

Compressed and Distributed Sensing of Population Activity for Multiscale Decoding of Motor Cortical Response Properties

Michael A. Shetliffe, Kevin Lorenz, Bryan Pietrzyk, and Karim G. Oweiss

Abstract— Estimating time varying response properties in the motor cortex is an essential task to decode motor neural activity for adequate and reliable neuroprosthetic control. Response properties are typically governed by the precision in neuronal firing and changes in the mean firing rate of individual neurons in relation to the preferred movement direction. This paper proposes a new approach for simultaneously estimating both types of information directly from a multiscale representation of neural data. The approach is based on exploiting the sparsity introduced in the data through a wavelet transformation to first derive distinctive features of neuronal spike waveforms in addition to precision in neuronal firing. Then, changes in mean firing rates are captured along an extended path of the representation to co-localize both response properties. We use spatiotemporal point process models to describe these properties and demonstrate the validity of the approach in decoding a 2D movement trajectory synthesized with 24 neurons directly from the sparsely represented data.

I. INTRODUCTION

NEUROPROSTHETIC control with cortical neural signals relies on decoding motor cortex activity before the actual muscular system response takes place [1]. The decoding process is a cascade of processing steps to be undertaken to enable the subject to experience a natural motor system behavior. These steps generally consist of simultaneously monitoring the activity of multiple cortical neurons through the detection of their extracellular spike waveforms in the recorded data, followed by waveform classification to obtain binary single-neuron spike trains. The functional interdependence between these binary spike trains determines the population code to be decoded for prosthetic limb control.

There is ample theoretical and experimental body of literature suggesting that improvements in motor output behavior are closely associated with the size of the cortical population recorded, the day-to-day stability of these recordings, and the accuracy of the decoding algorithms [2]. There is no apparent consensus in the scientific community on how many neurons are needed for adequate closed-loop

motor control, partly because it is unclear how many, if not all, motor cortex neurons are participating in the information representation [3]. Despite varying population sizes reported in the literature, tens to a few hundred neurons are generally needed to achieve above 90% accuracy [3]. It was also suggested that decoding the population activity needs to be achieved within the *preparatory* period, estimated to be approximately 100-200ms [1],[2], before the actual movement takes place. More complex motor tasks such as finger manipulation will require sampling potentially thousands of possibly non-adjacent neurons, which further places more stringent limitations on the amount of computations that can be performed within a cortical implant's resource constrained environment.

Population response properties of motor cortex neurons seem to heavily depend on two types of information: precision in neuronal firing, and relative changes in mean firing rate characteristics during preparatory and actual movement periods [4]. This information needs to be detected and extracted to determine which of the recorded neurons are task related. On one end, synchrony measured by coincidence in firing seem to be present at various epochs of the preparatory period [5], while on the other end, changes in mean firing rate translating into bursts of activity are also present during movement periods [1]. It is thus evident that the time scale used to quantify response properties is of crucial importance for decoding the activity.

From a signal processing viewpoint, these observations imply a heavily unbalanced distribution of signal power over time as well as temporal information of spike timing precision that cannot be neglected. This leads to a considerable strain on the entire system due to the largely variable data throughput. We have proposed an approach for reducing the strain on the telemetry system in [6], mainly relying on sparse representation of the multi-neuron data to overcome bandwidth limitations. In this paper, we show that this representation enables adequate estimation of both types of information in the population response, without the need for separate signal reconstruction and sorting modules typically needed to estimate these properties. We demonstrate both analytically and computationally that the compressed data contains all the information needed to drive the decoding algorithm for faster, reliable and ubiquitous motor control with implantable BMI systems.

Manuscript received February 16, 2007. This work was supported by NIH under Grant NIH-NINDS-NS047516.

M. A. Shetliffe, K. Lorenz and B. Pietrzyk are with the ECE Department at Michigan State University, East Lansing, MI 48824 (e-mail: shetliff@egr.msu.edu)

K. G. Oweiss is with the ECE Department and Neuroscience Program at Michigan State University, East Lansing, MI 48824 USA (phone: 517-432-8137; fax: 517-353-1980; e-mail: koweiss@msu.edu).

