



ISSN:2229-6107



**INTERNATIONAL JOURNAL OF
PURE AND APPLIED SCIENCE & TECHNOLOGY**

E-mail :
editor.ijpast@gmail.com
editor.ijpast@.in

www.ijpast.in

A Scalable Wavelet Transform VLSI Architecture for Real-Time Signal Processing in High-Density Intra-Cortical Implants

Mr. JADA LINGAIAH, P. SATHISH BABU, Mrs. GUNDUBOINA ANUSUYA, BONGU CHANDINI

Abstract—This paper describes an area and power-efficient VLSI approach for implementing the discrete wavelet transform on streaming multielectrode neurophysiological data in real time. The VLSI implementation is based on the lifting scheme for wavelet computation using the symmlet4 basis with quantized coefficients and integer fixed-point data precision to minimize hardware demands. The proposed design is driven by the need to compress neural signals recorded with high-density microelectrode arrays implanted in the cortex prior to data telemetry. Our results indicate that signal integrity is not compromised by quantization down to 5-bit filter coefficient and 10-bit data precision at intermediate stages. Furthermore, results from analog simulation and modeling show that a hardware-minimized computational core executing filter steps sequentially is advantageous over the pipeline approach commonly used in DWT implementations. The design is compared to that of a B-spline approach that minimizes the number of multipliers at the expense of increasing the number of adders. The performance demonstrates that *in vivo* real-time DWT computation is feasible prior to data telemetry, permitting large savings in bandwidth requirements and communication costs given the severe limitations on size, energy consumption and power dissipation of an implantable device.

Index Terms—B-spline, brain machine interface, lifting, micro-electrode arrays, neural signal processing, neuroprosthetic devices, wavelet transform.

INTRODUCTION

VLSI implementation of the discrete wavelet transform (DWT) has been widely explored in the literature as a result of the transform efficiency and applicability to a wide range of signals, particularly image and video [1], [2]. These implementations are generally driven by the need to fulfill certain characteristics such as regularity, smoothness and linear

Manuscript received August 16, 2006, revised December 11, 2006. This work was supported by the National Institutes of Health (NIH) under Grant NS047516. This paper was recommended by Associate Editor A. Van Schaik.

phase of the scaling and wavelet filters, as well as perfect reconstruction of the decomposed signals [3].

In some applications, it is desirable to meet certain design criteria for VLSI implementation to enhance the overall system performance. For example, minimizing area and energy consumption of the DWT chip is highly desirable in wireless sensor network applications where resources are very scarce. In addition to miniaturized size, minimizing power dissipation is strongly sought to minimize tissue

heating in some biomedical applications where the chip needs to be implanted subcutaneously.

In this paper, we deal primarily with the design of DWT VLSI architecture for an intracortical implant application. Motivated by recent advances in microfabrication technology, hundreds of microelectrodes can be feasibly implanted in the vicinity of small populations of neurons in the cortex [4], [5], opening new avenues for neuroscience research to unveil many mysteries about the connectivity and functionality of the nervous system at the single cell and population levels. Recent studies have shown that the activity of ensembles of cortical neurons monitored with these devices carry important information that can be used to extract control signals to drive neuroprosthetic limbs, thereby improving the lifestyle of severely

ASSISTANT PROFESSOR^{1,2,3}, STUDENT⁴

Department of ECE

Arjun College Of Technology & Sciences

Approved by AICTE & Affiliated to JNTUH

SPONSORED BY BRILLIANT BELLS EDUCATIONAL SCOTEY

paralyzed patients [6]–[8].

One particular challenge with the implant technology is the need to transmit the *ultra-high* bandwidth neural data to the outside world for further analysis. For example, a typical recording experiment with a 100 microelectrode array sampled at 25 kHz per channel with 12-bit precision yields an aggregate data rate of 30 Mbps which is well beyond the reach of state-of-the-art wireless telemetry. Other significant challenges consist of the need to fit circuitry within cm for the entire signal processing system, and operate the chip at very low power (no more than 8–10 mW) to prevent temperature rise above 1 C that may cause neural tissue damage. In previous studies, we have shown that the DWT enables efficient compression of the neural data while maintaining high signal fidelity [9]–[11]. To be implemented in an actual implanted device, chip size, computational complexity and signal fidelity must be balanced to create an optimal application-specific integrated circuit (ASIC) design tailored to this application. Generally speaking, the case of computing the DWT for high throughput streaming data has not been fully explored [12]. It has been argued that a lifting scheme [13] provides the fewest arithmetic operations and in-place computations, allowing larger savings in power consumption but at the expense of

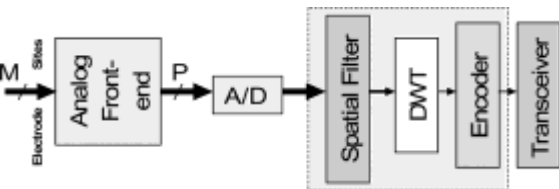


Fig. 1. Block diagram of an implantable neural system illustrating the mixed signal processing proposed.

