	0.466	0.424	0.537	0.413	0.613	0.740	0.674	0.559	0.820	0.573
Literature	non-negative	0.372	0.391	0.288	0.419	0.312	0.562	0.605	0.646	0.580	0.760	0.493
	con-oopsi	0.372	0.427	0.288	0.541	0.312	0.562	0.627	0.651	0.585	0.794	0.516
	BSD	0.396	0.435	0.291	0.525	0.305	0.565	0.625	0.649	0.601	0.800	0.519
	adaptive BSD	0.451	0.459	0.453	0.541	0.420	0.611	0.733	0.643	0.613	0.837	0.576
	non-negative * PSF	0.421	0.398	0.396	0.451	0.411	0.547	0.584	0.601	0.537	0.733	0.508
	BSD * PSF	0.440	0.439	0.368	0.529	0.391	0.549	0.628	0.641	0.572	0.788	0.535
	adaptive BSD * PSF	0.486	0.473	0.493	0.546	0.497	0.600	0.725	0.620	0.577	0.823	0.584
Train Set	non-negative	0.456	0.457	0.441	0.442	0.403	0.600	0.701	0.654	0.588	0.794	0.554
	con-oopsi	0.441	0.458	0.440	0.550	0.405	0.600	0.715	0.654	0.588	0.811	0.564
	BSD	0.446	0.458	0.449	0.539	0.398	0.598	0.734	0.641	0.611	0.803	0.568
	adaptive BSD	0.446	0.459	0.453	0.540	0.420	0.621	0.737	0.705	0.618	0.835	0.583
	non-negative * PSF	0.488	0.467	0.485	0.473	0.479	0.598	0.681	0.641	0.536	0.777	0.563
	BSD * PSF	0.497	0.470	0.491	0.543	0.475	0.599	0.717	0.645	0.582	0.802	0.582
	adaptive BSD * PSF	0.495	0.475	0.493	0.546	0.498	0.613	0.730	0.684	0.574	0.817	0.593
Ground Truth	non-negative	0.461	0.460	0.270	0.444	0.395	0.613	0.723	0.681	0.592	0.777	0.542
	con-oopsi	0.461	0.462	0.270	0.550	0.395	0.613	0.741	0.681	0.592	0.789	0.555
	BSD	0.446	0.456	0.269	0.539	0.390	0.610	0.747	0.671	0.609	0.790	0.553

TABLE III: Correlation scores for the SpikeFinder data (train sets), for various choices of inference algorithms, postprocessings (convolved with the Point-Spread function or not) and initial kernel values: Blind = no information at all; Literature = using parameters derived from the literature, see Supplementary Table VI; Train set = using parameter derived from the train set, see Table II; Ground-truth = using different parameters for each neuron, derived using the knowledge of the spike positions. The best score for each dataset and configuration is indicated in bold

Sparsity prior. As expected, introducing a sparsity prior reduces the false detection rate (see Figure 6) and significantly improves the correlation regardless of the kernel choice. At fixed kernel parameters, BSD and con-oopsi are virtually equivalent but the former is faster.

Kernel Inference. When no training data is available, adaptive kernel inference is critical for performance as was found in the simulated data. Even when training data is available, adaptive kernel inference is always equivalent or better to using fixed kernels derived from training data as it takes into account inter-neuron variability. The most important gains is found for dataset 8, where non-linear effects result in large variability of decay rates. The quantitative results are corroborated by Figure 7, which compares the rise and decay time constants and transient amplitude a inferred by the blind algorithm with their ground truth counterparts. Adaptive BSD successfully captures the rise and decay time constants variability both across experiments, and across neurons in a given experiment, notably for dataset 8 where the decay time can vary by more than two-fold. When the training data is of low quality as for data sets 3 and 5, blind inference actually outperforms the ground-truth. Importantly, although our best submission was obtained with the train set initial value, the blind and literature initialization follow closely: BSD almost does not require any training data. We note however that the correlations are slightly smaller than the values found using synthetic data with similar parameters (see Figure 3). This can be explained by artifacts, temporally correlated noise and non-linear effects. The inferred transient amplitudes also correlate with the ground-truth values, but less accurately. For the OGB-1 datasets, which have low signal-to-noise ratio, a is systematically overestimated because the algorithm frequently misses isolated spikes and can only detect bursts of consecutive spikes, see Figure 6 panel (c); it therefore confuses the latter with individual spikes. Such confusion is expected and inevitable: for instance, in dataset 4, only bursts of more than 10 spikes produce visually detectable transients. For dataset 8 (GCaMP6s), we find conversely that some transients amplitudes are largely underestimated, because the algorithm infers many small (but above threshold) spurious spikes due to artifacts or non-linear effects, and confuses them with real spikes.

Temporal resolution. We assess the temporal accuracy of the algorithms as follows: for each recording, we computed the empirical point-spread function (PSF) between the ground-truth spikes and the inferred spikes using Eqn. 45. An average PSF is then computed for each dataset by weighted average of the PSFs (normalized to max = 1,

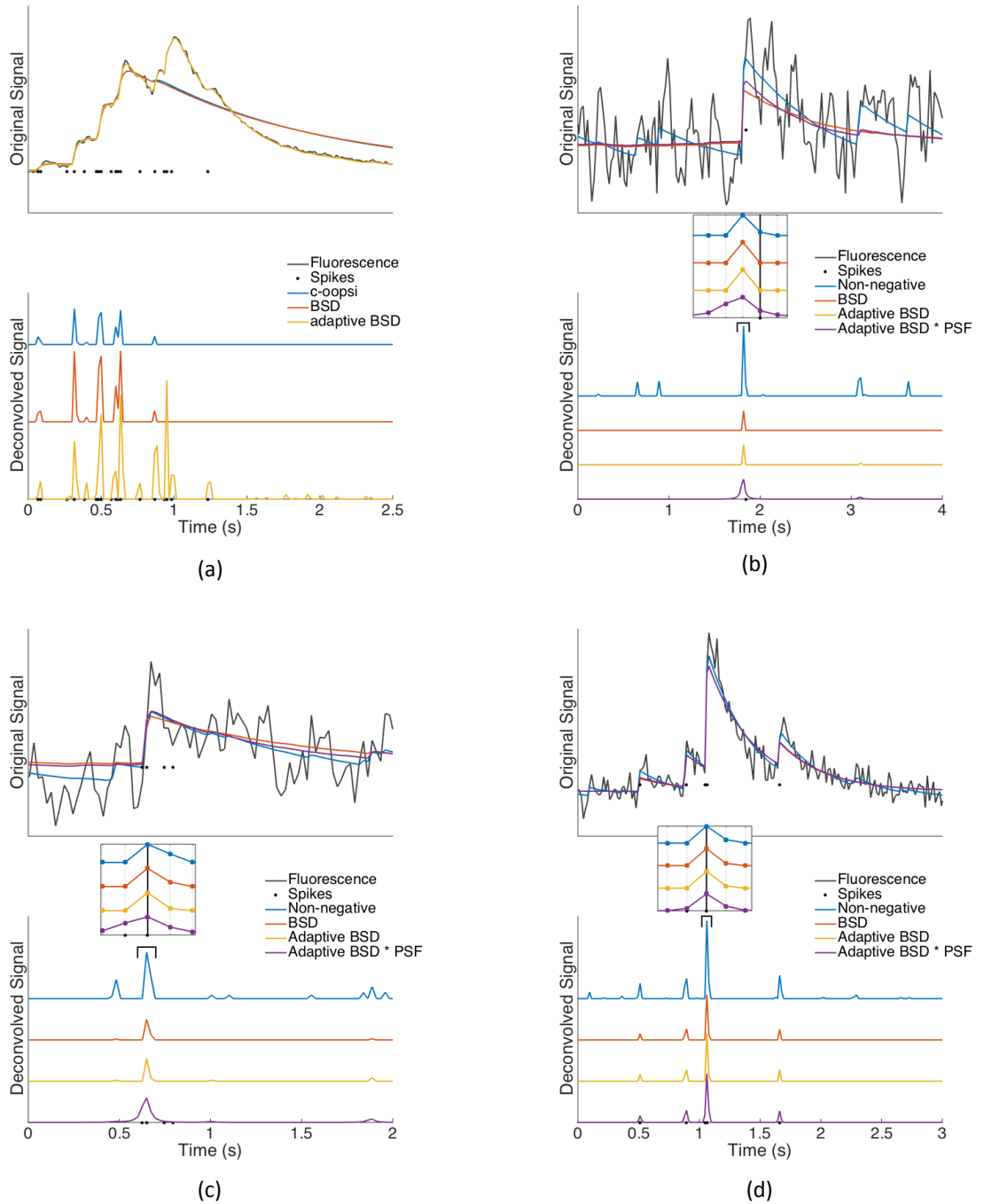


FIG. 6: **Selected examples from the SpikeFinder datasets** Four fluorescence traces from the SpikeFinder dataset and the corresponding inferred and ground truth spike trains for various algorithms. (a): GCaMP6s, $f = 60Hz$ (dataset 8), blind initial kernel. (b), (c): OGB-1 $f = 40Hz$ (dataset 1), train set initial kernel (d) GCaMP6f $f = 60Hz$ (dataset 7), train set initial kernel. For panels (b),(c),(d), insets show the same inferred spike trains, zoomed-in on a spike to illustrate timing errors.

then weighted by the number of spikes). We also predicted for each dataset a theoretical Point-Spread Function, using Eqn. 29 and either the parameters of Table II or the parameters inferred by BSD for each recording. Results are shown in Table IV and Figure 8. Empirical PSF and the fluorescence kernel for the last five datasets are displayed in Figure 8. The PSF have smaller width than the kernel, *i.e.* the temporal resolution is significantly improved by deconvolution; it is in general of order Δt (the inter-frame period) rather than τ_d (the calcium reporter decay time), see Table II. The results hold even in the absence of information about the kernel: in most datasets, the temporal resolution of blind adaptive BSD is very close if not equal to the one of BSD with ground-truth kernel parameters. Importantly, we find an excellent overall agreement between the empirical PSF and the predicted one using the train set parameters: except for datasets 3 and 5 which have unreliable ground-truth, the predicted widths are notably more accurate than naive estimates such as the rise time τ_r or the sampling rate Δt . The prediction is overly pessimistic for dataset 4 because the PSF is computed for an isolated spikes, whereas only spike bursts are detected in practice. The offsets are accurate for datasets 6-10 but not for datasets 1-5, probably due to delays between the electrophysiological and fluorescence recordings. When the inferred parameters are used, the prediction is less accurate owing to the difficulty to estimate the signal-to-noise ratio for single spikes.

Interestingly, we found that convolving the inferred spike train with the point-spread function computed using the inferred parameters significantly improves the accuracy for all algorithms and kernel initializations for the first five datasets, and has little to no impact for the last five datasets. Indeed at low signal-to-noise ratio, fluorescence transients may appear in advance (Fig. 6 b) or delayed (Fig. 6 c) with respect to the action potential and the corresponding point-spread function spans across several time windows. In contrast, datasets 6,7,8,10 have higher signal-to-noise ratio, hence thinner point-spread function (see Fig. 6 (d)) and convolution does not change significantly the prediction at $20Hz$. The small performance drop for dataset 8 can be explained by the overall underestimation of the signal-to-noise ratio.

Super-resolution We tested super-resolution for the datasets 6-10 which have the largest signal-to-noise ratios. The same pipeline is used, but BSD is applied with super-resolution upsamplings factors of 2 (datasets 6,7,8), 6 (dataset 9) and 3 (dataset 10). Figure 8 shows for each dataset the empirical point-spread function of the inferred spikes with (yellow) and without (red) super-resolution. Super-resolution reduces both the offset and the width of the point-spread function for all datasets, see Table IV. The most significant improvement is for dataset 9, which has high a signal-to-noise ratio and a relatively low sampling rate. Overall, super-resolution is relevant when the main source of temporal error is the sampling rate rather than the noise level.

