

OPTIMIZING SIGNAL CODING IN NEURAL INTERFACE SYSTEM-ON-A-CHIP MODULES

Karim G. Oweiss¹, David J. Anderson^{2,3}, Marios M. Papaefthymiou³

¹ECE Dept., Michigan State University, East Lansing, USA

²BME, ³EECS Depts., University of Michigan, Ann Arbor, USA

Abstract – With the recent technology advances in multichannel microprobe fabrication for neural recording, many limitations still encumber the associated processing and communication capabilities. Consequently, optimizing the information transfer from multiunit neural recording devices is strongly motivated to better understand the underlying mechanisms of neural information processing in the central nervous system in *real time*. We recently developed a framework aimed at processing neural signals with optimized computational complexity. In this work, we propose an associated coding solution aimed at implementing fast and efficient neural interface System-On-a-Chip (SOC) modules capable of exchanging neural information without compromising issues of communication bandwidth and signal fidelity. Preliminary results demonstrate that neural information processing and on-chip coding offers tremendous savings in communication costs compared to raw data transmission for off-chip analysis. Moreover, we demonstrate that bit rates in the order of 1bit/sample can be easily achieved with a 32 channel device and 25kHz sampling rate. Performance illustrations and experimental neural data examples are described in details.

Keywords – **Multichannel neural recording, array processing, wavelet transform, signal coding**

I. INTRODUCTION

Signal and information management of neural recordings from many closely spaced recording and stimulation electrode arrays in the brain poses many challenges in excess of those seen by the current generation of neural prosthetic devices such as cochlear implants [1]. Among these are: the sheer volume of the data, variable signal quality across members of the array and vulnerability to external Electro-Magnetic Interference (EMI) sources. In a static laboratory environment, with large bandwidth to the recording devices, data processing can be managed off-line using massive computing power, and human interaction at a rate slower than real-time. On the other hand, implantable/wearable and hopefully unobtrusive neural implant devices with high channel count interacting at the single neuron or cluster level of detail requires an ultra high bandwidth communication link. Therefore, pre-processing neural signals is sought to be implemented as close as possible to where the signal is acquired to reduce the computational and communication costs associated with neural implants¹. This implies that recording and stimulation

¹ Since neural signals mostly consist of short rapid transients, a rough calculation of the *useful* time where a neural signal is detectable in the acquired data can be calculated to yield a range between 10% (spontaneous activity) to 25% (stimulus driven activity) of the total recording time. This roughly means that 75%-90% of the time, the telemetry system is used to transmit *signal-absent* data.

of multiple neural cell populations needs to be implemented with limited power, bandwidth and in an autonomous manner. Consequently, there is an intrinsic need to develop neural interface SOC modules that optimize the information transfer from brain regions under investigation without compromising issues of bandwidth, detection, and classification accuracy.

Analysis of basic information rates from groups of nerve cells from a multichannel recording device has sparked the need for a coding solution that will allow the transfer of the information across the communication channel without significant loss. As an example, a 96 recording electrode device sampled at 25 kHz per channel leads to an aggregate rate 2.4 Msamples/sec. At 12 bits/word, the bit rate is nearly 29 Mbit/sec, which is well beyond the reach of present day biotelemetry developments. Nevertheless, transmission of 1-2 Mbit/second over a telemetry link from an implant to an extracutaneous device is feasible with the current technology [2,3], which yields a telemetry channel capacity at one bit per sample. This requires a smart way of compressing the basic information in the neural data to adhere to this capacity limit. Basic information in multiunit neural data consists of spike event timing and spike waveform classification result. To transmit this information within the available channel capacity, the coder may require large amounts of computation and probably time delays. Accordingly, an appropriate compromise among power, required bandwidth, computational load, memory, and delay time is needed to yield the best reliable design with highest signal fidelity. These parameters are linked in a very complex way and an optimum solution is beyond the scope of this effort. Nevertheless, some facts listed below can relax some of the parameter constraints to ease the design approach:

1. 32 physical channels from a closely-spaced device is a manageable number of interconnects and signal conditioning components in the volume over the tissue being considered [4].
2. A neural prosthesis system can tolerate delays of *several* milliseconds [5].
3. Several hundred thousand transistors can be operational in a package the size of a standard 14mm burr hole used by neurosurgeons.
4. 25kHz samples/sec quantized to 8 bits is a usual and logical sampling scheme if scaling is controlled.