longer critical path than that of convolution-based ones [13]. Recent work by Huang *et al.* [14] focused on analyzing DWT architectures with respect to tradeoffs between critical path and internal buffer implementations. Such critical path can be shortened using pipelining with additional registers or using a so-called *flipping* structure with fixed number of registers [15]. The *B-spline* approach [16], on the other hand, requires fewer multipliers than lifting, replacing them with adders that may permit a smaller chip area [17]. Nonetheless, most of the reported hardware approaches focus on computational speed and do not adequately address severe power and area constraints. By comparing with other implementations of the DWT in this paper, we demonstrate that the appropriate compromise among power, size and speed of computations is achieved with a sequential implementation of integer arithmetic lifting approach. The paper is organized as follows. In Section II, the classical single channel one-dimensional (1-D) DWT and lifting DWT are introduced. Section III describes the motivation for integer lifting DWT and approaches to efficiently map the algorithm to hardware for a single channel, single level DWT decomposition. In Section IV, proposed architectures for integer lifting are described and analysed. Section V describes hardware considerations of the proposed architecture for multiple channels and multiple L levels of decomposition, and Section VI describes performance comparisons and overall results.

THEORY

A typical state-of-the-art implantable neural interface systems depicted in Fig. 1 contains an analog front end consisting of pre-amplification, multiplexing and A/D

conversion prior to extra-cutaneous transmission. An analog front end integrated onto a 64-electrode array would occupy 4.3 mm in 3 m technology and would dissipate 0.8 mW of power [5]. This traditional approach is not well suited for wireless data transmission due to power demands associated with the resulting large data throughput. In the proposed approach, the power and chip area of the analog front end is reduced by using contemporary mixed-signal VLSI design approaches and more modern fabrication processes (e.g., 0.18 m), allowing advanced signal processing to take place within the implanted system without significant increase in the chip size. Power- and area-efficient implementations of the spatial filter, the DWT, and the encoder blocks would provide on-chip signal processing and data compression, enabling wireless transmission by reducing bandwidth requirements. In this paper, we only discuss VLSI implementation of the DWT block, mostly contained in the short transients -or spikes- above the noise level that result from the activity of an unknown number of neurons. It can be observed that the sparsity introduced by the DWT compaction property enables very few “large” coefficients to capture most of the spikes’ energy, while leaving many “small” coefficients attributed to noise. This property permits the later ones to be thresholded [19], yielding the denoised signal shown.

For near-optimal data compression, a wavelet basis needs to be selected to best approximate the neural signal waveform with the minimal number of data coefficients. A compromise between signal fidelity and ease of hardware implementation has to be made. A near-optimal choice was proposed in [9] from a compression standpoint and demonstrated that the *biorthogonal* and the *symmlet4* wavelet functions are advantageous over other wavelet basis families for processing neural signals. From a hardware implementation viewpoint, the *symmlet4* family has much smaller support size for similar number of vanishing moments compared to the biorthogonal basis [20]. In addition, they can be implemented in operations.

A. Single Channel Lifting-Based Wavelet Transform

The lifting scheme [12] illustrated in Fig. 3 is an alternative approach to computing the DWT. It is based on three steps: First, splitting the data at level j into *even* and *odd* samples and , respectively; Second, predicting the odd samples from the even samples such that the prediction error becomes the high pass coefficients ; and third, updating the even samples with to obtain the approximation coefficients . This process is repeated times. At an arbitrary prediction and update step , the prediction and update filters $T_n(z)$ and , respectively, are obtained by factorizing the wavelet filters details and of the next level.

A lifting factorization of the *symmlet4* wavelet basis amounts to the following filtering steps:

$$\begin{aligned}
\text{symlet2: } f_1[i] &= f_0[i] + C_3 \times g_1[i] + C_4 \times g_1[i+1] \\
\text{symlet3: } g_2[i] &= g_1[i] + C_4 \times f_1[i] + C_5 \times f_1[i-1] \\
\text{symlet4: } a[i] &= f_1[i] + C_5 \times g_2[i] + C_7 \times g_2[i-1]
\end{aligned} \tag{7}$$

Alternatively, a B-spline approach for DWT computation [16] is based on factorizing the filters as

$$H(z) = h_a \tag{8}$$

where $Q(z)$ and $R(z)$ are known as the distributed parts, h_a and l_n are normalization factors [17], and are the orders of the B-spline parts, respectively. For the *symmlet4*, this factorization can be expressed as

$$H(z) \tag{9}$$

where the coefficients B_1 through B_6 are listed in Table II. Since the B-spline parts in both filters can be expressed as

$$+ z^{-4} \tag{10}$$

they can be typically implemented using simple shifting and addition. The polyphase decomposition similar to lifting can therefore be performed on the distributed parts $R(z)$ and $Q(z)$ [16]. This is achieved by splitting the distributed parts into odd and even components $R_e(z)$ and $R_o(z)$, $Q_e(z)$ and $Q_o(z)$, respectively. For example, the low-pass even distributed part can be represented as $R_e(z) = B_2 z^{-2}$, and likewise for the remaining components. The benefit in the B-spline method is a reduction in the number of floating point multiplications at the expense of more additions [17]. Table III compares the computational requirements of lifting and B-spline DWT implementations along with traditional convolution. In B-spline, four x4 multiplications are replaced by shifts and two x6 multiplications are replaced by shifts and additions. Relative to lifting, B-spline requires two fewer multiplications at the expense of ten more additions for one level of decomposition. Nevertheless, as the detailed low-power/area DWT implementation below will show, any benefit to B-spline is diminished for multilevel multichannel decomposition.