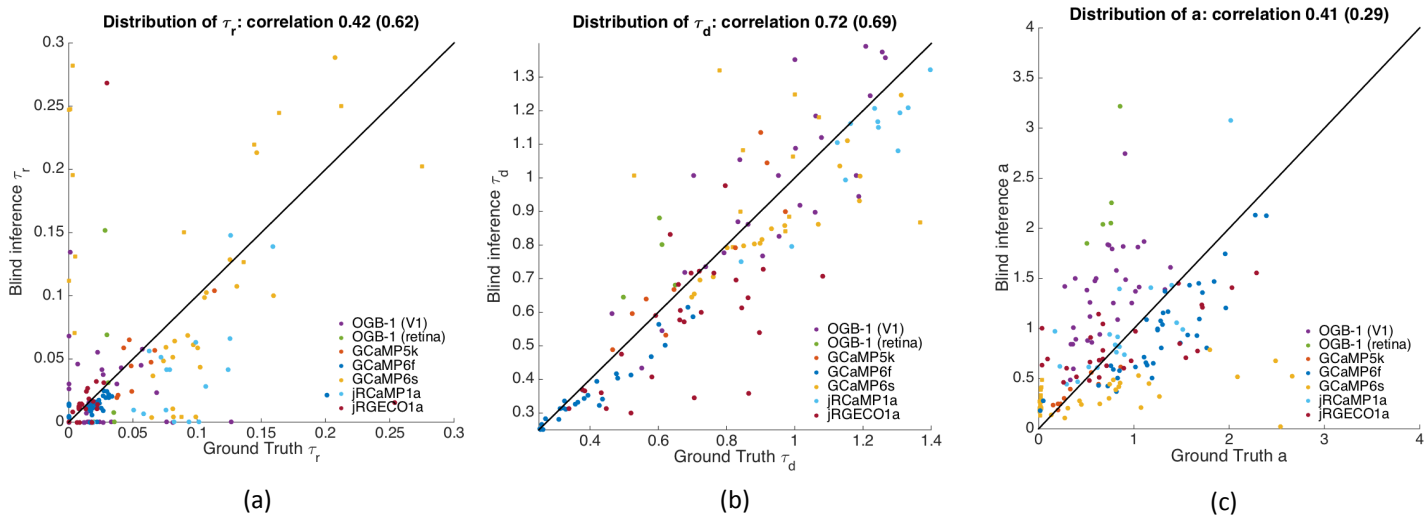


FIG. 7: **Evaluation of the generative models parameter inference for the ten SpikeFinder datasets** For each recording, we derive 'ground-truth' generative model parameters τ_r, τ_d, a by a temporal regression of the fluorescence against the spikes measured by electrophysiology, see Section I G 3. The values are compared to blind estimates obtained using adaptive BSD. (a) Rise time. (b) Decay time. (c) transient amplitude a (each fluorescence trace is normalized to unit variance). Recordings from datasets 3 and 5 are represented as squares, and the others as disks. The correlation coefficient are computed either using all datasets or excluding datasets 3 and 5 (values in parenthesis), for which the ground-truth is inaccurate due to offset between electrophysiological and fluorescence recordings.

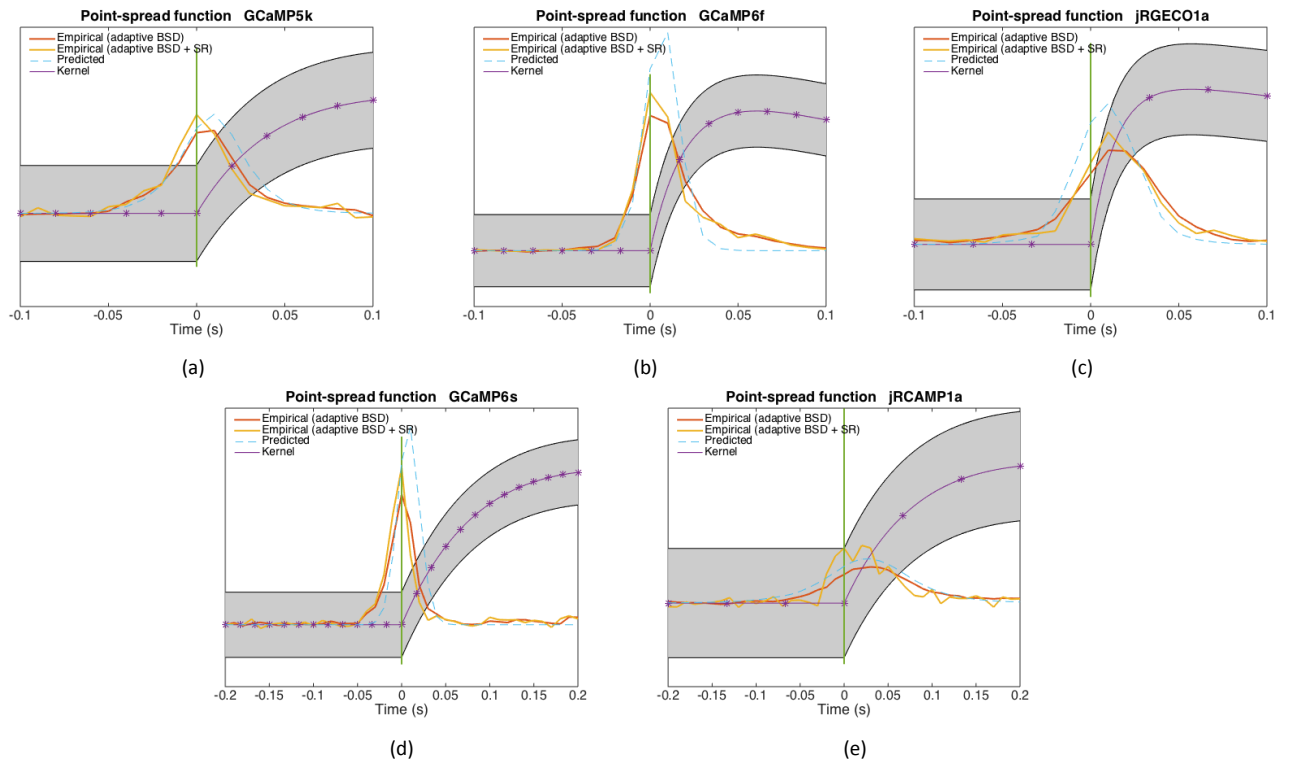


FIG. 8: **Temporal accuracy of BSD for datasets 6-10 of SpikeFinder.** For each dataset, spikes are inferred with adaptive BSD either using super-resolution or not. After resampling to $100Hz$, empirical point-spread functions are computed and averaged for each dataset. Dashed blue curves denote the predicted point-spread function using Eqn. 29 and Table II. The corresponding fluorescence kernel is displayed for comparison. Asterisks indicate fluorescence measurements.

Comparison with other submissions Our algorithm compares positively with other algorithms that have been tested during the competition. According to the current leaderboard, it performs similarly as Team 6 (OASIS [17]), and is above Team 7 (Suite2p [35]) of the original contest. Both teams also used a non-negative deconvolution framework, followed by a reconvolution by a PSF-like kernel. Both teams used the training set to determine hyperparameter values: Team 6 for the convolution kernel, sparsity prior and PSF and Team 7 for the convolution kernel and the PSF; for Team 7 the sparsity prior $\lambda = 0$ was used as in the non-negative deconvolution baseline presented here. In contrast, our approach does not require any training set; the fluorescence kernel, sparsity prior and PSF are determined automatically for each neuron; it can therefore be applied to any experimental configuration.

Dataset	1	2	3	4	5	6	7	8	9	10
Calcium Indicator	OGB-1	OGB-1	GCaMP6s	OGB-1	GCaMP6s	GCaMP5k	GCaMP6f	GCaMP6s	jRCAMP1a	jRGECO1a
BSD (train set)	-1 ± 36	-36 ± 66	115 ± 194	-34 ± 96	142 ± 117	4 ± 19	3 ± 12	-4 ± 16	14 ± 39	13 ± 22
adaptive BSD (blind)	-1 ± 37	-42 ± 67	3 ± 136	-54 ± 111	31 ± 72	5 ± 19	5 ± 13	5 ± 19	35 ± 46	14 ± 23
adaptive BSD (train set)	-2 ± 36	-37 ± 66	98 ± 175	-31 ± 96	121 ± 93	5 ± 19	6 ± 13	1 ± 14	32 ± 46	16 ± 23
adaptive BSD + SR (train set)	/	/	/	/	/	2 ± 16 (-16%)	4 ± 11 (-22%)	-3 ± 11 (-15%)	19 ± 30 (-37%)	14 ± 19 (-14%)
BSD (ground truth)	-8 ± 33	-40 ± 67	15 ± 129	-32 ± 87	132 ± 115	4 ± 19	3 ± 12	-6 ± 18	11 ± 42	14 ± 22
BSD + SR (ground truth)	/	/	/	/	/	0 ± 16 (-16%)	0 ± 11 (-15%)	-9 ± 16 (-2%)	-13 ± 25 (-36%)	10 ± 16 (-28%)
Predicted (train set)	9 ± 43	14 ± 66	0 ± 19	23 ± 167	0 ± 17	9 ± 19	7 ± 9	8 ± 12	28 ± 49	9 ± 17
Predicted (blind)	5 ± 24	3 ± 46	3 ± 25	1 ± 60	2 ± 19	9 ± 16	5 ± 10	8 ± 18	24 ± 50	8 ± 19

TABLE IV: **Empirical and predicted point-spread functions for the SpikeFinder challenge.** For each dataset and each algorithm, the point-spread function is fitted with a Gaussian distribution of mean μ and standard deviation σ . the latter are displayed as $\mu \pm \sigma$ in ms. For super-resolution BSD, we also indicate in parenthesis the relative gain in mean square temporal error $\sqrt{\mu^2 + \sigma^2}$ compared to regular BSD.

D. Light-sheet Imaging of Zebrafish

Compared to standard fluorescence microscopy techniques, such as confocal or two-photon epifluorescence microscopy, light-sheet imaging allows for a parallelization of the recording, yielding ~ 100 -fold increase in data-throughput [2–4]. When applied to zebrafish larvae, this enables simultaneous recording of the quasi-entirety of the neurons ($\sim 100,000$ units) at typically 1 brain/second. The BSD algorithm might prove to be particularly useful for such experiments, as the size of individual datasets precludes supervision. Furthermore, the gain in speed with respect to con-oopsi should also be beneficial as it may allow one to carry out the spike inference on the fly.

To illustrate this latter claim, we test con-oopsi and BSD inference algorithms, as well as MLspike, one of the top performing algorithms of the SpikeFinder contest [13] on a typical whole-brain recording, consisting of 1,800 successive volumetric stacks sampled at 1 stack/second, each of them comprising 20 z-sections. The experiment is performed on a 5 dpf larva expressing the GCaMP5 reporter panneurally. After segmentation, 255463 fluorescence traces encompassing the brain volume are processed independently. The baseline is computed as described before and the spike deconvolution is then carried out using both BSD and con-oopsi on an Intel Xeon Phi (28 cores) computer. In line with our observations of Section II-D2, we find that BSD achieves a 7-fold increase in speed compared to con-oopsi, and a X-fold increase compared to MLspike, see table V. Under these experimental conditions, the computation time with BSD matches the duration of the experiment itself (20 minutes), and is thus compatible with real-time spike inference. Importantly, the computation time per voxel is fairly stable with BSD, whereas some voxels use up to 200 times more time to be processed than others with con-oopsi.

These brain-scale simultaneous recordings allow one to compute the correlation of neuronal pairs activity, which might then be used to extract information regarding the large-scale functional organization of the brain. In this context, we examine whether the correlation statistics of the spike-inferred signals may be significantly different from the one computed using the raw DF/F signals. For this purpose, we use a 2D recording acquired at 20 frame/second for 20 minutes in a 5dpf-old zebrafish larva expressing the genetically encoded indicator GCaMP3 (elavl3:GCaMP3). Automatic segmentation allowed us to identify 8082 individual neurons or neuropil regions of similar area, and the inference is then carried out on the ROI-averaged fluorescence traces. The rise and decay times are inferred for all neurons (see Annex H). The average values of these two time-constants are then used to perform spike inference.