The following sections will address how to optimize the processing architecture to deliver both the timing and waveform information within the one bit per sample objective.

II. NEURAL SIGNAL CONDITIONING

Due to the nature of the neural signal environment, some conditioning needs to be performed before coding the signal can take place. A recent array processing framework relying on the compaction property of the Discrete wavelet Transform (DWT) [6] has shown great promise in overcoming many challenges in processing neural signals with simple computational blocks [7]. In this context, it was shown that spatial filtering implemented either in the time domain [8] or the multiresolution domain [9,10] greatly reduces the correlated noise usually observed in closely spaced electrode site recordings. Therefore, a key processing block for signal conditioning is a spatial filter that performs *prewhitening* of the neural data to aggregate the energy spread across many physical channels to a few neural channels that can be individually filtered and coded in subsequent processing stages. The block diagram in Figure 1 shows the architecture of the signal processing and coding system considered. The schema utilizes standard components that have been implemented on either the Michigan merged micro-machining and CMOS processes [4], a mixed signal chip or standard digital electronics.

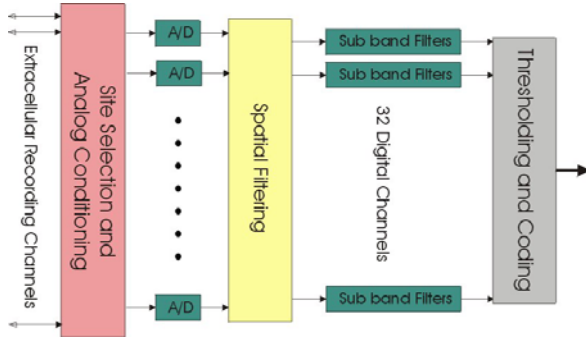


Figure 1: The salient features of the architecture are alternate cross channel processes (site selection, spatial filtering, and block coding) and per channel processes (amplification, scaling, A/D and digital filtering). Depending on the complexity of the overall system chosen, more or less of the blocks will exist.

To gain some insight on the theoretical side, the multichannel neural data can be modeled using the additive noise array model expressed as

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z} \quad (1)$$

where $\mathbf{Y} \in \mathfrak{R}^{M \times N}$ denotes N time samples of an array of M physical microprobe channels, $\mathbf{X} \in \mathfrak{R}^{M \times N}$ the *noise-free* matrix of observations to be detected (linear mixture of multiple neural sources) and $\mathbf{Z} \in \mathfrak{R}^{M \times N}$ denotes a zero mean additive Gaussian noise component with a nontrivial spatial covariance matrix $\mathbf{R}_Z \in \mathfrak{R}^{M \times M}$. The Maximum Likelihood (ML) estimate of the spatial covariance matrix of \mathbf{Y} is expressed as

$$\mathbf{R}_Y = E[\mathbf{Y}\mathbf{Y}^T] \cong \frac{1}{N-1} \sum_{n=1}^N \mathbf{y}[n]\mathbf{y}^T[n] \quad (2)$$

where $\mathbf{y}[n] \in \mathfrak{R}^{M \times 1}$ denotes the n^{th} snapshot of the array.

Using Eigen-Value Decomposition (EVD) [11], \mathbf{R}_Y can be spectrally factored to yield

$$\mathbf{R}_Y = \mathbf{U}_Y \mathbf{D}_Y \mathbf{U}_Y^T \quad (3)$$

where the orthonormal $M \times M$ matrix \mathbf{U}_Y comprises the eigenvectors spanning the column space of \mathbf{R}_Y , and the $M \times M$ diagonal matrix \mathbf{D}_Y contains the rank ordered eigenvalues of \mathbf{R}_Y , $\delta_1 > \delta_2 > \dots > \delta_M$. The matrix \mathbf{U}_Y corresponds to the spatial filter that prewhitens the data matrix \mathbf{Y} , i.e.,

$$\tilde{\mathbf{Y}} = \mathbf{U}_Y^T \mathbf{Y} \quad (4)$$

consists the matrix of *focused* physical channels that spatially compacts the neural signal energy into a fewer number of channels. These are subsequently subband filtered using the DWT operator \mathbf{W}_j in the j^{th} subband as