B. Hardware Considerations

Power and area requirements of the DWT hardware are determined largely by the complexity of the computational circuitry and the required memory. To systematically reduce hardware requirements, we have explored different options to reduce computation and memory requirements at the algorithm level and analyzed their impact on signal integrity to determine an optimal approach. We summarize below two key ideas that contribute largely to the reduction of circuit complexity and memory requirements that are discussed in subsequent sections, while more details of this analysis are further provided in Section V.

1) *Integer Approximation:* Fixed-point integer approximation limits the range and precision of data values but greatly reduces the computational demand and memory requirements for processing and storage. To explore the potential of utilizing integer approximation in the proposed

system, we observed that neural signal data will be entering the system through an A/D converter and will thus inherently be integer valued within a pre-scribed range. The data is first scaled to obtain data samples within a 10-bit integer precision. The integer approximation is then computed for the scaled data. The integer-to-integer transformation [22] involves rounding-off the result of the lifting filters and S

$$L(z) = (1 + z^{-1})^2 \cdot \frac{1}{2} \cdot l_n$$

that are used to filter odd and even data samples, respectively. The last step that requires scaling by $1/2$ and G^{-1} is omitted. Hence, the dynamic range of the transform at each level will now change by G . As our results will demonstrate (Section V), the minimized circuit complexity

$$L(z) = (1 + z^{-1})^4 (1 + B_4 z^{-1} + B_3 z^{-2} + B_2 z^{-3})$$

associated with integer representation should be well suited to this application provided that data precision is sufficient to maintain signal integrity.

Quantization of the Filter Coefficients: Rounding-off wavelet filter coefficient values to yield a fixed point integer precision format can further reduce the computation and memory requirements. Implementing lifting-based wavelet transform with only integer computational hardware requires the filter coefficients be represented as integers along with the sampled data. Tables I and II show the scaled filter coefficients $C_1 - C_8$ (x16) and $B_1 - B_6$ (x2) for the *symmlet4* basis. These coefficients are further quantized into integer values. The level of quantization has a significant impact on the complexity of computational hardware. We quantified the effect of the roundoff and quantization errors on the signal fidelity as a function of multiplier complexity [21]. Our results (Section V) demonstrate that 6 bits (5 bits 1 sign bit) coefficient quantization can adequately preserve signal integrity.

I. SINGLE-CHANNEL SINGLE-LEVEL HARDWARE DESIGN

In a first-order analysis, the area of a CMOS integrated circuit is proportional to the number of transistors required, and power consumption is proportional to the product of the number of transistors and the clocking frequency. Through transistor-level custom circuit design, circuit area and power consumption can be further reduced, with significant improvement in efficiency over field-programmable gate arrays (FPGA) or standard cell ASIC implementations.

Parallel execution of the DWT filter steps using a pipelined implementation is known to provide efficient hardware utilization and fast computation. In fact, a vast majority of the reported hardware implementations for lifting-based DWT rely on pipeline structures [20], [23], [24]. However, these circuits target image and video applications where speed has highest priority and the wavelet basis is chosen to optimize signal representation. A different approach is required to meet the power and area constraints imposed by implantability requirements, the low bandwidth of neural signals, and the type of signals observed. Two promising integer lifting DWT implementations, a pipeline approach and a sequential scheme, have been optimized and compared for the *symmlet4* factorization and data/coefficient quantization described above. Furthermore, the hardware requirements for lifting DWT

have been compared to a B-spline implementation to verify the advantage of lifting in the application at hand.

A. Computation Core Design

To begin, notice that the arithmetic operations in the lifting scheme in (7) have a noticeable regularity that permits any arbitrary step to be defined as

$$O_{ij} = a_i \cdot f^j + b_j \cdot O_{i-1}^j \quad (11)$$

where O_{ij} , a_i , b_j , and f^j take the values of O_1^1 , O_2^1 , O_3^1 , and O_4^1 in (7), and a_i and b_j are the quantized filter coefficients given in Table I. The regularity of this repeated operation indicates that an optimized integer DWT implementation would include a hardware unit specifically designed to evaluate (11). By tailoring this circuit to the near-optimal data and coefficient bit width described above, a single computation core (CC) suitable for all lifting filter steps in (7) can be obtained.

Fig. 4 describes a CC block that was custom designed to minimize transistor count and power consumption while supporting up to 10-bit data and 6-bit filter coefficients, both in signed integer formats. The CC employs a simple hardwired shifting operation to

$$1: O_1^3 = f^3 + a_1 \cdot f^1$$

remove the x16 scaling factor from the quantized coefficients. It generates a 10-bit output and an overflow error bit, though the lifting scheme should inherently maintain results within 10-bit magnitude. Several multiplier topologies were experimentally compared to define the most efficient option for

6 10-bit operations. A Wallace tree multiplier with modified Booth recoding was implemented along with a custom 3-term adder optimized for power rather than speed. The fixed x16 scaled integer coefficients were modified for Booth recording before being stored in on-chip ROM to eliminate the need for

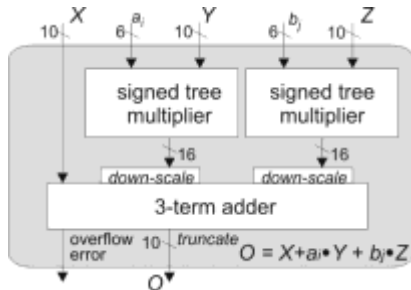


Fig. 4. Customized computation core for integer-lifting wavelet transform using binary scaled filter coefficients.

an on-chip encoder. The resulting circuit very efficiently implements steps 2-4 of (7) and can also compute steps 1 and 5 using a control signal that shuts off the unused multiplier to eliminate unnecessary power consumption.