Figure 9a displays the time-averaged image of the brain section. Fluorescent traces and associated inferred spike trains for 5 representative neurons located in various brain regions are shown in Figure 9b. As expected, the deconvolved spike trace appear much sparser and less noisy than the original fluorescent signal. The pair-wise correlations, corrected for uniform coherent noise, are then computed for both the raw DF/F signal and the inferred spike traces. We find the correlation distribution to be much more peaked after deconvolution (Figure 9c) which reflects in the more uniform appearance of the associated correlation matrix (Figure 9d).

This difference may have two possible origins. First, it may reflect the gain in temporal precision brought along by the spike inference, which may reduce the correlation of neuronal pairs that tend to discharge coherently (due to common inputs for instance), but with a slight systematic time-lag. A second explanation is related to the denoising property of the inference. In light-sheet imaging, the noise tends to display significant spatial correlation. This is notably due to the motion of small absorbing objects such as red cells that project elongated shadows and produce characteristic streaking features. Provided that these artifacts have characteristic timescales distinct from the spike-induced fluorescent transient, they are not interpreted as actual spike by BSD. This latter interpretation is confirmed

by the fact that the highly negatively correlated pairs in the raw fluorescence signals are mostly confined within thin bands aligned along the beam direction (Figure 9e). For the same neuronal pairs, the correlation value computed from the inferred signal is thus largely reduced (Figure 9f).

Algorithm	Total run time	Average run time per voxel
con-oopsi	124 minutes	0.38s (min: 0.31 s, max: 110s)
MLspike	120 minutes	0.37s
BSD	18 minutes	0.051 s

TABLE V: Time for performing deconvolution on voxelated data

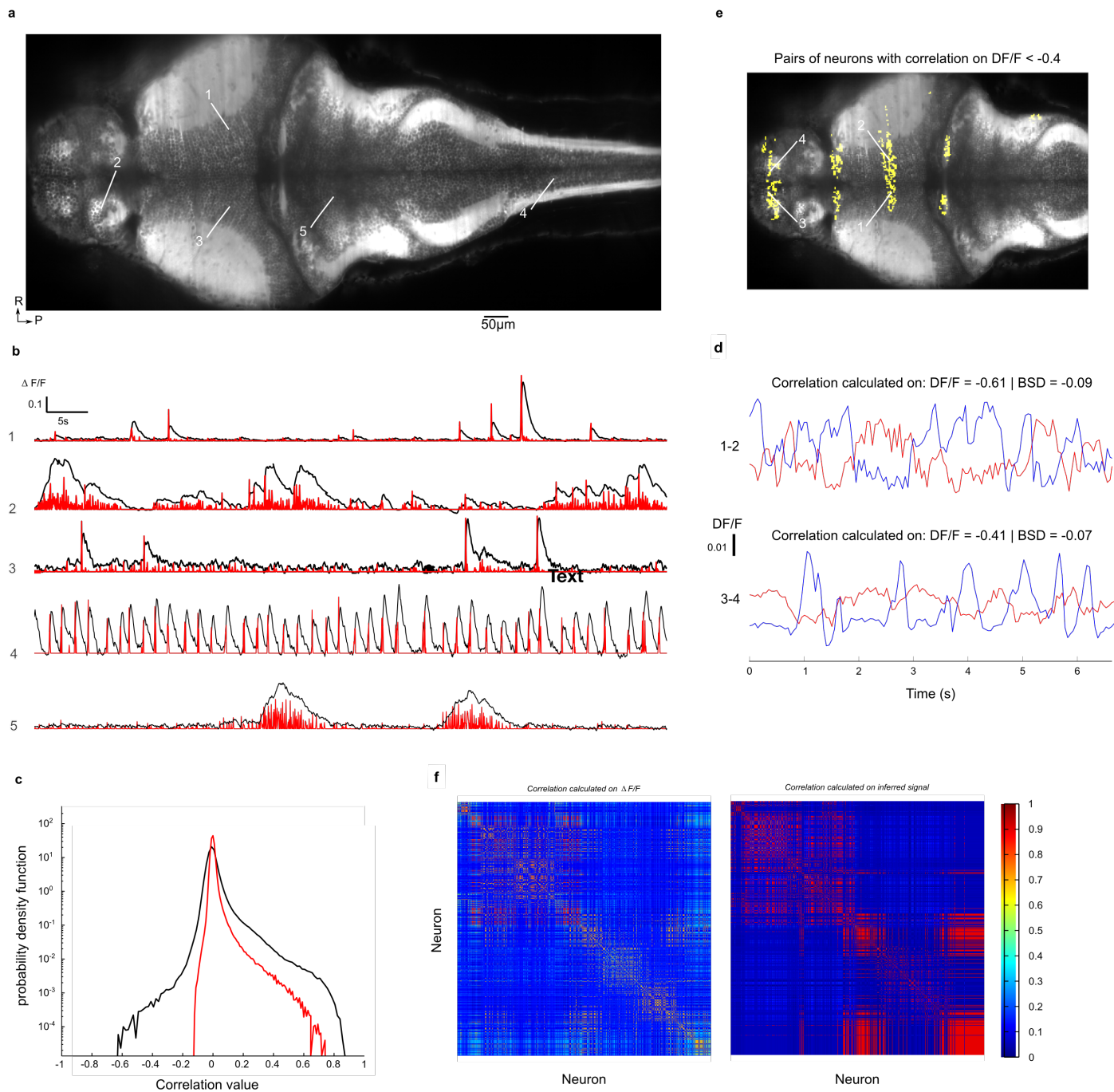


FIG. 9: (a) Bottom: Individual traces of 5 neurons recorded at 20Hz from 6dpf larvae, in black curve DF/F, in red resulting signal from BSD deconvolution algorithm. Top: Time-averaged image of a brain slice of the larva, the white arrows give the location of the 5 neurons. (b) Distribution of pair-wise correlations of DF/F (black) and signal after BSD deconvolution. The data were obtained from a 20Hz, 20 min long experiment on a 6dpf larva. (c) Time-averaged image of a brain slice of the larva. In yellow, pairs of neurons that display a correlation on DF/F inferior to -0.4. (d) Top: Pair of neuron DF/F traces that display a pair-wise correlation calculated on DF/F of -0.61 and a pair-wise correlation calculated after BSD deconvolution of -0.09. Down: Pair of neuron DF/F traces that display a pair-wise correlation calculated on DF/F of -0.41 and a pair-wise correlation calculated after BSD deconvolution of -0.07. (e) Correlation matrix computed from DF/F. (f) Correlation matrix computed from the signal after deconvolution

Discussion

The last few years have seen the release of numerous spike inference algorithms [32]. Their increasing complexity make them poorly interpretable and their performance, beyond the specific conditions for which they have been optimized, can thus be difficult to predict. This is particularly true for supervised machine learning approaches. Although they can offer excellent results on datasets on which they have been trained [14] - which requires the availability of a ground truth, i.e. a simultaneous electrophysiological measurement of the actual spike train - they become less reliable when generalized to other experimental conditions and/or calcium reporters.

In this context, forward generative models, such as non-negative sparse deconvolution, offer a more robust and tractable solution to this problem. These models are based on explicit hypothesis regarding the form of the fluorescence response kernel, the statistics of the spike train and the noise signal, and are thus less prone to systematic bias. However, their performance can still be very sensitive to the way the different model parameters are set, as exemplified by the unreliable results that we obtained with the oopsi algorithm. Due to a paucity of theoretical understanding of the expected performances, these failures are generally impossible to anticipate. The broad implementation of inference methods in functional imaging laboratories will thus depend on the robustness of these algorithms as much as on their optimal performance. Neuroscientists need algorithms that are not only efficient and fast, to accommodate the rapidly growing size of calcium imaging datasets, but that also provide them with a reliable way to assess the quality of the inferred signals.

Here we introduced a novel non-negative sparse algorithm, named Blind Sparse Deconvolution (BSD), which was designed to specifically address these issues. This fully unsupervised algorithm features state-of-the-art computational speed, accuracy and adaptability while incorporating a theoretically-grounded framework to derive estimates of the expected deconvolution performance in terms of temporal accuracy and precision-recall of the inferred spike train. These information may be used before recording as guidelines for experimental design, allowing one to choose, in a given experimental context and for a given calcium reporter, the recording rate that will provide the optimal temporal resolution. They also can be used *a posteriori* to estimate error rates and thus provide bounds on the reliability of the inferred spike trains.

One of the main assets of BSD, compared to other generative models, owes to the fact that most model parameters are analytically derived, in particular the sparsity prior. This allows a tractability of the algorithm, but also a gain in speed compared to other approaches (such as constrained-oopsi) that require recursive evaluations of the sparsity prior.

The main output of the algorithm is a continuous signal that approximates the spike-evoked calcium influx for each recorded neuron. This signal can be directly used to characterize *e.g.* the tuning properties of a sensory-responsive neuron through correlation with the sensory input. BSD also provides automatic theoretically-grounded thresholding and thus enables to binarize the signal into active and non-active periods. Although such binarization comes at a cost of a loss of information, as revealed by the cross-correlation with the actual spike train, it can be necessary for the implementation of graphical circuit inference approaches, such as RBM [36] or Ising models, [37–39]. These models aim at interpreting the collective dynamics of large neuronal ensembles by inferring effective interactions between neurons using the measured pairwise firing statistics. In this particular context, the knowledge of the temporal PSF, as offered by BSD, is also very beneficial as it indicates the minimal time bin over which the pairwise correlations can be robustly evaluated.

In calcium imaging, the temporal resolution is generally thought to be limited by the recording rate. However, at high enough SNR, one may in principle overcome this limit. BSD thus introduces temporal super-resolution, which was shown to significantly increase the temporal accuracy of inferred spikes on real data, yielding a temporal resolution better than the recording period. As both calcium reporters and imaging methods will gain in sensitivity and speed, this capability may help to reveal spatio-temporal short term dynamics, such as activity propagating waves, or to investigate the role of spike-timing in neural coding.

Although we demonstrated that BSD provides consistent results over a large spectrum of reporters and experimental conditions, this algorithm would benefit from embedding more features that would address specific conditions. In particular, we noticed that the performance tends to degrade when neurons display sustained periods of bursting activity. First, the moving percentile method may provide inconsistent estimate of the slowly drifting baseline in this case. [We have shown however that the iterative baseline estimation allows to partially correct this issue]. Second, kernel inference and SNR estimate become less reliable as bursts of activity are mistaken for individual spiking events. One way around would consist in constraining the inferred parameters based on their values estimated on neurons exhibiting sparser activity. A second aspect, ignored in the present implementation, is the non-linearity of the fluorescence response of usual calcium reporters. Since the predicted signal is continuous, this could be accounted for a posteriori in a straightforward way provided that an experimental characterization of the fluorescence vs spike-rate relationship is available. Importantly, the modularity of the different features introduced in BSD - sparsity parameter estimation, iterative kernel iteration and PSF estimation - makes it straightforward to incorporate them

into complete calcium imaging packages such as CaIman [40] Suite2P [35] that address other challenges of calcium imaging processing, such as spatial filtering.

Availability

The spike inference program BSD and its companion program for evaluating its accuracy are implemented in MATLAB and both available at <https://github.com/jertubiana/BSD>. By design, it is straightforward to infer spike trains and evaluate a posteriori the precision-recall and temporal accuracy for each recording. Then, users can optionally convolve the inferred spikes with the predicted point-spread function to account for the uncertainty on the spike location. Tutorial scripts are provided for all use cases and the scripts for generating all figures of this article and reproducing the SpikeFinder experiments are also made available.

Acknowledgements

J.T. acknowledges partial support by a fellowship from the Edmond J Safra Center for Bioinformatics at Tel Aviv University. This work was funded by the Human Frontier Science Program (RGP0060/2017) and the Fondation pour la Recherche Médicale (SPF201809007064).