II. THEORY

A. Spatiotemporal Point Process Model

The continuous-time intensity function $\lambda_p(t,x,y,z)$ underlying the point process of a single neuron p can be used to calculate the expected number of events, N_p , encountered during an interval $[T_a, T_b]$

$$E[N_p] = \int_{T_a}^{T_b} \int \int \lambda_p(t,x,y,z) dx dy dz dt \quad (1)$$

Typically, the measured N_p is used in the decoding task by binning the data equally and counting the number of events within each bin. The bin width, $T_b - T_a$, plays an important role in assessing the two types of information in the response: 1) When $T_b - T_a$ is small such that only one event at most is expected in a given bin (precision firing); and 2) $T_b - T_a$ is large permitting subtle changes in $\lambda_p(t,x,y,z)$ to reflect the neuron's relative changes in its mean firing rate as an indicator of its functional role in the population activity during the motor task. One can use the following model

$$\lambda_p(t,x,y,z) = \beta_p(x-x_p, y-y_p, z-z_p) + \alpha_p \Psi(x-x_p, y-y_p, z-z_p, t-\tau_p) \quad (2)$$

to describe the rate function, where β_p , α_p and τ_p denote the background rate, strength and onset of the evoked response Ψ of neuron p located at position (x_p, y_p, z_p) , respectively. For the scope of this paper, we assume that the dependence of the intensity function on the spatial coordinates is stationary, i.e., the rate model exhibits only temporal nonstationarity. This is expressed in the parameters of the *phasic* response α_p and τ_p that are governed by the kinematics of the movement.

B. Measurement Model and Sparse Representation

Following a spike detection stage, the observations usually consist of the time of occurrence of the N_p events (e.g. when the observed signal surpasses a threshold) and the entire spike waveform to be sorted. Assuming that the p^{th} neuron spike waveform is represented by a vector of length N_p , denoted $s_p = [s_p(1) \ s_p(2) \ \dots \ s_p(N_s)]$, the goal is to convert the $N_p \times (m_p N_s \times 1)$ *non-binary* valued observation vector (where m_p is the number of channels on which spikes from neuron p appear) to $N_p \times 1$ *binary* valued vector to quantify the neuron's firing properties¹. Since typically $M_p \ll N_s$, reducing $N_s \rightarrow 1$ is more challenging and also provides a complete solution to the problem when $m_p = 1$. In [6], it was shown that a sparsity transformation operator (e.g., a wavelet transform) can significantly reduce the number of coefficients representing each spike waveform to some $N_c \ll N_s$. In the ideal case, N_c should converge to 1. Given a basis $\phi_j: j=1, \dots, N_s$, the spike is represented by the transform coefficients $w_p^j = \langle s_p, \phi_j \rangle$, such that a considerable sparsity is obtained if the condition

$$\|w\|_q = \left(\sum_j |w_j|^q \right)^{1/q} \leq K \text{ is satisfied, where } K > 0, \text{ and}$$

$0 < q < 2$ [7]. A q approaching zero implies that more sparsity is involved, while $q = 2$ implies no sparsity at all. The objective is to identify the most significant transform coefficient(s) of every spike in such a way to be able to simultaneously preserve the unique features of that spike while minimizing N_c to permit response properties (precise timing and rate) to be directly estimated from the sparsely represented data.

For simplicity, assume that the reference point for marking an event presence is taken as the first sample of the spike waveform. The sampled non-binary spike train of neuron p can be expressed as

$$s_p = \sum_{i \in \{N_p\}} \sum_{k=0}^{N_s-1} s_p[kT_s] \delta[i - kT_s + N_r] \quad (3)$$

Where N_r expresses all the refractory and rebound effects of the neuron, $\delta(\cdot)$ is the Dirac delta function, and T_s is the sampling period. The set $\{N_p\}$ contains all the event times of neuron p . For a given basis ϕ_j , the transformed spike can be expressed as

$$s_p^j = \sum_{i \in \{N_p\}} \sum_{k \in \{N_c^j\}} w_p^j[kT_s^j] \delta[i - kT_s^j + N_r^j] \quad (4)$$

where the set $\{N_c^j\}$ comprises the indices of the most important transform coefficients $N_c^j \ll N_s$ needed to reconstruct the spike waveform with error ε , $T_s^j = T_s / 2^j$. In [8], it was suggested that the $N_c^j \cong \varepsilon^{-(q-2)/2q}$ largest coefficients are kept. Since the transform preserves temporal information, the N_c^j transform coefficients are typically centered around the event time. The height and spread of the coefficients around $i \in \{N_p\}$ will depend on two important characteristics of the signal: first, the magnitude of the coefficients contain all the information about the degree of correlation of ϕ_j and s_p (i.e. features of the spike waveform); and second, the coefficient spread is dependent on the basis kernel support as well as the interspike interval.