B. Real-Time Integer DWT Processing Architectures

To identify the most efficient architecture for executing the entire set of lifting equations in real time on a continuous flow of input data samples, let us first re-define the filter equations in (7) with a more hardware-friendly notation. Building on the concept of a fixed three-term computation core described above, the notation in (11) can be used to rewrite (7) at a specific cycle,

$$\begin{aligned}
 \text{Step 1: } & O_1^1 = f^1 + a_1 \cdot f^0 \\
 \text{Step 2: } & O_2^1 = O_1^1 + a_2 \cdot O_1^0 + b_1 \cdot O_1^0 \\
 \text{Step 3: } & O_3^1 = O_2^1 + a_3 \cdot O_2^0 + b_2 \cdot O_2^0 \\
 \text{Step 4: } & O_4^1 = O_3^1 + a_4 \cdot O_3^0 + b_3 \cdot O_3^0 \\
 \text{Step 5: } & O_5^1 = O_4^1 + a_5 \cdot O_4^0 + b_4 \cdot O_4^0
 \end{aligned} \quad (12)$$

where O_1^1 and f^1 are the input data pair of samples, the outputs of steps 1-5 are O_1^1 - O_5^1 , coefficients C_1 - C_5 have been replaced by a_1 - a_5 and b_1 - b_4 to indicate the CC input to which they will be applied, and the superscripts represent the computation cycle in which the data value was generated. The 2nd and 3rd terms in step 2 have been swapped to maintain a regular data flow described further below. Steps 2 and 5 require data from future computation cycles. Thus, in order to compute the five filter steps in real time, where all inputs must be available from prior computations, execution must span *three* computation cycles. During cycle t the following five steps can be executed in real time:

$$\begin{aligned}
 1: & O_1^3 = f^3 + a_1 \cdot f^1 + b_1 \cdot O_1^2 \\
 2: & O_2^3 = O_1^3 + a_2 \cdot O_1^2 + b_2 \cdot O_2^2 \\
 3: & O_3^3 = O_2^3 + a_3 \cdot O_2^2 + b_3 \cdot O_3^2 \\
 4: & O_4^3 = O_3^3 + a_4 \cdot O_3^2 + b_4 \cdot O_4^2 \\
 5: & O_5^3 = O_4^3 + a_5 \cdot O_4^2 + b_5 \cdot O_5^2
 \end{aligned} \quad (13)$$

Notice that each step in (13) relies only on previously calculated data, provided these steps are performed *sequentially*. Having rearranged the terms in step 2 of (7), the output of each step in (13) becomes the 2nd term input to the subsequent step, which is useful for efficient hardware implementation. Notice also that most of the data values needed are generated within the same

cycle; only the four values in (13) with boldface type (two are repeated twice) are generated in a previous cycle. Thus, if the filter steps are implemented sequentially, only *four* storage/delay registers are required.

Although (13) does allow real time computation of the filter steps in sequence, dependencies within the steps in (13) preclude parallel execution necessary for a pipeline implementation. To make each filter step dependent only on data from prior cycles, execution must span *seven* data samples. During cycle t the following sequence could be computed without any dependency on current or future cycle results:

$$\begin{aligned}
 O_1^7 &= f^7 + a_1 \cdot f^1 \\
 O_2^7 &= f^3 + a_2 \cdot O_1^6 + b_1 \cdot O_1^5
 \end{aligned}$$

$$\begin{aligned}
 O3^4 - O1^4 + a_3 O2^4 + b_3 O2^3 \\
 O4^3 - O2^3 + a_4 O3^3 + b_4 O3^2 \\
 O5^1 \quad O4^2.
 \end{aligned} \tag{14}$$

Here, the second term of each computation relies on the output from the preceding step during the previous computation cycle. In a pipeline, these four second-term data inputs could be held in memory with one-cycle delay. The first and third terms require seven additional data values from prior cycles, one of which is needed twice, resulting in six independent values. One of the values (f^5 in step 2) needs a two-cycle delay, requiring an extra delay register. Thus, a total of 11 storage/delay registers would be required to hold all of the necessary values from prior cycles for a pipeline implementation.