'Declarations of interest: none

-
- [1] T. F. Holekamp, D. Turaga, and T. E. Holy, "Fast three-dimensional fluorescence imaging of activity in neural populations by objective-coupled planar illumination microscopy," *Neuron*, vol. 57, no. 5, pp. 661–672, 2008.
 - [2] T. Panier, S. A. Romano, R. Olive, T. Pietri, G. Sumbre, R. Candelier, and G. Debrégeas, "Fast functional imaging of multiple brain regions in intact zebrafish larvae using selective plane illumination microscopy," *Frontiers in neural circuits*, vol. 7, 2013.
 - [3] M. B. Ahrens, M. B. Orger, D. N. Robson, J. M. Li, and P. J. Keller, "Whole-brain functional imaging at cellular resolution using light-sheet microscopy," *Nature methods*, vol. 10, no. 5, pp. 413–420, 2013.
 - [4] S. Wolf, W. Supatto, G. Debrégeas, P. Mahou, S. G. Kruglik, J.-M. Sintes, E. Beaupaire, and R. Candelier, "Whole-brain functional imaging with two-photon light-sheet microscopy," *Nature methods*, vol. 12, no. 5, pp. 379–380, 2015.
 - [5] E. Yaksi and R. W. Friedrich, "Reconstruction of firing rate changes across neuronal populations by temporally deconvolved ca2+ imaging," *Nature Methods*, vol. 3, no. 5, pp. 377–383, 2006.
 - [6] T. Sasaki, N. Takahashi, N. Matsuki, and Y. Ikegaya, "Fast and accurate detection of action potentials from somatic calcium fluctuations," *Journal of neurophysiology*, vol. 100, no. 3, pp. 1668–1676, 2008.
 - [7] J. T. Vogelstein, B. O. Watson, A. M. Packer, R. Yuste, B. Jodynski, and L. Paninski, "Spike inference from calcium imaging using sequential monte carlo methods," *Biophysical journal*, vol. 97, no. 2, pp. 636–655, 2009.
 - [8] J. T. Vogelstein, A. M. Packer, T. A. Machado, T. Sippy, B. Babadi, R. Yuste, and L. Paninski, "Fast nonnegative deconvolution for spike train inference from population calcium imaging," *Journal of neurophysiology*, vol. 104, no. 6, pp. 3691–3704, 2010.
 - [9] B. F. Grewe, D. Langer, H. Kasper, B. M. Kampa, and F. Helmchen, "High-speed in vivo calcium imaging reveals neuronal network activity with near-millisecond precision," *Nature methods*, vol. 7, no. 5, pp. 399–405, 2010.
 - [10] Y. Mishchenko, J. T. Vogelstein, and L. Paninski, "A bayesian approach for inferring neuronal connectivity from calcium fluorescent imaging data," *The Annals of Applied Statistics*, pp. 1229–1261, 2011.
 - [11] E. A. Pnevmatikakis, J. Merel, A. Pakman, and L. Paninski, "Bayesian spike inference from calcium imaging data," in *Asilomar Conference on Signals, Systems and Computers, 2013*, 2013.
 - [12] H. Lütcke, F. Gerhard, F. Zenke, W. Gerstner, and F. Helmchen, "Inference of neuronal network spike dynamics and topology from calcium imaging data.," *Frontiers in neural circuits*, vol. 7, p. 201, jan 2013.
 - [13] T. Deneux, A. Kaszas, G. Szalay, G. Katona, T. Lakner, A. Grinvald, B. Rózsa, and I. Vanzetta, "Accurate spike estimation from noisy calcium signals for ultrafast three-dimensional imaging of large neuronal populations in vivo," *Nature Communications*, vol. 7, 2016.
 - [14] L. Theis, P. Berens, E. Froudarakis, J. Reimer, M. R. Rosón, T. Baden, T. Euler, A. S. Tolias, and M. Bethge, "Benchmarking spike rate inference in population calcium imaging," *Neuron*, vol. 90, no. 3, pp. 471–482, 2016.
 - [15] M. A. Picardo, J. Merel, K. A. Katlowitz, D. Vallentin, D. E. Okobi, S. E. Benezra, R. C. Clary, E. A. Pnevmatikakis, L. Paninski, and M. A. Long, "Population-level representation of a temporal sequence underlying song production in the zebra finch," *Neuron*, vol. 90, no. 4, pp. 866–876, 2016.

- [16] E. A. Pnevmatikakis, D. Soudry, Y. Gao, T. A. Machado, J. Merel, D. Pfau, T. Reardon, Y. Mu, C. Lacefield, W. Yang, *et al.*, “Simultaneous denoising, deconvolution, and demixing of calcium imaging data,” *Neuron*, vol. 89, no. 2, pp. 285–299, 2016.
- [17] J. Friedrich, P. Zhou, and L. Paninski, “Fast online deconvolution of calcium imaging data,” *PLoS computational biology*, vol. 13, no. 3, p. e1005423, 2017.
- [18] J. Friedrich, W. Yang, D. Soudry, Y. Mu, M. B. Ahrens, R. Yuste, D. S. Peterka, and L. Paninski, “Multi-scale approaches for high-speed imaging and analysis of large neural populations,” *bioRxiv*, p. 091132, 2016.
- [19] A. Kazemipour, J. Liu, P. Kanold, M. Wu, and B. Babadi, “Efficient Estimation of Compressible State-Space Models with Application to Calcium Signal Deconvolution,” oct 2016.
- [20] I. Selesnick, “Sparse deconvolution (an mm algorithm).,” *Connexions*, 2012.
- [21] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, “Online dictionary learning for sparse coding,” in *Proceedings of the 26th annual international conference on machine learning*, pp. 689–696, ACM, 2009.
- [22] W. Freeman, F. Durand, Y. Weiss, and A. Levin, “Understanding and evaluating blind deconvolution algorithms,” 2009.
- [23] M. J. Rust, M. Bates, and X. Zhuang, “Sub-diffraction-limit imaging by stochastic optical reconstruction microscopy (storm),” *Nature methods*, vol. 3, no. 10, pp. 793–796, 2006.
- [24] B. Huang, W. Wang, M. Bates, and X. Zhuang, “Three-dimensional super-resolution imaging by stochastic optical reconstruction microscopy,” *Science*, vol. 319, no. 5864, pp. 810–813, 2008.
- [25] M. Fernández-Suárez and A. Y. Ting, “Fluorescent probes for super-resolution imaging in living cells,” *Nature Reviews Molecular Cell Biology*, vol. 9, no. 12, pp. 929–943, 2008.
- [26] M. Heilemann, S. Van De Linde, M. Schüttelpelz, R. Kasper, B. Seefeldt, A. Mukherjee, P. Tinnefeld, and M. Sauer, “Subdiffraction-resolution fluorescence imaging with conventional fluorescent probes,” *Angewandte Chemie International Edition*, vol. 47, no. 33, pp. 6172–6176, 2008.
- [27] S. W. Hell, “Microscopy and its focal switch,” *Nature methods*, vol. 6, no. 1, pp. 24–32, 2009.
- [28] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly, “Compressed sensing mri,” *IEEE signal processing magazine*, vol. 25, no. 2, pp. 72–82, 2008.
- [29] M. Pachitariu, C. Stringer, and K. D. Harris, “Robustness of spike deconvolution for neuronal calcium imaging,” *Journal of Neuroscience*, vol. 38, no. 37, pp. 7976–7985, 2018.
- [30] B. A. Wilt, J. E. Fitzgerald, and M. J. Schnitzer, “Photon Shot Noise Limits on Optical Detection of Neuronal Spikes and Estimation of Spike Timing,” *Biophysical Journal*, vol. 104, no. 1, pp. 51–62, 2013.
- [31] T.-W. Chen, T. J. Wardill, Y. Sun, S. R. Pulver, S. L. Renninger, A. Baohan, E. R. Schreiter, R. A. Kerr, M. B. Orger, V. Jayaraman, *et al.*, “Ultrasensitive fluorescent proteins for imaging neuronal activity,” *Nature*, vol. 499, no. 7458, pp. 295–300, 2013.
- [32] P. Berens, J. Freeman, T. Deneux, N. Chenkov, T. McColgan, A. Speiser, J. H. Macke, S. C. Turaga, P. Mineault, P. Rupprecht, *et al.*, “Community-based benchmarking improves spike rate inference from two-photon calcium imaging data,” *PLoS computational biology*, vol. 14, no. 5, p. e1006157, 2018.
- [33] J. Akerboom, T.-W. Chen, T. J. Wardill, L. Tian, J. S. Marvin, S. Mutlu, N. C. Calderón, F. Esposti, B. G. Borghuis, X. R. Sun, *et al.*, “Optimization of a gcamp calcium indicator for neural activity imaging,” *The Journal of Neuroscience*, vol. 32, no. 40, pp. 13819–13840, 2012.
- [34] H. Dana, B. Mohar, Y. Sun, S. Narayan, A. Gordus, J. P. Hasseman, G. Tsegaye, G. T. Holt, A. Hu, D. Walpita, *et al.*, “Sensitive red protein calcium indicators for imaging neural activity,” *Elife*, vol. 5, p. e12727, 2016.
- [35] M. Pachitariu, C. Stringer, M. Dipoppa, S. Schröder, L. F. Rossi, H. Dalgleish, M. Carandini, and K. D. Harris, “Suite2p: beyond 10,000 neurons with standard two-photon microscopy,” *Biorxiv*, p. 061507, 2017.
- [36] S. Cocco, R. Monasson, L. Posani, S. Rosay, and J. Tubiana, “Statistical physics and representations in real and artificial neural networks,” *Physica A: Statistical Mechanics and its Applications*, vol. 504, pp. 45–76, 2018.
- [37] L. Meshulam, J. L. Gauthier, C. D. Brody, D. W. Tank, and W. Bialek, “Collective behavior of place and non-place neurons in the hippocampal network,” *Neuron*, vol. 96, no. 5, pp. 1178–1191, 2017.
- [38] S. Cocco, R. Monasson, L. Posani, and G. Tavoni, “Functional networks from inverse modeling of neural population activity,” *Current Opinion in Systems Biology*, vol. 3, pp. 103–110, 2017.
- [39] L. Posani, S. Cocco, and R. Monasson, “Integration and multiplexing of positional and contextual information by the hippocampal network,” *PLoS computational biology*, vol. 14, no. 8, p. e1006320, 2018.
- [40] A. Giovannucci, J. Friedrich, P. Gunn, J. Kalfon, B. L. Brown, S. A. Koay, J. Taxisidis, F. Najafi, J. L. Gauthier, P. Zhou, *et al.*, “Caiman an open source tool for scalable calcium imaging data analysis,” *Elife*, vol. 8, p. e38173, 2019.
- [41] T. Hendel, M. Mank, B. Schnell, O. Griesbeck, A. Borst, and D. F. Reiff, “Fluorescence changes of genetic calcium indicators and ogb-1 correlated with neural activity and calcium in vivo and in vitro,” *Journal of Neuroscience*, vol. 28, no. 29, pp. 7399–7411, 2008.
- [42] The choice of the convention $\mathcal{K}_{ij} = K[\Delta t(i - j + 1)]$ instead of $\mathcal{K}_{ij} = K[\Delta t(i - j)]$ ensures that $\mathcal{K}_{ij} > 0 \forall j \geq i, \mathcal{K}_{ij} = 0 \forall j < i$. Thus, N_i is the count of spikes occurring *after* measurement F_{i-1} and *before* measurement F_i
- [43] Indeed, $\exists \alpha, \beta, \forall t > 0, \sum_{t_i} K(t - t_i) = \sum_{t_j > 0} K(t - t_j) + \alpha e^{-\frac{t}{\tau_r}} + \beta e^{-\frac{t}{\tau_d}}$. We assume here that $\alpha = \beta = 0$ but we could treat them as unknown variables to be inferred
- [44] We go beyond this approximation in Section I-F, when discussing super-resolution
- [45] We approximate the FNR as the probability that the Dirac solution at $i = i_0$ is zero; both probabilities are not strictly equal because there is a small probability that this solution is zero but other solutions at different time steps are non-zero
- [46] In [7], the authors assume $\tau_r = 0$ and that a is fixed

[47] The spikes occurring between measure $i - 1$ and measure i are the $N_{(i-1)+r} \forall r \in [1, s]$