$$\tilde{\mathbf{Y}}_j = \mathbf{W}_j \{\tilde{\mathbf{Y}}\} \quad j = 0, 1, \dots, J \quad (5)$$

where J denotes the total number of subbands in the DWT stage [10]. Equation (5) yields a temporal energy compaction into a few sparse coefficients. The thresholding property of the DWT enables discarding small coefficients, presumably attributed to noise and very weak signals, thus reducing the amount of information that will be pipelined to the coder at the last stage as illustrated in Fig.1. Denoting the thresholding matrix operator by H , the thresholded wavelet expansion can be expressed as

$$\bar{\mathbf{Y}}_j = H \{\tilde{\mathbf{Y}}_j\} \quad j = 0, 1, \dots, J \quad (6)$$

Equation (6) describes the prewhitened, thresholded wavelet coefficients that need to be coded for transcutaneous transmission from the neural interface *SOC*.

III. CODER OPTIMIZATION

At this stage, two major alternatives exist for compressing the neural data for transcutaneous transmission:

1. Transmission of the processed version of the raw data given in (6) using a *lossy* compression scheme. This is contingent upon having a large bit rate to transmit event timing with full waveform data that can be classified *off-chip*,
2. Detection of neural events in the processed data in (5), and transmitting only the event times across the channel with some side information about classification [12]. This is contingent upon sophisticated detection and classification algorithms *on-chip*, and yields very low transmission rates.

For the purpose of this paper, the former was chosen because the range of compromise is greater, that is, it offers the transmission of evoked potentials in parallel with the

spike data with no special handling and it can take better advantage of information distributed across channels.

IV. METHODS AND PERFORMANCE

A 32-channel device requires a spatial filter U_y with 25×10^6 multiplication per second for prewhitening the energy of distributed spike to a single neural channel. This intensity of computation requires partitioning the 32x32 matrix down to a more manageable load. As an example, choosing eight 4x4 spatial filters reduces the computation to only 3.2×10^6 multiplications per second. This could provide local prewhitening over electrodes that are closely spaced but would not provide any global effects such as artifact reduction or evoked potential extraction at the sensor level. These signals could be derived with a simple row partition and canceled out of all the channels by subtraction. This process could be accomplished by analog referencing as well. This is out of the scope of this paper and is currently under investigation. As far as wavelet basis choice, we have gained experience processing neural data with several wavelet families and have determined that for this compression problem, the ‘symlet4’ family [6] carried down 4 decomposition levels provides satisfactory results and has a modest computational cost of about 2.2×10^6 multiplications/sec.

IV-i Waveform Coding

Figure 2 below shows more details of the coding algorithm. The data are blocked by 16 samples per channel. That will yield one A4 sample/channel/block. The D1 band is disregarded because at a sample rate of 25kHz the spike train is over sampled leaving only noise in the D1 band.

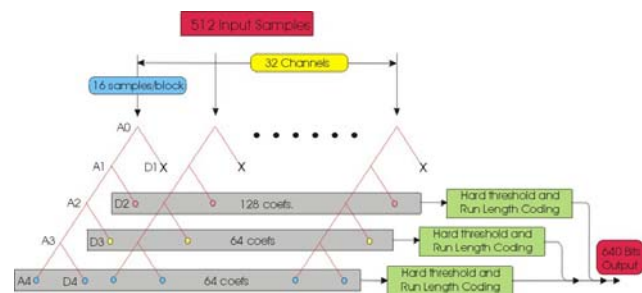


Figure 2: 16 samples from each of 32 channels are encoded using a wavelet decomposition, hard denoising algorithm [9] and run length coding. The resulting compression transmits data at just above 1 bit per sample.

There are three stages to the coding process:

- First the wavelet coefficients are grouped across channels in bands with the A4 and D4 bands being combined because they are small and have similar statistics.
- Next the coefficients are thresholded at levels determined by the noise power on each band.

- Lastly the zero run lengths and the remaining coefficients are coded respectively by count and sign/magnitude with a simple binary quantizer.

The coding efficiency achieved is just slightly more than one bit per sample, which is well within the channel capacity. The minimum delay is .64 msec. which is also well within the tolerance for neuroprosthetic devices. Increasing the buffer length to more blocks of 16 samples would tend to average the bit rate over modulated signals bringing the maximum down and slightly increasing the run length code efficiency. Typical input output of the coding system is shown below in Figure 3. Reduction in wavelet domain vs. time domain is about 30:1 for this data set. That reduction ratio holds the bit rate for 32 channels to below 1Mbit/sec.