C. Pipeline Design

The integer DWT filter equations in (14) can be implemented simultaneously in a pipeline structure that permits real time, continuous signal processing to take place. Fig. 5(a) illustrates a pipeline structure designed around the customized three-term computation core from Fig. 4. The output of each of the five filter stages is held by a darkly shaded pipeline register, and other registers provide the necessary delays. By clocking all of the registers out of phase from the CC blocks, continuous operation is provided. The computation latency is seven cycles, due to the five pipeline stages and the two delay cycles built into (14). The temporal latency for detail and approximation results is 14 samples because each computation cycle operates on a pair of data samples. The overall pipelined computational node consists of five CC blocks, 15 10-bit registers, and an 8 6b coefficient ROM. An additional delay phase could be added at the $O4$ output to synchronize the latency of the detail and approximation outputs.

D. Sequential Design

Although the pipeline structure achieves fast integer DWT processing via a large hardware overhead, it is very resource-efficient and thus well suited for low-power, single channel, neural signal processing. However, as discussed below, scaling the pipeline for multiple data channels and/or multiple decomposition levels begins to break down the efficiency of the pipeline structure. An alternative approach is to process each of the filter steps (or pipeline stages) sequentially using a single CC

(O4)
x

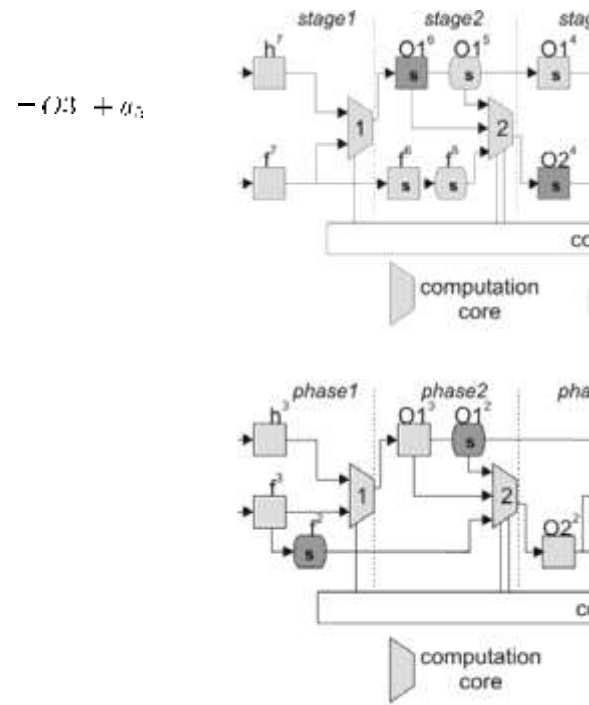


Fig. 5. (a) Pipeline structure for integer-lifting wavelet transform with data notations to match filter equations in (11) at a single point in time. (b) Sequential structure over five operation phases for comparison to the pipeline structure.

block and a fraction of the registers required by the pipeline. This approach takes advantage of the low bandwidth of neural signals that permits the CC to be clocked much faster than the input data sampling frequency (typically in the range of 25–40kHz).

Sequential processing of the integer DWT filter steps can be achieved using (13), where each stage depends only on data from previous cycles or from same-cycle outputs generated in a preceding step. The simplicity of data dependencies relative to the pipeline structure can be observed from Fig. 5(b), which illustrates the sequential structure in a format comparable to the pipeline. Here, each section of the circuit represents a temporal phase rather than a physical stage. An important observation is that significantly fewer registers are needed because the inputs of subsequent phases rely largely on preceding outputs from the same computation cycle. Therefore, it can be shown that the overall sequential DWT circuit can be efficiently implemented with six 10-bit registers to manage data flow between computation cycles, a single CC block, an 8 6b coefficient ROM, and a simple control block to direct data from memory to the appropriate CC input during each phase of operation. Sequential execution has a computation latency of two cycles, and the temporal latency for detail and approximation results is four samples.

E. Analysis and Comparison

As stated above, the sequential approach requires only one CC unit and six 10-bit memory registers compared to five CC units and 15 registers for the pipeline circuit. The sequential design does, however, require additional multiplexers and control logic to redirect data and coefficients to CC inputs, which are not necessary in the inherently

hardware-efficient pipeline design. This added circuitry will make the critical path of the sequential circuit longer than that of the pipeline structure. Furthermore, to maintain the same throughput, the sequential design must be operated at five times the clock rate of the pipeline. Because data is processed in a real-time streaming mode, neither approach requires a large input data buffer.

Both architectures have been thoroughly analyzed to determine which approach is best suited to the power and area requirements of an implantable neural signal processor. To first validate that both approaches can achieve the application speed requirements, a custom computation core has been implemented in CMOS, and analog simulations show the critical path delay is 6.5 ns in 0.5- μm technology. Thus, approximately 6000 computation cycles could be performed within a

custom circuit layouts shows that a single value for reasonably approximates all of the integer DWT blocks, especially for comparing two similar circuits. Conservative values of 80 μm^2 per transistor for 0.5- μm technology and 5 μm^2 per transistor for 0.13- μm technology have been selected to estimate the required chip real estate.