[48] We have $\lim_{\tau \rightarrow \pm\infty} PSF_\tau > 0$ due to false positives

[49] This estimator is more stable with respect to noise and finite sample size than directly computing the first and second moments of the PSF

Annex A: Stability of the single spike solution and the half-spike problem

In Section II, we have not studied the stability of the single-spike solution. We study it here, and discuss when it is a global optimum. Assuming $\hat{N}_i = \delta_{i,i_0} \max \left\{ a - \frac{\lambda}{\sum_i K^2(t)} + \sigma \frac{\sum_i K(t)\epsilon_i}{\sum_i K^2(t)}, 0 \right\} > 0$ and looking for the stability of the solution, w.r.t the other coordinates, we find:

$$-\frac{\partial \mathcal{L}}{\partial N_{i_0+\delta}} = \sigma \sum_i \{K(\Delta t(i - i_0 - \delta)) - K(\Delta t(i - i_0) \cos(\theta(\delta \Delta t)))\} \epsilon_i - \lambda(1 - \cos(\theta(\delta \Delta t))) \quad (47)$$

$$P \left(-\frac{\partial \mathcal{L}}{\partial N_{i_0+\delta}} > 0 \right) = \Phi \left[\frac{\lambda \tan \frac{\theta(\delta \Delta t)}{2}}{\sigma \|K\|} \right] \quad (48)$$

Where $\Phi(x) = \int_x^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$. Therefore, the Dirac solution is stable only if the above probability is small enough for all values of δ . Far away from the spike $\delta \rightarrow \infty$, the angle $\theta_\delta \rightarrow \frac{\pi}{2}$ and we recover $P = \Phi \left[\frac{\lambda}{\sigma \|K\|} \right]$, as in the spikeless signal. On the other hand, the smaller δ , the smaller θ_δ and the probability is higher. For $\lambda = \lambda_{BSD}$ and low noise, the above probability reduces to $P = \Phi \left[z_1 \tan \frac{\theta_\delta}{2} \right]$; the Dirac solution can become unstable. In practice, the result depends on the level of noise: for low σ , the optimum remains close to the Dirac solution, whereas for high noise, we can find 'half-spikes' solutions, of the form $N_i = \frac{\alpha n}{2} (\delta_{i,i_0} + \delta_{i,i_0 \pm 1})$

Annex B: Impact of kernel parameters mismatch on inference

We systematically studied the bias in spike inference that arises when the estimated time constants τ_r and τ_d differ from their true values, τ_r^0 and τ_d^0 . As illustrated in Figure 10 a-d, inferring the spikes with an incorrect convolution kernel leads to systematic errors. The nature of the error depends on the kernel mismatch:

- $\tau_r < \tau_r^0, \tau_d = \tau_d^0$ (Figure 10a): the inferred spikes are split in two, to compensate for the smaller rise time than expected for a single spike.
- $\tau_r > \tau_r^0, \tau_d < \tau_d^0$ (Figure 10b): the inferred spikes are in advance, to compensate for the faster rise of the fluorescence signal.
- $\tau_r = \tau_r^0, \tau_d < \tau_d^0$ (Figure 10c): the inferred spikes exhibit 'echos' to compensate for the slower than expected decay of F .
- $\tau_r = \tau_r^0, \tau_d > \tau_d^0$ (Figure 10d): the inferred subsequent spikes are 'screened' (lower amplitude) to compensate for the slower than expected signal decay.

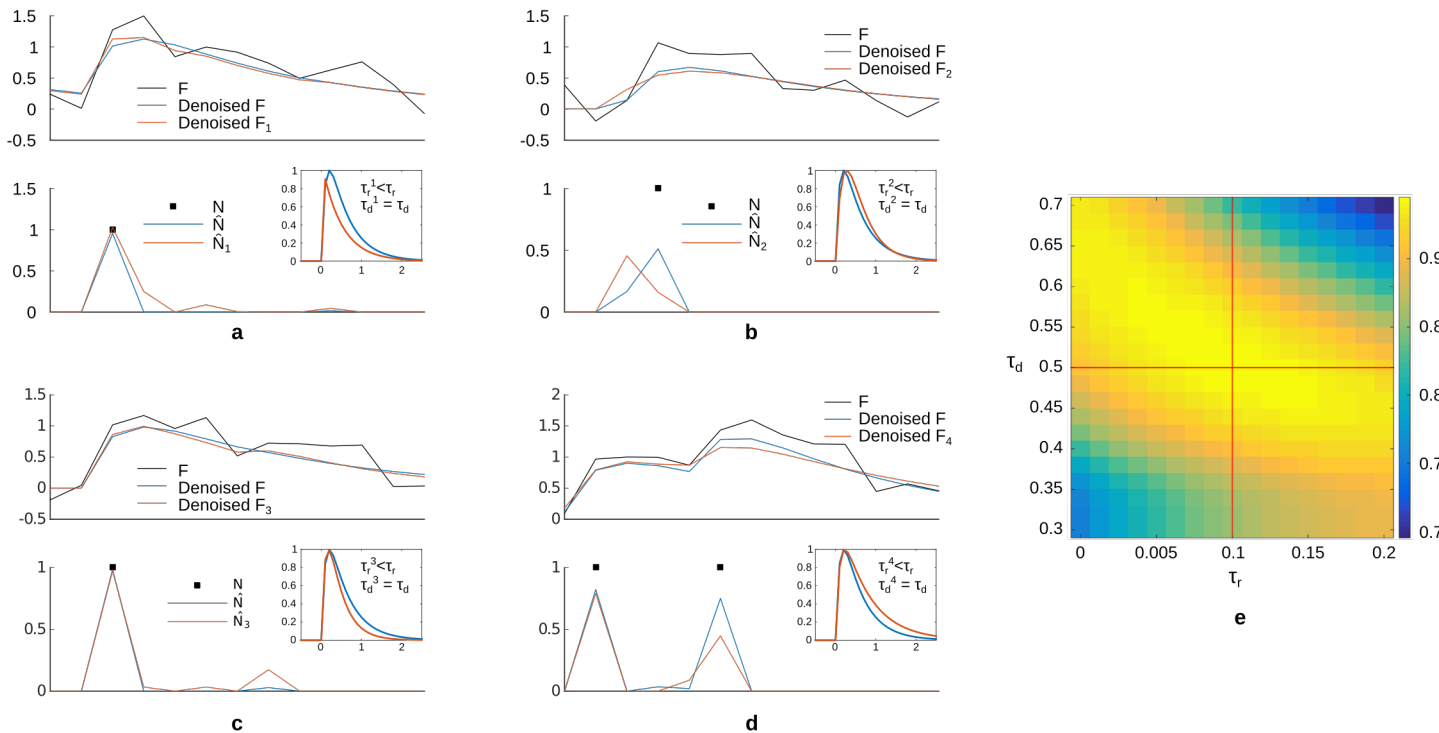


FIG. 10: Left: Example results of spike inference on synthetic data with mismatched convolution kernels. For each of the four figures, a fluorescence signal is generated with a kernel K^0 ; inference is performed with true parameters (blue curves) and with mismatched parameters (red curves). The true and mismatched kernels K^0 and K are depicted (insets). Systematic errors appear in the spike timings. Right: Area Under Curve classification performance with time tolerance $\delta t = 0s$, as a function of the rise and decay time constants. The parameters used to generate the signal, depicted in red, are typical of a GCaMP6 reporter.

We quantified how a kernel misestimation degrades the decoding performance by evaluating the relative reduction in precision-recall (area under curve) for various offsets of τ_r and τ_d (Figure 10e). Interestingly, some direction of the mismatch vector can be less deleterious: when both $\tau_r > \tau_r^0, \tau_d < \tau_d^0$ or $\tau_r > \tau_r^0, \tau_d < \tau_d^0$, the loss in performance remains modest. These findings motivate the use of parameter refinement.

Annex C: Kernel inference: proof of convergence and fast algorithm

We prove here that for isolated spikes and small noise, the cost function $\mathcal{L}(\mathbf{N}, \mathbf{K}) = \frac{1}{2} \|\mathbf{F} - \mathcal{K}\mathbf{N}\|^2 + \lambda(\mathbf{K})\mathbf{1}^T \mathbf{N}$ admits solution $K = K^0$ as local minimum. Denoting $\hat{\mathbf{N}} = \arg \min_{\mathbf{N} \geq 0} \mathcal{L}(\mathbf{N}, \mathbf{K})$

For a signal with a single spike $F_i = aK[\Delta t(i - i_0 + 1)] + \sigma\epsilon_i$, if the noise is small and K is close enough to K_0 , we have: $\hat{N}_i = an\delta_{i,i_0}$, $\lambda = z\sigma\|K\|$. Optimizing over n yields:

$$\begin{aligned}
 n &= \max \left\{ \frac{\cos \phi_K \|K^0\|}{\|K\|} + \frac{\sigma}{a\|K\|} (\tilde{\epsilon}_1 - z), 0 \right\} \\
 \mathcal{L}(\hat{\mathbf{N}}, \mathbf{K}) &= \frac{1}{2} a^2 \|K^0\|^2 (1 - \cos^2 \theta_K) + za\sigma \|K^0\| \cos \theta_K \\
 &\quad + \frac{\sigma^2}{2} \left(\sum_i \epsilon_i^2 - \tilde{\epsilon}_1^2 - z^2 \right) + \sigma^2 z \tilde{\epsilon}_1 - a\sigma \|K^0\| \sqrt{2(1 - \cos \phi_K)} \tilde{\epsilon}_2
 \end{aligned} \tag{49}$$

Where:

$$\begin{aligned}
\|K\| &= \sqrt{\sum_l K(l\Delta t)^2} \\
\|K^0\| &= \sqrt{\sum_l K^0(l\Delta t)^2} \\
\cos \phi_K &= \frac{\sum_i K^0 [(i - i_0 + 1)\Delta t] K [(i - i_0 + 1)\Delta t]}{\|K\| \|K^0\|}
\end{aligned} \tag{50}$$

Note that we recover Eqn. 17 when $K = K^0$. For a signal of multiple isolated spikes $F_i = a \sum_l K [\Delta t(i - i_l + 1)] + \sigma \epsilon_i$, with $|i_l - i'_l| \gg \frac{\tau_d + \tau_r}{\Delta t}$, a similar solution $\hat{N}_i = \sum_l a n_l \delta_{i,i_l}$ can be derived, and \mathcal{L} is self averaging:

$$\mathcal{L}(\hat{\mathbf{N}}, \mathbf{K}) \propto \langle \mathcal{L}(\hat{\mathbf{N}}, \mathbf{K}) \rangle \propto \frac{1}{2} a^2 \|K^0\|^2 (1 - \cos^2 \phi_K) + z a \sigma \|K^0\| \cos \phi_K + \text{Constant} \tag{51}$$

Hence, the function depends on \mathbf{K} only through $\cos \theta_K$. One can check that when $\frac{z\sigma}{a\|K\|} < 1$, the minimum is reached at $\cos \phi_K = 1$, *i.e.* $K = K^0$. This concludes the proof. Although we can not prove more about the radius of convergence, good convergence was achieved in practice after starting from the initialization.