For the basis to capture the intrinsic features in the spike waveform simultaneously with the response properties, it is important to consider the spike projection onto the entire library provided by $\phi_j: j=1, \dots, N_s$ [9]. For small j , i.e., fine time scales, the precise timing of the spikes is captured efficiently since the basis is well localized in time. For large j , i.e. coarse time scales, an estimate of the average rate can be obtained since the basis has the longest time support (well localized in frequency). Therefore, the multiscale sparse representation provides a repertoire of information that can be used to obtain the response properties of each neuron.

¹ The case of non closely-spaced array is a special case in which $m_p = 1, \forall p = 1 \dots P$, where P is the population size.

C. Estimation of Response Properties

The objective is to estimate both response properties from single trial sparsely represented data to allow real time decoding of arm trajectory. Consider the example in Fig. 1 in which $\beta_p=10$ events/s, $\alpha_p=25$ events/s and $\tau_p=400$ ms. Outside the range where the rate is rapidly changing (i.e., ~400ms to 800ms), the variability in the individual responses observed in Fig. 1 is typically caused by irregularly spaced events that we consider “point process noise”. This accounts for distinct neural state, habituation, adaptation, etc... This noise is typically cancelled out by averaging across trials.

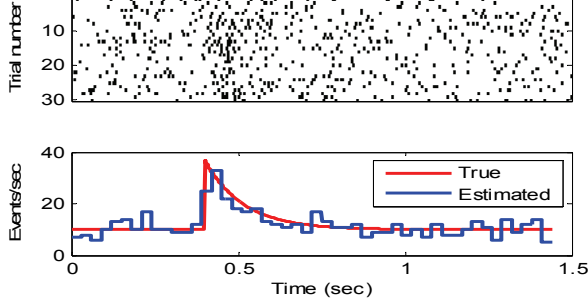


Fig. 1: Intensity function and associated point process for 30 trials. Estimate obtained using a fixed 30 ms bin width averaged across trials.

However, one can observe that if the point process underlying the spiking activity of neuron p is expressed as

$$\begin{aligned} \mathbf{b}_p &= \sum_{k=1}^{N_p} \delta(i-k) = \sum_{k=1}^{N_{p_i}} \delta(i-k) + \sum_{k=1}^{N_{p_u}} \delta(i-k) \\ &= \mathbf{e}_p + \mathbf{z}_p \end{aligned} \quad (5)$$

where both \mathbf{e}_p and \mathbf{z}_p are binary (containing N_{p_i} and N_{p_u} events) the occurrence of 1's in \mathbf{e}_p will be very “regular” in the sense that rapid changes in $\lambda_p(t)$ can be reliably estimated from observing \mathbf{e}_p without trial averaging. On the other hand, the occurrence of 1's in \mathbf{z}_p can be regarded as *spurious* events that do not positively contribute to estimating these rapid variations in a single trial and therefore their effect is reduced by averaging across trials. This is most pronounced when $\lambda_p(t)$ is very smooth and approaches a constant. Using (5) in (4), the regularity term is expressed as

$$\mathbf{e}_p^j = \sum_{i \in \{N_{p_i}\}} \sum_{k \in \{N_{p_i}^j\}} \mathbf{w}_p^j [kT_s^j] \delta[i - kT_s^j + N_r^j] \quad (6)$$

as j becomes large, the time resolution T_s^j becomes coarser, enabling the large coefficients \mathbf{w}_p^j to spread their energy across multiple time indices. Smaller interspike intervals are expected on average during the rapid change epoch of $\lambda_p(t)$, which will translate into a “bump” (or a “ridge”) around the location of the sudden change in the rate. On the other hand, the process noise term,

$$\mathbf{z}_p^j = \sum_{i \in \{N_{p_u}\}} \sum_{k \in \{N_{p_u}^j\}} \mathbf{w}_p^j [kT_s^j] \delta[i - kT_s^j + N_r^j] \quad (7)$$

will be smoothed out for large j , since events in the set $\{N_{p_u}\}$ are more sparsely distributed than those in $\{N_{p_i}\}$ across time.