TABLE IV
CHARACTERISTICS OF SINGLE-LEVEL, SINGLE-CHANNEL INTEGER DWT HARDWARE FOR PIPELINE AND SEQUENTIAL CONFIGURATIONS AT TWO TECHNOLOGY NODES

	# transistors	0.5 μm CMOS		0.13 μm CMOS	
		area [mm ²]	power [μW]	area [mm ²]	power [μW]
Pipeline	23,336	1.867	15.10	0.117	0.71
Sequential	7,931	0.634	18.55	0.040	0.87

Although absolute power consumption is inherently difficult to estimate, for the purpose of comparing the two design alternatives, dynamic power can be determined as

$$P = \frac{1}{2} VDD^2 f_s \sum_i N_i s_i \quad (16)$$

where VDD is the supply voltage and f_s is the data sampling frequency (nominally 25 kHz). The parameter k accounts for the average output load capacitance, the average number of transistors per output transition, and the average output transitions per clock cycle. This parameter is a function of both fabrication process and circuit topology and has been derived empirically as 3 and 0.75 fF for 0.5- μm and 0.13- μm technology, respectively. The variable s_i is the clock rate scaling factor relative to f_s for each block such that the clocking frequency of each circuit block is $s_i f_s$. For example, in the pipeline configuration, the computation core will be clocked only every other cycle, i.e., $1/2 f_s$, so that the first of the pair of samples to be processed can be acquired in the idle cycle. Correspondingly, because the sequential configuration must be clocked at five times the rate of the pipeline, it will have an average clocking rate of $5 f_s$. In the pipeline approach, all of the blocks are

nominally 25-kHz sampling frequency for neural signals. This indicates that speed is not a critical design constraint and that circuit optimization can focus on chip area and power consumption.

Using custom design techniques, the chip area, A , required to implement both approaches will be roughly proportional to the number of transistors in the circuit

$$A = \sum_i T_i N_i \quad (15)$$

where T_i is the area per transistor and N_i is the number of transistors in the i th circuit block. Empirical observations of several

clocked at the same frequency, except the coefficient memory, that is static in both designs. In the sequential implementation, one of the multipliers is idle during two of the five stages, so we estimate the sequential CC clock scaling factor to be 2. Similarly, in the sequential controller, most of the circuits are clocked at $5 f_s$ while others are clocked at f_s , so we estimate its clock scaling factor to be 2 as well.

Table IV lists the total number of transistors in each approach along with the area and power estimated from (15) and (16) for both 0.5- μm and 0.13- μm technology. As expected, the pipeline computation unit requires nearly three times the area of the sequential approach and would occupy about 21% of the chip area on a 3.3 mm chip in 0.5- μm technology or 5% of a 1.5 mm chip in a 0.13- μm process. The power model predicts that the sequential approach will consume only 23% more power than the pipeline. The larger power consumption of the sequential approach can be attributed to its requirement for a more complex controller and the need to move more data around within the single computation core. Overall, these results show a tradeoff between area and power consumption between the two approaches.

F. Lifting Versus B-Spline

As an alternative to lifting, the B-spline method was investigated because it permits a reduction in the number of floating

point multiplications at the expense of more additions. However, as demonstrated above, for implantable applications, integer processing is preferred. Table III shows that B-spline saves two multiplications at the cost of 10 additions per cycle compared to lifting. Designs using Verilog synthesized to a custom library have shown that, for a pipeline implementation, B-spline requires significantly less 24-bit floating point hardware, but for integer processing (with 10-bit data and 6-bit coefficients) B-spline saves only 6% compared to lifting [25]. Furthermore, B-spline can not be as efficiently implemented in a sequential structure, where lifting has been shown to require only 53% of the B-spline hardware resources for integer DWT. While B-spline

implementations do have slightly less delay, speed is not a design constraint. Relative memory requirements are a

more important issue in multichannel implementations as we show next.

II. MULTILEVEL AND MULTICHANNEL INTEGER DWTIMPLEMENTATION

A. Hardware Design

In implantable neuroprosthetic applications where a typical microelectrode array has many electrodes integrated on a single device, there is a strong need to support integer DWT computations with multiple levels of decomposition for multiple signal channels pseudo-simultaneously (i.e., within one sampling period). The lifting scheme and the two integer DWT implementations described above have been chosen because of their ability to scale to an arbitrary number of channels and levels. Considering that both of the single channel, single level, integer DWT approaches discussed above require a substantial portion of a small chip, it is unreasonable to pursue a hardware intensive solution that utilizes a “copy” of the circuit for each channel and level. This would dramatically increase circuit area beyond limitations for implantable systems. Given the available computation bandwidth of the CC block, the more appropriate solution is to scale the clocking frequency as needed to sequentially compute filter equations for multiple channels and/or levels. Although clock scaling will still cause power to increase with channel and level, the circuit area required will be minimized and the power density can be held within the acceptable application limits.

Both the pipeline and sequential architectures can be scaled to multiple channels and/or levels by reusing the computational node hardware and increasing the clocking frequency to complete all computations within the input sample period. In both approaches, registers within the computational node hold data necessary for the next cycle’s calculation. To sequentially reuse the computational node, some register values for a specific

channel/level must be saved so they will be available when that channel/level is next processed in a future cycle. Fig. 6 shows the multichannel, multilevel, implementations of the pipeline and sequential configurations.