In practice, the optimization with respect to K can be performed efficiently using standard temporal regression tricks. Observe that:

$$\begin{aligned}
\frac{1}{2} \|F - KN\|^2 &= \frac{1}{2} (F^T F - 2F^T \mathcal{K} N + N^T \mathcal{K}^T \mathcal{K} N) \\
&= \frac{1}{2} (F^T F - 2 \text{Trace} [\mathcal{K} N F^T] + \text{Trace} [(\mathcal{K}^T \mathcal{K}) N N^T]) \\
&= \frac{1}{2} \left\{ \sum_{i=1}^T F_i^2 - 2 \sum_{l=0}^{\infty} K(\Delta t(l+1)) \left(\sum_{i=1}^{T-l} F_{i+l} N_i \right) \right. \\
&\quad + \sum_{l=-\infty}^{\infty} \left(\sum_{i=\max(1,1-l)}^{\min(T,T-l)} N_{i+l} N_i \right) \left(\sum_{j=-\infty}^{\infty} K[\Delta t j] K[\Delta t(j+l)] \right) \\
&\quad \left. - \sum_{j=T+1}^{\infty} \left(\sum_i K[(j-i+1)\Delta t] N_i \right)^2 \right\}
\end{aligned} \tag{52}$$

To go from the second line to the third line, we used the translation invariance property of \mathcal{K} , the causality of \mathcal{K} ($\mathcal{K}_{ij} = 0 \forall j \geq i$) and wrote $\sum_{l=1}^T \mathcal{K}_{li} \mathcal{K}_{lj} = \sum_{l=-\infty}^{\infty} \mathcal{K}_{li} \mathcal{K}_{lj} - \sum_{l=T+1}^{\infty} \mathcal{K}_{li} \mathcal{K}_{lj}$. Hence, $\mathcal{L}(\mathbf{N}, \mathbf{K})$ depends on N and F only through:

- the sums $S_1 = \sum_{i=1}^T N_i$ and $S_2 = \sum_{i=1}^T F_i^2$
- the unnormalized cross-correlation between fluorescence and inferred spikes $X(l) = \sum_{i=1}^{T-l} F_{i+l} N_i$.
- the unnormalized autocorrelation function of the inferred spikes $A(l) = \sum_{i=\max(1,1-l)}^{\min(T,T-l)} N_{i+l} N_i$.
- the boundary term $\sum_{j=T+1}^{\infty} \left(\sum_{i=1}^T K[(j-i+1)\Delta t] N_i \right)^2$

The first three terms can be precomputed in $\mathcal{O}(T)$ once for all, and the second and third up to a cutoff $l_{\max} \sim \lfloor 5 \frac{\tau_r + \tau_d}{\Delta t} \rfloor$, such that $K(l_{\max}) \ll 1$. The last one can be computed in $\mathcal{O}(l_{\max})$, by noting that after T , the convolved spikes is a double exponential, with coefficients depending on the $\sim l_{\max}$ last time bins. Overall, the cost function can be evaluated in $\mathcal{O}(l_{\max})$ and optimized efficiently.

Annex D: Detailed computations for the point spread function estimation

We assume a noisy single spike signal, $F_i = aK[\Delta t i - t_0] + \sigma \epsilon_i$, where we write formally $t_0 = \Delta t(i_0 - 1 + r_0)$, with $r_0 \in [0, 1[$; *i.e.* the spike is emitted before measurement i_0 . The likelihood becomes:

$$N_i = an\delta_{i,i_0+\delta}$$

$$\begin{aligned} \mathcal{L}(n, \delta) &= \frac{1}{2} \sum_i F_i^2 + \frac{a^2}{2} \left\{ -2n \sum_i K[\Delta t(i - i_0 + 1 - r)] K[\Delta t(i - i_0 + 1 - \delta)] \right. \\ &\quad \left. + n^2 \sum_i K[\Delta t(i - i_0 + 1 - \delta)]^2 \right\} - a\sigma \sum_i K[\Delta t(i - i_0 + 1 - \delta)] \epsilon_i + \lambda an \\ n^\delta = \arg \max_{n \geq 0} \mathcal{L}(n, \delta) &= \frac{\min \left\{ \sum_i K[\Delta t(i - i_0 + 1 - r)] K[\Delta t(i - i_0 + 1 - \delta)] + \frac{\sigma}{a} \sum_i K[\Delta t(i - i_0 + 1 - \delta)] \epsilon_i - \frac{\lambda}{a}, 0 \right\}}{\sum_i K[\Delta t(i - i_0 + 1 - \delta)]^2} \\ \mathcal{L}(n^\delta, \delta) &= \frac{1}{2} \sum_i F_i^2 + \frac{\min \left\{ \sum_i K[\Delta t(i - i_0 + 1 - r)] K[\Delta t(i - i_0 + 1 - \delta)] + \frac{\sigma}{a} \sum_i K[\Delta t(i - i_0 + 1 - \delta)] \epsilon_i - \frac{\lambda}{a}, 0 \right\}}{2 \sum_i K[\Delta t(i - i_0 + 1 - \delta)]^2} \end{aligned} \quad (53)$$

In the last expression, the term $\rho_{r,\delta} = \sum_i K[\Delta t(i - i_0 + 1 - r)] K[\Delta t(i - i_0 + 1 - \delta)]$ can be computed analytically for all δ and r and is independent of i_0 ; the term $\sum_i K[\Delta t(i - i_0 + 1 - \delta)]^2$ is the usual $\|K\|^2$ and the term involving noise can be rewritten by introducing new, correlated gaussian noises:

$$\begin{aligned} \tilde{\epsilon}_\delta &= \sum_i K[\Delta t(i - i_0 + 1 - \delta)] \epsilon_i \\ &< \tilde{\epsilon}_\delta > = 0 \\ < \tilde{\epsilon}_\delta \tilde{\epsilon}_{\delta'} > &= \sum_i K[\Delta t(i - i_0 + 1 - \delta)] K[\Delta t(i - i_0 + 1 - \delta')] \\ \mathcal{L}(n^\delta, \delta) &= \frac{1}{2} \sum_i F_i^2 + \frac{\min \left\{ \rho_{r,\delta} + \frac{\sigma}{a} \tilde{\epsilon}_\delta - \frac{\lambda}{a}, 0 \right\}}{2 \|K\|^2} \end{aligned} \quad (54)$$

For a given r and noise realization, we can thus compute the optimal δ - and by Monte Carlo averaging, we obtain an estimate of the probability distribution $P(\delta|r)$. To obtain a point spread function in continuous time, it is then transformed into a continuous piecewise-constant probability density through: $P^c(\delta^c \in \mathbb{R}|r) = \frac{P(\lfloor \delta^c \rfloor | r)}{\Delta t}$

And the overall point spread function is obtained by averaging over r , yielding:

$$R(\delta^c) = \int_{r=0}^1 P^c(\delta^c + r|r)$$

In practice, R and r are computed over a discrete grid of the form $k \frac{\Delta t}{s}$.

For the super-resolution case, the computation is almost the same; the only difference being that we reconstruct the spikes with a thinner resolution.

Annex E: Proof of unbiased estimation for super-resolution

We show here that the choice $\lambda_j = z\sigma \|K_j\|$ is best suited for an unbiased (in time) reconstruction of the spikes. We consider again the single-spike setting, with a single spike of a at position $k_0 = (i_0 - 1)s + r_0$, for which $F_i = aK\left[\Delta t(i - i_0) + \frac{\Delta t(2-r_0)}{s}\right] + \sigma \epsilon_i$.

We now look for optima of $\mathcal{L}(\mathbf{N}, \mathbf{K})$ of the form $\hat{N}_{(i-1)s+r} = an\delta_{i,i_0+\Delta}\delta_{r,r_0+\delta}$. Note that instead of doing this computation, we can simply observe that it is a special case of Annex B, using reference kernel $K^0(t) \equiv K(t)$, and measurement kernel $K(t) \equiv K(t - \Delta\Delta t - \frac{\delta\Delta t}{s})$.

$$\begin{aligned}
n &= \max \left\{ \frac{\|K_{r_0}\|}{\|K_{r_0+\delta}\|} \cos \theta_{\Delta,\delta} + \frac{\sigma}{a} \|K_{r_0+\delta}\| (\tilde{\epsilon}_1 - z) \right\} \\
\mathcal{L}(\hat{\mathbf{N}}, \mathbf{K}) &= \frac{1}{2} \|K_{r_0}\|^2 (1 - \cos^2 \theta_{\Delta,\delta}) + za\sigma \|K_r\| \cos \theta_{\Delta,\delta} \\
&\quad + \frac{\sigma^2}{2} \left(\sum_i \epsilon_i^2 - \tilde{\epsilon}_1^2 - z^2 \right) + \sigma^2 z \tilde{\epsilon}_1 - a\sigma \|K_r\| \sqrt{2(1 - \cos \theta_{\Delta,\delta})} \tilde{\epsilon}_2
\end{aligned} \tag{55}$$

Where:

$$\begin{aligned}
\|K_r\| &= \sqrt{\sum_l K \left[l\Delta t + \frac{\Delta t(s+1-r)}{s} \right]^2} \\
\cos \theta_{\Delta,\delta} &= \frac{\sum_i K \left[i\Delta t + \frac{\Delta t(s+1-r)}{s} \right] K \left[(i+\Delta)\Delta t + \frac{\Delta t(s+1-r')}{s} \right]}{\|K_r\| \|K_{r'}\|}
\end{aligned} \tag{56}$$

And $\tilde{\epsilon}_1, \tilde{\epsilon}_2$ are gaussian noises of variance unity (see Annex D). Thus, the optimum over Δ, δ is with highest probability $\delta = \Delta = 0$, and the estimator is unbiased. Note that this result is expected: using the equivalence with a LASSO regression developed in Sec. 5, we know that the coefficients (here, the spikes) are correctly estimated with a uniform λ only when the features (Here, K) are normalized to unity $\sum_i K_{ij}^2 = 1 \forall j$.

Annex F: Kernel inference in the super-resolution setting

Since the convolution matrix \mathcal{K} is not fully translation invariant in the super-resolution setting, the estimation of the kernel is slightly different. For the initial estimation, Eqn. 30 becomes:

$$\begin{aligned}
A_F(l) - \sigma^2 \delta_{l,0} &= a^2 \nu \frac{\Delta t}{s} \sum_{k=1}^{Ts} K \left[\Delta t i - \frac{\Delta t}{s} (k-1) \right] K \left[\Delta t (i+l) + \frac{\Delta t}{s} (k-1) \right] \\
&= a^2 \nu \Delta t \sum_{j=1}^T \frac{1}{s} \left(\sum_{r=1}^s K \left[\Delta t (i-j-1) + \frac{\Delta t(s+1-r)}{s} \right] K \left[\Delta t (i+l-j-1) + \frac{\Delta t(s+1-r)}{s} \right] \right) \\
&\approx a^2 \nu \Delta t \sum_{m=-\infty}^{\infty} \frac{1}{s} \left(\sum_{r=1}^s K \left[\Delta t m + \frac{\Delta t(s+1-r)}{s} \right] K \left[\Delta t (m+l) + \frac{\Delta t(s+1-r)}{s} \right] \right)
\end{aligned} \tag{57}$$

For $s > 1$, this formula is different from Eqn. 30. It can be shown (see Annex G) that the right-hand side has a well-defined limit when $s \rightarrow \infty$, ie in the continuous setting.

Similarly, the iterative kernel update Eqn. 52 is different:

$$\begin{aligned}
\frac{1}{2} \|F - KN\|^2 &= \frac{1}{2} \left\{ \sum_{i=1}^T F_i^2 - 2 \sum_{l=-\infty}^{\infty} \sum_{r=1}^s K \left[\Delta t l + \frac{s+1-r}{s} \right] \left(\sum_{i=1}^T F_{i+l} N_{s(i-1)+r} \right) \right. \\
&\quad \left. + \sum_{l=-\infty}^{\infty} \sum_{r=1}^s \sum_{r'=1}^s \left(\sum_{m=-\infty}^{\infty} K \left[\Delta t m + \frac{\Delta t(s+1-r)}{s} \right] K \left[\Delta t (m+l) + \frac{\Delta t(s+1-r')}{s} \right] \right) \left(\sum_{i=1}^T N_{(i+l-1)s+r} N_{(i-1)s+r'} \right) \right\}
\end{aligned} \tag{58}$$

The sparsity penalty becomes:

$$\lambda^T N = \sum_{j=1}^{sT} \lambda_j N_j = z\sigma \sum_{r=1}^s \sqrt{\sum_{m=-\infty}^{\infty} K \left[m\Delta t + \frac{(s+1-r)\Delta t}{s} \right]} \left(\sum_{i=1}^T N_{(i-1)s+r} \right) \tag{59}$$

Hence, $\mathcal{L}(\mathbf{N}, \mathbf{K})$ now depends on \mathbf{F} and \mathbf{N} through the following quantities:

- the sum $S_2 = \sum_{i=1}^T$
- the sums vector $S_1(r) = \left(\sum_{i=1}^T N_{(i-1)s+r} \right)$
- the cross-correlation matrix $X(l, r) = \left(\sum_{i=1}^T F_{i+l} N_{s(i-1)+r} \right)$
- the autocorrelation tensor $A(l, r, r') = \left(\sum_{i=1}^T N_{(i+l-1)s+r} N_{(i-1)s+r'} \right)$

Altogether, the cost function can be evaluated relatively fast. Note that the complexity of the kernel optimization is now $\mathcal{O}(l_{\max} s^2)$.