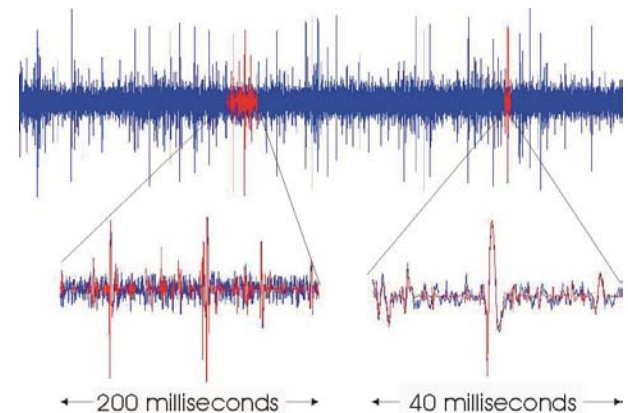


Figure 3: Four seconds of spike data (upper trace) from a chronic rat cortical preparation have been compressed and reconstructed (red) preserving most of the detail of the spikes.

IV-ii Event time Coding

The usual reference for an event time is the time since the previous event time or *run length*. The coding of run lengths depends strongly on the statistics of the intervals, which are related directly to neural properties. Some neurons discharge with nearly regular intervals that are slowly rate modulated while others are Poisson distributed and can tolerate rapid rate modulation. A simple and efficient method for coding run lengths is to count the number of items in the run with a finite register. When the register overflows, a zero value is assigned to that word, making it a marker and the counting starts in successive words until the run length is reached. There is an optimum word length for the counter and that can be determined easily assuming the exponential distribution of run lengths. Figure 4 gives two views of bit rates of exponentially distributed intervals. The fundamental limit of information in bits per interval as a function of event rate is [13]

$$r = -(\log_2(e^\lambda - 1) + \log_2(e^{-\lambda})/(1 - e^{-\lambda}))$$

or $\cong \log_2(e/\lambda)$ (7)

It is not surprising that the optimum word length for the coder is closely related to $\log_2(\text{average run length})$.

This makes it easy to calculate the best word length and the overall bit rate needed for placing each event at its correct time location given an event rate. For rates below .02 to .02, word lengths 7, 8, or 9 bits are very satisfactory using less than .2 bits per data sample for the intervals leaving $\geq .8$ per sample for the numerical value of the data. The algorithm complexity is very low.

V. CONCLUSION

A design methodology for optimizing neural information coding for transmission in neural interface SOC modules has been proposed. This consists the first thrust towards the development of highly sophisticated implantable/wearable neural chips interfaced with small populations of neural cells. The main challenge in the design of the neural interface chips will be to achieve the required level of processing throughput with the minimum possible power dissipation and chip area. The area of the compression chip should not exceed 0.5cm^2 , one third of the approximately 1.5cm^2 available on the platform. The chip will be capable of simultaneously processing 32 recording channels, each at the rate of 25KHz. To that end, the chip hardware should be operating at an approximate rate of 150MHz. Using 16-bit arithmetic, the estimated power dissipation of the proposed programmable chip is expected to be in the 40mW range, yielding a power density below $100\text{mW}/\text{cm}^2$. These issues will greatly benefit from the optimization methodology proposed in this work and are currently under investigation.

VI. REFERENCES

[1] Bell T. E., Wise K. D., and Anderson D. J., "A Flexible Micromachined Electrode Array for a Cochlear Prosthesis," *Sensors and Actuators, A*, 66, pp. 63-69, April 1998
 [2] Wolpaw J. et. al., "Brain-Computer Interface Technology: A Review of the First International Meeting," *IEEE Trans. on*