1) *Multichannel Considerations:* In scaling the system to multiple data channels, the computation clock rate is scaled by the number of channels and a new memory block is added to save critical register data for each channel. For the pipeline, the 11 registers must be stored, while for the sequential circuit only

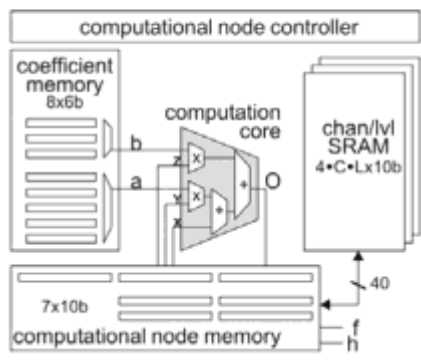
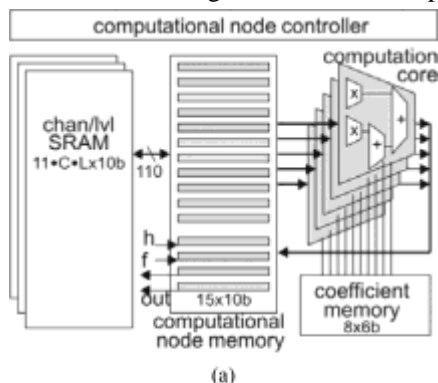


Fig. 6. Multilevel, multichannel implementations of (a) pipeline structure and (b) sequential structure.

four registers need to be saved. These registers are marked with “s” in Fig. 4. An on-chip SRAM can be interfaced to the computational node to store register values, and the size of the SRAM will grow linearly with the number of channels. Note for comparison that a sequential B-spline implementation requires eight register values to be stored.

2) *Multilevel Considerations:* When expanding the DWT to multiple levels, notice that each level of dyadic DWT decomposition introduces only half the number of computations as the previous level. More explicitly, the number of results, n , per number of samples, N , for an arbitrary level j can be expressed as

$$\frac{N}{n} = \sum_{j=1}^L \left(\frac{1}{2}\right)^j \quad (17)$$

which is always less than twice the number of samples.

Consider also that, to process multichannel input pairs, before each computation cycle the system must implement one idle cycle, wherein the first input of the pair is stored for each channel. Thus, if the level-one computations are executed in, say, the even cycles, the higher level computations can be executed in the odd cycles [26] while input samples (one of the pair) are being stored for the next level-one computation. This is illustrated in Fig. 7. If we define the

B. Area and Power Modeling

For multiple channels/levels, the need to copy the entire set of pipeline registers to memory effectively negates one of the primary advantages of the pipeline over the sequential approach. On the other hand, the sequential processing circuit is inherently designed to swap new data in/out each clock cycle. To quantitatively compare these two approaches, circuit models have been developed to describe the power and area for each option as a function of the number of channels and the number of decomposition levels. The following models assume the hardware (including control logic) has been scaled to manage multiple channels and levels, though they are still valid for single channel, single level implementations.

analysis since both approaches consist of similar arithmetic and memory blocks.

Using (16), a general expression for power consumption as a function of channels and levels, which is valid for both approaches being considered, is given by

$$P = 2 \cdot k \cdot VDD^2 \cdot f_s \sum_i \left(N_{oi} + (C-1) \cdot L \cdot N_{ci} + (L-1) \cdot C \cdot N_{li}(Cf_i \cdot Lw_i \cdot s_i) \right) \quad (20)$$

where Cf_i is the channel clock frequency scaling factors, Lw_i is a level usage factor, and all other variables are as previously

usage rate, $U(L)$, as the average number of cycles for a single computation to occur, then for the first decomposition level the usage rate is one half, i.e., $U(1) = 0.5$, and the computational hardware is idle during the other half of the cycles. Moreover, $U(L)$ approaches 1.0 as the number of levels increase, i.e.,

$U(L) = 1 - 2^{-L}$ Fig. 7. Sequential processing scheme for multilevel, multichannel computation. At the top of this sequence, one DWT result is available at each decomposition level. With the four levels shown, one idle computation cycle will occur every 16 cycles.

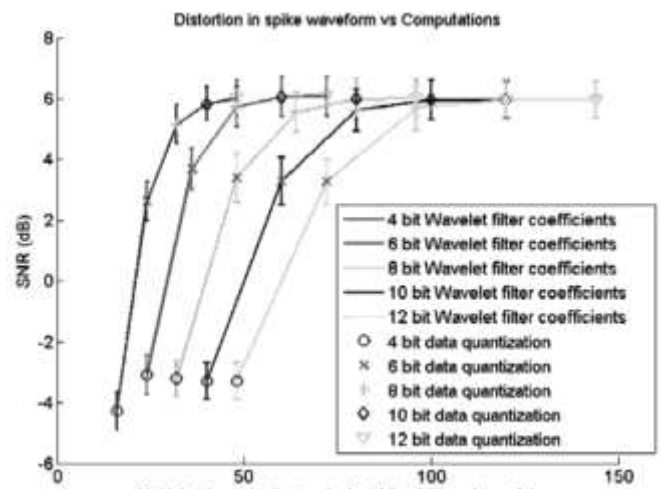
As the number of levels increases, the usage rate will increase toward maximum utilization without increasing computation frequency. For each level of decomposition beyond the first, one memory block per channel is required to store values held in the computational node registers. The registers to be stored are the same as those described in the multichannel case above.