Annex G: Various explicit formulas for the double exponential kernel

Various useful formulas for blind sparse deconvolution are consigned, here for double exponential kernels.

Kernel normalization. We normalize K such that $\max_{t \geq 0} K(t) = 1$. This gives:

$$K(t) = \frac{1}{M(\tau_r, \tau_d)} \left[e^{-\frac{t}{\tau_d}} - e^{-\frac{t}{\tau_r}} \right] \mathbf{1}_{t \geq 0}$$

$$M(\tau_r, \tau_d) = \left(\frac{\tau_r}{\tau_d} \right)^{-\frac{\tau_r}{\tau_d - \tau_r}} - \left(\frac{\tau_r}{\tau_d} \right)^{\frac{\tau_d}{\tau_d - \tau_r}} \quad (60)$$

Kernel norms. The L_1 and L_2 norms are computed as follow:

$$\lambda_d = e^{-\frac{\Delta t}{\tau_d}}$$

$$\lambda_r = e^{-\frac{\Delta t}{\tau_r}}$$

$$\|K\| \equiv \sqrt{\sum_{i=-\infty}^{\infty} K[\Delta t i]^2} = \frac{1}{M(\tau_r, \tau_d)} \sqrt{\frac{\lambda_d^2}{1 - \lambda_d^2} - \frac{2\lambda_d \lambda_r}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^2}{1 - \lambda_r^2}} \quad (61)$$

$$\|K\|_1 \equiv \sum_{i=-\infty}^{\infty} K[\Delta t i] = \frac{1}{M(\tau_r, \tau_d)} \left(\frac{\lambda_d}{1 - \lambda_d} + \frac{\lambda_r}{1 - \lambda_r} \right)$$

Kernel norms for super-resolution. The L_1 and L_2 norms for a spike emitted at time $(j-1)s+r$, $r \in [1, s]$ are given by:

$$\|K_r\| = \frac{1}{M(\tau_r, \tau_d)} \sqrt{\sum_i K \left[\Delta i + \frac{\Delta t(s+1-r)}{s} \right]^2} = \sqrt{\frac{\lambda_d^{\frac{2s+1-r}{s}}}{1 - \lambda_d^2} - 2 \frac{(\lambda_d \lambda_r)^{\frac{s+1-r}{s}}}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^{\frac{2s+1-r}{s}}}{1 - \lambda_r^2}} \quad (62)$$

$$\|K_r\|_1 = \frac{1}{M(\tau_r, \tau_d)} \frac{\lambda_d^{\frac{s+1-r}{s}}}{1 - \lambda_d} + \frac{\lambda_r^{\frac{s+1-r}{s}}}{1 - \lambda_r}$$

Kernel overlaps Useful for assessing temporal uncertainty and for kernel estimation

$$\cos \theta(\delta \Delta t) \equiv \frac{\sum_{i=-\infty}^{+\infty} K[i\Delta t] K[(i+\delta)\Delta t]}{\|K\|^2} = \frac{\frac{\lambda_d^{2+\delta}}{1 - \lambda_d^2} - \frac{(\lambda_d^{\frac{\delta}{s}} + \lambda_r^{\frac{\delta}{s}}) \lambda_d \lambda_r}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^{2+\delta}}{1 - \lambda_r^2}}{\frac{\lambda_d^2}{1 - \lambda_d^2} - \frac{2\lambda_d \lambda_r}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^2}{1 - \lambda_r^2}} \quad (63)$$

Boundary term The estimation of the kernel involves the computation of the following boundary term:

$$\sum_{j=T+1}^{\infty} \left(\sum_{i=1}^T K[(j-i+1)\Delta t] N_i \right)^2 = \frac{1}{M(\tau_r, \tau_d)^2} \sum_{j=T+1}^{\infty} \left(\lambda_d^{j-T} \left[\sum_i \lambda_d^{T-i+1} N_i \right] - \lambda_r^{j-T} \left[\sum_i \lambda_r^{T-i+1} N_i \right] \right)^2 \quad (64)$$

$$= \frac{1}{M(\tau_r, \tau_d)^2} \left(\frac{\lambda_d^4 (\sum_i \lambda_d^{T-i} N_i)^2}{1 - \lambda_d^2} - \frac{2(\lambda_d \lambda_r)^2 (\sum_i \lambda_r^{T-i} N_i) (\sum_i \lambda_d^{T-i} N_i)}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^4 (\sum_i \lambda_r^{T-i} N_i)^2}{1 - \lambda_r^2} \right)$$

Kernel overlaps for super-resolution Useful for assessing temporal uncertainty and for kernel estimation

$$\begin{aligned}
A_K(l, r_1, r_2) &= \sum_i K \left[\left(i + l + \frac{s+1-r_1}{s} \right) \Delta t \right] K \left[\left(i + \frac{s+1-r_2}{s} \right) \Delta t \right] \\
&= \frac{1}{M(\tau_r, \tau_d)^2} \left(\frac{\lambda_d^{l + \frac{2(s+1)-r_1-r_2}{s}}}{1 - \lambda_d^2} - \frac{\lambda_d^{l + \frac{s+1-r_1}{s}} \lambda_r^{\frac{s+1-r_2}{s}} + \lambda_r^{l + \frac{s+1-r_1}{s}} \lambda_d^{\frac{s+1-r_2}{s}}}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^{l + \frac{2(s+1)-r_1-r_2}{s}}}{1 - \lambda_r^2} \right)
\end{aligned} \tag{65}$$

In particular:

$$\begin{aligned}
\frac{1}{s} \sum_{r=1}^s A_K(0, r, r) &= \frac{1}{M(\tau_r, \tau_d)^2} \left(\frac{\phi_s(\lambda_d^2)}{1 - \lambda_d^2} - \frac{\phi_s(\lambda_d \lambda_r)}{1 - \lambda_d \lambda_r} + \frac{\phi_s(\lambda_r^2)}{1 - \lambda_r^2} \right) \\
\phi_s(x) &= \frac{1 - x}{s(-1 + x^{-\frac{1}{s}})}
\end{aligned} \tag{66}$$

Boundary-term for super-resolution

$$\begin{aligned}
\sum_{i=T+1}^{\infty} \left(\sum_{j=1}^{sT} K \left[i \Delta t - \frac{(j-1) \Delta t}{s} \right] N_j \right)^2 &= \frac{1}{M(\tau_r, \tau_d)^2} \sum_{i=T+1}^{\infty} \left(\lambda_d^{i-T} \left[\sum_i \lambda_d^{\frac{Ts-j+1}{s}} N_j \right] - \lambda_r^{i-T} \left[\sum_i \lambda_r^{\frac{Ts-j+1}{s}} N_j \right] \right)^2 \\
&= \frac{1}{M(\tau_r, \tau_d)^2} \left(\frac{\lambda_d^{\frac{4}{s}} \left(\sum_j \lambda_d^{\frac{Ts-j}{s}} N_j \right)^2}{1 - \lambda_d^2} - \frac{2(\lambda_d \lambda_r)^{\frac{2}{s}} \left(\sum_j \lambda_r^{\frac{Ts-j}{s}} N_j \right) \left(\sum_j \lambda_d^{\frac{Ts-j}{s}} N_j \right)}{1 - \lambda_d \lambda_r} + \frac{\lambda_r^{\frac{4}{s}} \left(\sum_j \lambda_r^{\frac{Ts-j}{s}} N_j \right)^2}{1 - \lambda_r^2} \right)
\end{aligned} \tag{67}$$

Annex H: Heterogeneity in rise and decay time constants in Zebrafish

Application of BSD to zebrafish data yields heterogeneous distributions of rise and decay times. This means that different regions show different patterns of fluorescence bursts. We see that the heterogeneities have a spatial structure: for instance neurons in the spinal chord tend to have longer rising time constants than neurons in the hindbrain, and neuropil regions have longer decay time. The two possible explanations are that the spike patterns are different in these regions (*e.g.*, regular vs sparse spike trains), and/or that the expression of GCaMP is significantly different. Overall, they motivate the use of heterogeneous time constants.

Annex I: Drawback of approximating Poisson prior to exponential prior in MAP

In order to see why approximating a Bernoulli or Poisson distribution with an exponential approximation with same mean, as is done in oopsi can be problematic for signal reconstruction, we consider the following single-variable inference problem:

$$y = aN + \epsilon \tag{68}$$

Where $\epsilon \sim \mathcal{N}(\mu, \sigma^2)$, $N \sim \text{Bernoulli}(\nu)$. Given y , we wish to estimate N by MAP using either the exact or approximate prior. In the first case, $P(N|y)$ writes:

$$p(N = 1|y) = \frac{e^{-\frac{(y-a)^2}{2\sigma^2}} \nu}{e^{-\frac{(y-a)^2}{2\sigma^2}} \nu + e^{-\frac{y^2}{2\sigma^2}} (1-\nu)} \tag{69}$$

Such that:

$$N^* = \arg \max p(N|y) = \begin{cases} 1 & \text{if } y > \frac{a}{2} - \frac{\sigma^2}{a} \log \frac{\nu}{1-\nu} \\ 0 & \text{Otherwise} \end{cases} \tag{70}$$

In the second case, $P(N|y)$ writes:

$$P(N|y) \propto e^{-\frac{(y-aN)^2}{2\sigma^2} - \nu a N} \quad (71)$$

Such that:

$$N^* = \arg \max p(N|y) = \frac{1}{a} \max(y - \frac{\sigma^2}{\nu}, 0) \quad (72)$$

For typical values such as $a = 1$, $\sigma = 0.2$, $\nu = 0.01$, we get thresholds at respectively 0.32 and 4. Clearly, all spikes ($N = 1$) would be missed using the MAP with approximate exponential prior. On the other hand, the MAP with an exact prior is an unbiased estimate i.e. such that $\langle N^* \rangle = \nu$, and correctly reproduces the average spike rate.