III. RESULTS

As a representative example of the results, we simulated a population response of 24 neurons involved in encoding a 2D trajectory of arm movement. The response properties were synthesized using a cosine tuning model of preferred directions similar to [1], [10]. No precise timing of spike occurrences was incorporated in this model, so only large j was used to estimate the relative changes in the mean firing rate in single trials. Over the entire library of basis used, we used the mean square reconstruction error between the true rate function and the estimated one in each basis as an indicator of performance

$$E[\varepsilon_p] = \left\| \lambda_p(t) - \mathbf{B} \left(I \left(\langle s_p, \psi_1 \rangle, \dots, \langle s_p, \psi_{N_s} \rangle \right) \right) \right\|_2 \quad (8)$$

where $\mathbf{B}(\cdot)$ is an algorithmic operator that determines the intensity function given a point process, $I(\cdot)$ is an information operator that determines the *most informative* basis to estimate the rate of a given neuron spike waveform. A rate estimate is then obtained as the inverse transform, with the coefficients of all nodes set to zero except those in the node(s) determined to be most informative.

The resultant normalized errors evaluated across the basis library for a selection of neurons from a particular trajectory are shown in Fig. 2. Here we observe that the rate estimate improves as we move down the wavelet tree (increasing j) and the longer basis support is better able to estimate the average firing rate. When we progress beyond level 12 (node 23), however, the basis support becomes too large and is unable to capture the characteristics of the underlying rate. In these experiments, it was generally found that using node 23 alone as the rate estimator yielded the best performance. Assuming that the target binary spike train resolution is Δ (this accounts for action potential duration + refractory and rebound effects), and that the basis filter support is B , then the number of dyadic wavelet levels required to convert the nonbinary spike waveform train to a binary event train with resolution Δ is

$$L = \log_2 \frac{\Delta}{B * T_s} \quad (9)$$

With $B=8$ (*Symmetlet 4* filter length), and a sampling period of 40 μ s, simple calculations show that level 12 basis (nodes 23 and 24) capture the intrinsic features of the rate function over a temporal resolution of approximately 436 ms (with $\Delta = 3$ ms). This resolution is shown below to be very adequate to capture the time constant of rate functions of neurons involved in encoding arm movement trajectory.

Fig. 3 shows actual rate estimates obtained for a particular neuron whose actual spike train is shown in Fig. 3(d). The effect of approximating the rate using a basis of insufficient support (i.e. at a lower decomposition level) is readily apparent from Fig. 3(a). Fig 3(b) shows the “best” rate estimate obtained from node 23 (level 12), and Fig 3(c) shows that there is no benefit gained from including the coefficients from node 24 (level 12) in the rate estimate. All

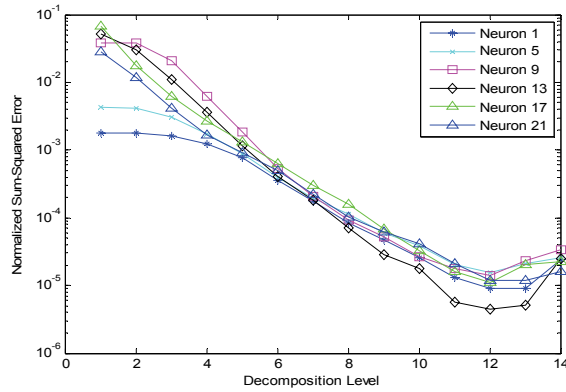


Fig. 2. Normalized error for rate estimators taken from various decomposition levels shown for a selection of neurons from the 24-neuron population.

rate estimates have been normalized in magnitude, and any negative components have been removed to give a physically meaningful rate estimate at all samples of the trial.

Using the inverse of the cosine tuning model across the entire population of 24 neurons, we “decoded” the actual 2D trajectory that was used to generate the original neural responses to obtain the estimates shown in Fig. 4.

IV. CONCLUSION

We have demonstrated the ability of a multiscale multi-neuron data representation to provide single trial estimates of neuronal firing rates. The response properties were found to be best represented at time scales proportional to the time constants of the underlying firing rates. These findings not only eliminate the need to conduct multiple trials to obtain a meaningful population response, but they also allow the response properties to be directly determined from a multiscale decomposition process in which temporal compression is performed along the data transmission path. Our current work is aimed at incorporating other information of response properties such as precise temporal firing in the decoding model as well as comparing the performance of other decoding models to the cosine tuning model used in this paper.