Rehab. Eng., vol.8, no.2, pp. 164-173, June 2000
 [3] Takeuchi S. and Shimoyama I., "An RF-telemetry System with Shape Memory Alloy Microelectrodes for Neural Recording of Freely Moving Insects," *Proc. of the 1st Int. IEEE-EMBS Conf. on Microtechnologies in Medicine & Biology*, pp. 491-496, Oct. 12-14, Lyon, France, 2000
 [4] Kim, C. and Wise, K. D., "A 64-Site Multiplexed Low-Profile Neural Probe with On-Chip CMOS Circuitry," *Digest 1994 IEEE Symposium on VLSI Circuits*, pp 23-27. June 1994
 [5] Anderson, D.J., et. al., "Batch-fabricated thin film electrodes for stimulation of the central auditory system," *IEEE Trans. on BME*, vol.36, pp. 693-704, 1989
 [6] Mallat, S., *A Wavelet Tour of Signal Processing*, Academic Press, 2nd edition, pp. 413: 1999.
 [7] Oweiss, K.G., and Anderson, D.J., " A Unified Framework for Advancing Array Signal Processing Technology of Multichannel Microprobe Neural Recording Devices," *Proc. of the 2nd IEEE Conf. on Microtechnology in Medicine and Biology*, pp. 245-250, May 2002
 [8] Bierer, S. M., and Anderson, D. J., "Multi-channel Spike Detection and Sorting using an Array Processing Technique," *Neurocomputing*, vol. 26-27, pp. 947-956, 1999
 [9] Oweiss, K. G., and Anderson, D. J., "A New Approach to Array Denoising," *Proc. of the IEEE 34th Asilomar Conference on Signals, Systems and Computers*, pp. 1403-1407, November 2000.
 [10] Oweiss, K. G., and Anderson, D. J., "A new technique for blind source separation using subband subspace analysis in correlated multichannel signal environments," *Proc. of ICASSP'2001*, pp. 2813-2816, May 2001
 [11] Moon, T.K., and Stirling, W. C., *Mathematical Methods and Algorithms for Signal Processing*, Prentice Hall, New Jersey, 1st edition, 2000
 [12] Oweiss, K.G., and Anderson, D.J., " A Multiresolution Generalized Maximum Likelihood Approach for the Detection of Unknown Transient Multichannel Signals in Colored Noise with Unknown Covariance," *Proc. of ICASSP*, vol. 3, pp.2993-2996, May 2002
 [13] Cover, M., and Thomas, Joy A.: *Elements of Information Theory*, John Wiley & Sons, 1991

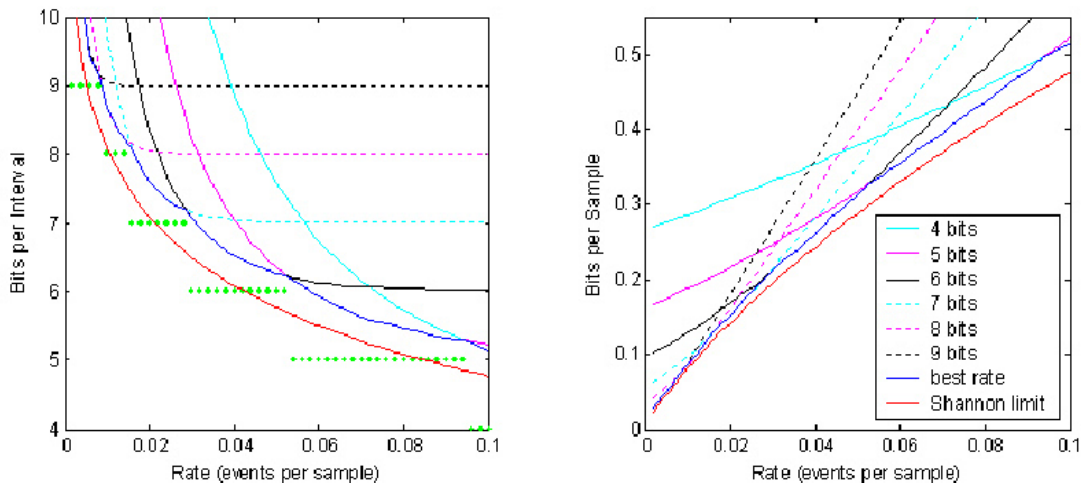


Figure 4: For exponential distributions of event intervals, Shannon theory predicts the optimum number of bits per interval is $\cong \log_2(e/\lambda)$. The run length algorithm with an optimum word length (best rate) exceeds the Shannon limit [13] by nearly a bit per interval while if the word length is badly chosen, the error grows. The data are shown for both bits/event and bits/sample. For the event rate ranges experienced for neural data, the overall bit rate is under .2 bits/sample.