Annex J: Supplementary figures and tables

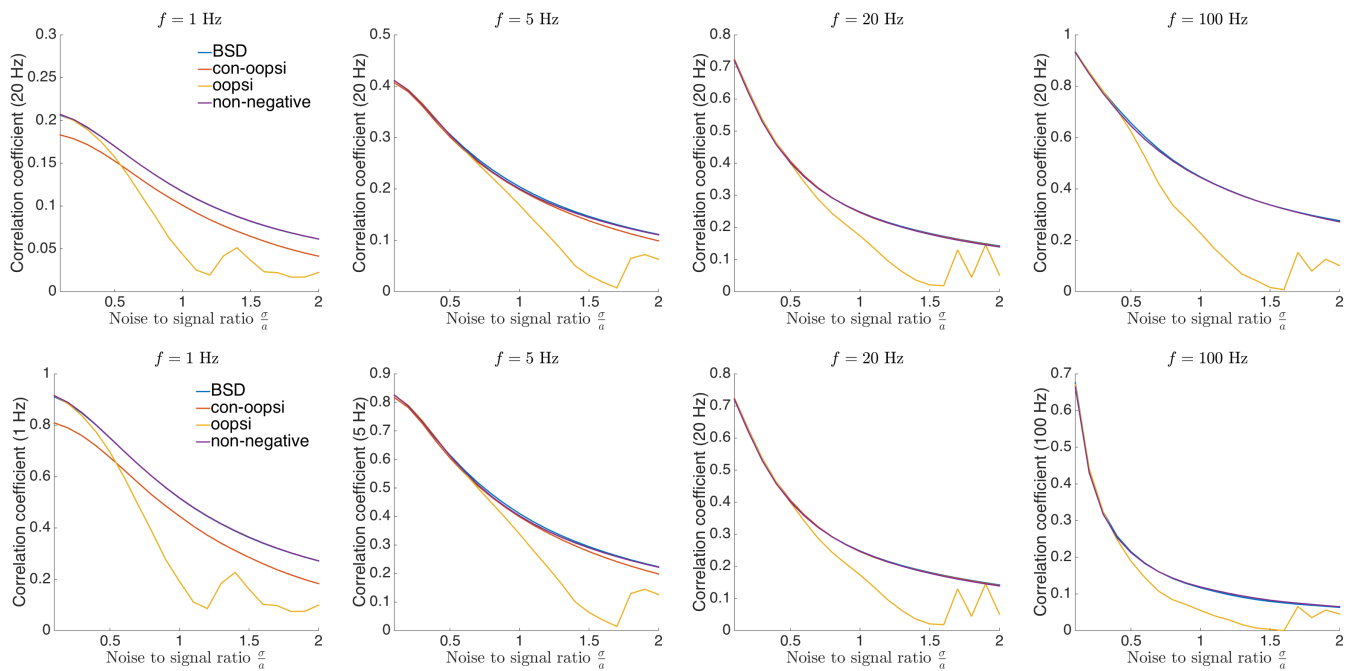


FIG. 11: Correlation between original and ground-truth spike train for various sampling frequencies, evaluation frequency as function of the signal-to-noise ratio; same as Figure 2 for a spiking rate of 1Hz

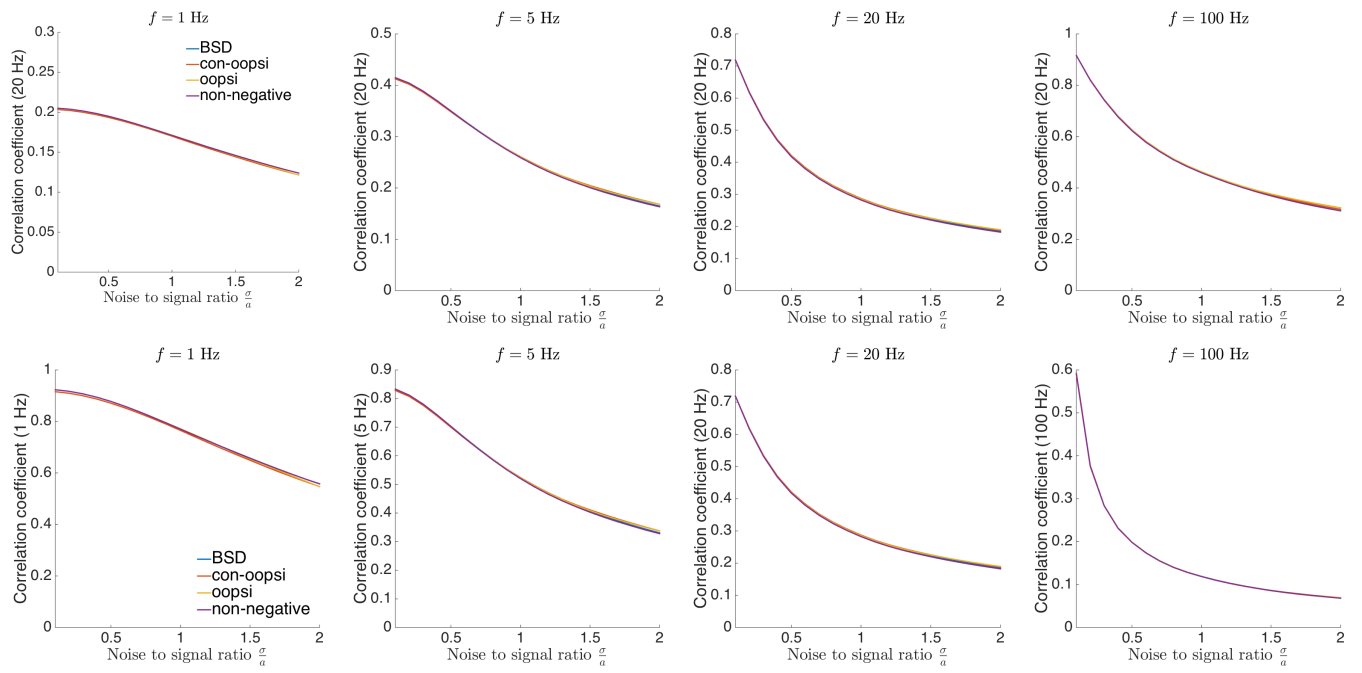


FIG. 12: Correlation between original and ground-truth spike train for various sampling frequencies, evaluation frequency as function of the signal-to-noise ratio; same as Figure 2 for a spiking rate of 5Hz

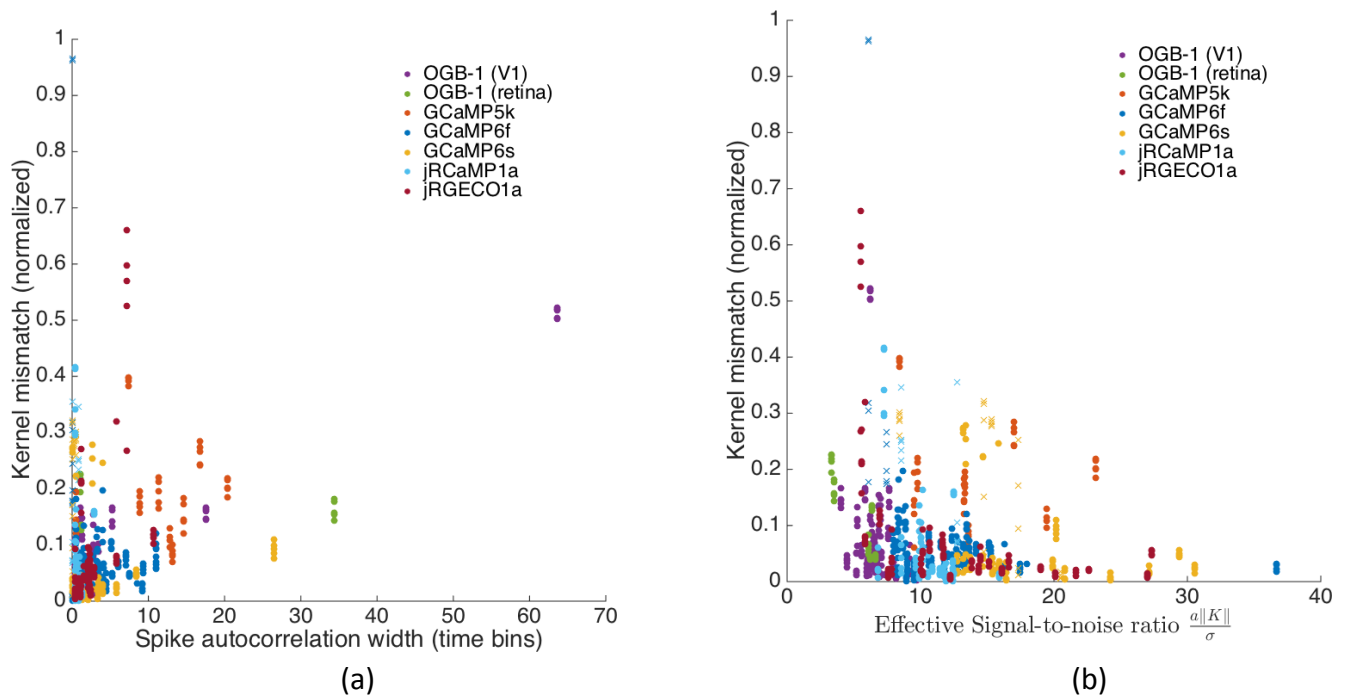


FIG. 13: **Causes of kernel mismatch for adaptive BSD** For each synthetic spike recording generated for Figure 3, we measure the kernel inference error as: $\sqrt{1 - \frac{(\sum_i K_{\text{true}}(i\Delta t)K_{\text{inferred}}(i\Delta t))^2}{(\sum_i K_{\text{true}}(i\Delta t)^2)(\sum_i K_{\text{inferred}}(i\Delta t)^2)}}$. Then, we compare it against (a) the effective signal-to-noise ratio $\frac{a\|K\|}{\sigma}$ (b) the spike 'burstiness'. We measure the later by the width of the spike autocorrelation function divided by the duration of the kernel transient $\tau_r + \tau_d$. The lower the effective SNR and the wider the autocorrelation (in units of kernel transient duration), the larger the mismatch. Additionally, spike trains with fewer than 20 spikes (shown as crosses) frequently result in large kernel mismatch.

Calcium Indicator	Rise time	Decay time	Source
OGB-1	0.010	0.40	Figure 1,2 of [41]
GCaMP5k	0.021	0.27	Supp. Table 3 of [31]
GCaMP6s	0.098	0.38	Supp. Table 3 of [31]
GCaMP6f	0.023	0.10	Supp. Table 3 of [31]
jRCaMP1a	0.0065	1.38	Figure 2 of [34]
jRGECO1a	0.0062	0.32	Figure 2 of [34]

TABLE VI: Rise and decay time constants for a single action potential as defined by Eqn. 2, derived from the literature. When the reference provides the time of peak and half-time decay, the rise and decay constant are obtained by numerical inversion

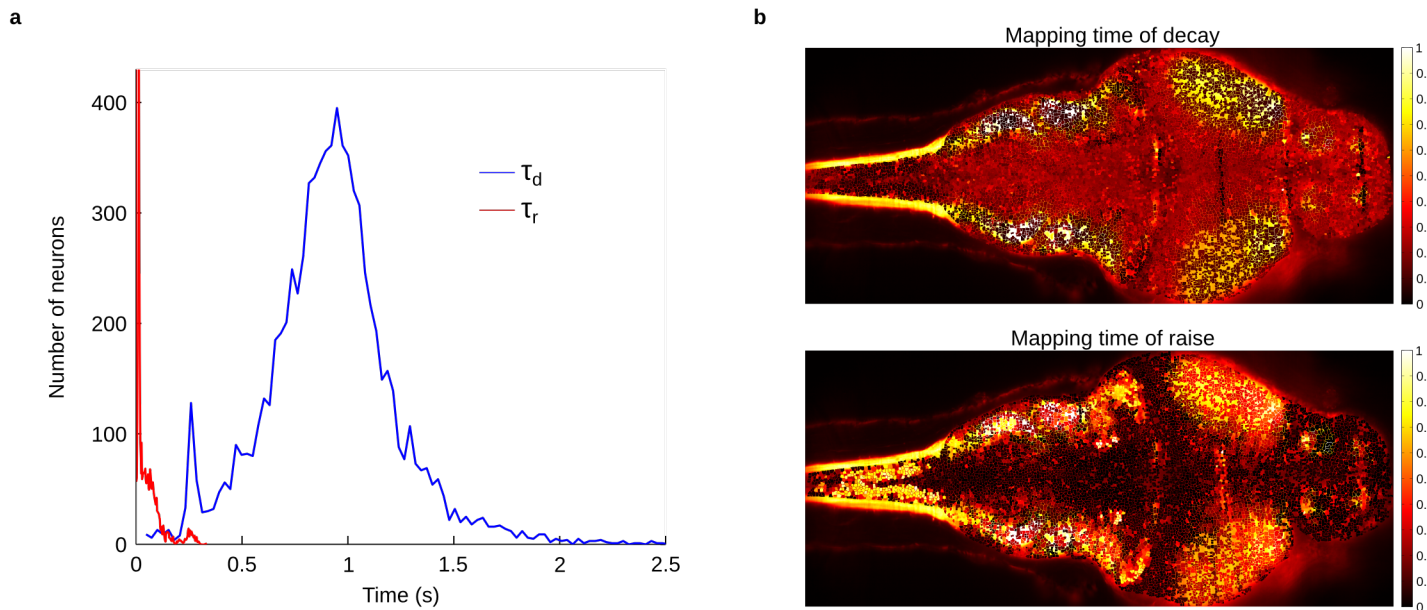


FIG. 14: (a) Distribution of rise and decay time. (b) Mapping of rise and decay time across a neurons