REFERENCES

- [1] D. Taylor, S. I. Helms Tillery, A. B. Schwartz, “Direct Cortical Control of 3D Neuroprosthetic Devices” *Science*, June 2002: Vol. 296. no. 5574, pp. 1829 – 1832, DOI: 10.1126/science.1070291
- [2] J. Wessberg, C. R. Stambaugh, J. D. Kralik, P. D. Beck, M. Laubach, J. K. Chapin, J. Kim, S. J. Biggs, M. A. Srinivasan, and M. A. L. Nicolelis, “Real-time prediction of hand trajectory by ensembles of cortical neurons in primates,” *Nature*, vol. 408, pp. 362–365, 2000.
- [3] M. Serruya, N. Hatsopoulos, M. Fellows, L. Paninski, J. Donoghue, “Robustness of neuroprosthetic decoding algorithms,” *Biological Cybernetics*, vol. 8, no. 3, 219-228, March 2003
- [4] D. Cohen and M. Nicolelis, “Reduction of Single-Neuron Firing Uncertainty by Cortical Ensembles during Motor Skill Learning,” *J. Neuroscience*, vol.24, no. 14, pp. 3574-3582, 2004.
- [5] F. Grammont, A. Riehle, “Spike synchronization and firing rate in a population of motor cortical neurons in relation to movement direction and reaction time,” *Biol. cybern.* vol. 88, no5, pp. 360-373, 2003.

- [6] K. Oweiss, “A systems approach for data compression and latency reduction in cortically controlled brain machine interfaces,” *IEEE Trans BME*, 53(7):1364-77, July 2006.
- [7] Donoho, D. L. “Unconditional bases are optimal bases for data compression and for statistical estimation,” *Appl Comput Harmonic Anal* 1: 100–115, December 1993.
- [8] Donoho, D. L. “Compressed Sensing,” *IEEE Trans Information Theory*, 52(4): 1289-1306, April 2006.
- [9] K. Oweiss, D. Anderson, “Spike Superposition Resolution in Multichannel Extracellular Recordings: A Novel Approach,” *Handbook of Neural Engineering*, ch 22, Edited by Metin Akay, pp. 369-381, IEEE/Wiley, 2007.
- [10] Paninski, L., Fellows, M. R., Hatsopoulos, N. G., Donoghue, J. P. “Spatiotemporal tuning of motor cortical neurons for hand position and velocity,” *Journal of Neurophysiology*, vol. 91, pp.515-532, 2004.

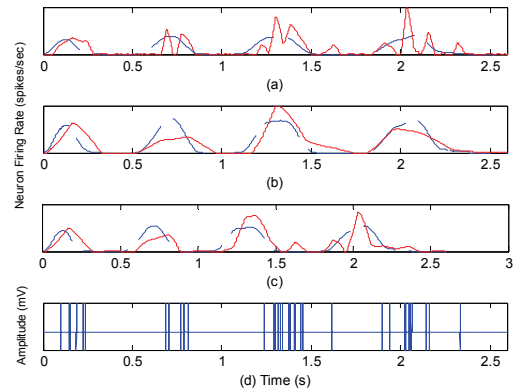


Fig. 3. Examples of rate estimators (solid red lines) for a particular neuron’s firing rate using (a) Nodes 21 and 22, (b) Node 23 and (c) Nodes 23 and 24. The actual rate is plotted with a dashed line in each

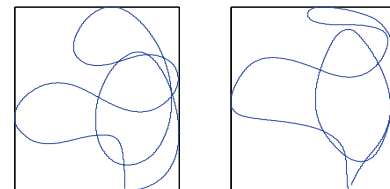
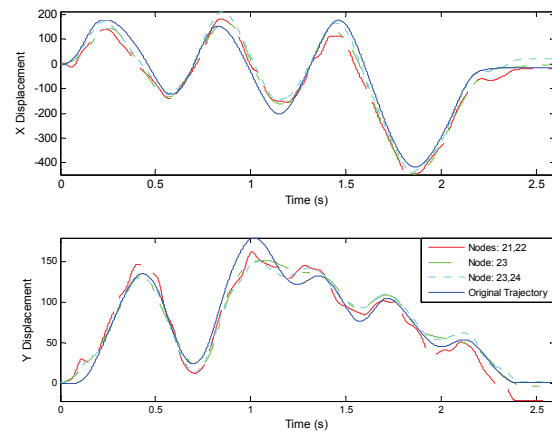


Fig. 4. Reconstructed trajectories (x and y components at top and middle) using several different node combinations for the rate estimate. The solid line shows the actual trajectory. Original 2D trajectory is shown at bottom left, and single-trial reconstructed 2D trajectory from sparsely represented data at bottom right.