- 0.5

A general expression for calculating the area of both the pipeline and the sequential approaches as a function of channels and levels is:

$$A = A_t \sum_i \left(N_{oi} + (C-1) \cdot L \cdot N_{ci} + (L-1) \cdot C \cdot N_{li} \right) \quad (19)$$

where A_t is the technology-dependent, empirically-derived average area per transistor, N_{oi} is the number of transistors that remain constant with level and channel in the i th circuit block, and N_{li} are the number of transistors that scale with channel and level, respectively, C is the number of channels, and



defined. Recall that the clock scaling factor was chosen to accommodate the fact that, in single level designs, every other cycle was idle while the data pair was being collected. To maintain a consistent definition of variables in multilevel implementations, which utilize the idle cycles to process all higher levels, the factor of 2 is introduced at the beginning of (20).

Both the pipeline and sequential architectures have been developed to define the model parameters given in Table V, which are valid for $T_c \gg 1$ and $C_c \gg 1$. The computational node circuitry, including control logic, has been scaled up to manage an arbitrary number of levels and channels, with negligible per channel/level increase in complexity. Thus, only data memory increases with the number of channels. Clocking frequency of the computational node circuits must scale with channel, while each memory block is only accessed once per cycle regardless of the number of channels. The controller frequency scales linearly with channel but is assumed to remain constant with level. For all other circuit blocks, the usage rate $U(L)$ accounts for inactive computation cycles.

V. RESULTS AND DISCUSSION

A. Signal Integrity

We have assessed the effects of data and filter coefficient ap-proximations on the quality of the signals obtained after recon-struction. We quantified the performance in terms of the com-plexity of hardware required to implement (7) and illustrated the results in Fig. 8. The wavelet filter coefficients were quantized to different resolutions ranging from 4 to 12 bits, with the 6-bit values given in Table I. The data was also quantized in the same range. The *effective* signal-to-noise ratio (eSNR), defined as the log ratio in dB of the peak spike power to the background noise power is illustrated in Fig. 8(a) versus multiplier complexity in equivalent bit addition/sample for an average input SNR of 6dB. These results demonstrate that, with sufficient precision, the use of integer computations does not result in significant signal degradation as quantified by the observed output SNR. Specifically, with quantization of filter

coefficients to 6 bits and data to 10 bits, the output SNR is within 1% of its average input value. In Fig. 8(b), the spectrum of the residual quantization and round-off noise is also illustrated to demonstrate the loss in the signal power-spectral density in different cases. In the case of 4-bit quantization of the filter coefficients, the residual noise frequency content is closest to that of the original signal in the low frequency range (subband 0–1 kHz), indicating that some signal loss may have occurred in that band. On the other hand, Fig. 8. (a) Effect of round off and quantization errors on the signal fidelity as a function of multiplier complexity. (b) Power-spectral density of the original data and the residual noise for integer approximated data and quantized wavelet filter coefficients for various bit widths. (c) Example spike waveforms obtained in each case.

filter quantization of 6 bit or higher results in residual noise that consists of high frequency components above 8 kHz, which is outside

the frequency range of neural spike trains and local field

TABLE V

MODEL PARAMETERS FOR AREA AND POWER CALCULATIONS

	N_o	N_c	N_l	C_f	L_u	s
Pipeline						
Controller	524	0	0	C	1	0.5
Comp. Block	2184 8	0	0	C	U(L)	0.5
Coeff. Mem.	964	0	0	C	U(L)	0
Data Mem.	0	1562	1562	1	U(L)	0.5
Sequential						
Controller	744	0	0	C	1	2*
Comp. Block	6125	0	0	C	U(L)	2*
Coeff. Mem.	1062	0	0	C	U(L)	0
Data Mem.	0	572	572	1	U(L)	0.5

* accounts for average usage at 2.5 clock scale

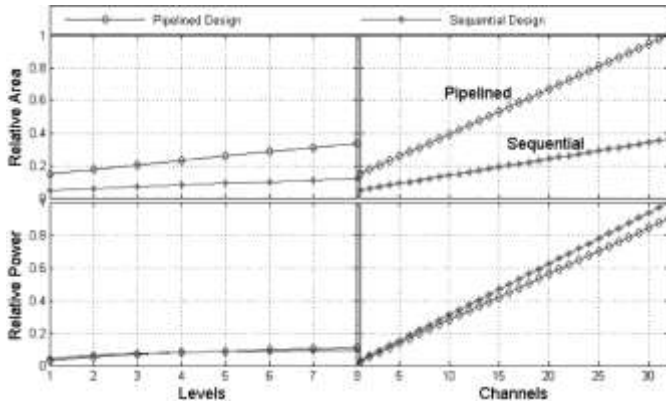


Fig. 9. Comparison of multichannel/multilevel pipeline and sequential integer DWT approaches: relative chip area and relative power consumption versus number of levels and channels.

potentials (LFPs) [27]. A representative example of spike wave-forms in each case is illustrated in Fig. 8(c) to demonstrate the very negligible effect of this process on the quality of the average spike waveform. Taking these results all together, it is clear that the choice of 6/10-bit coefficient/data quantization offers the best compromise among multiplier complexity and signal fidelity as concluded earlier.

We should emphasize that perfect reconstruction of signals off chip may not be always needed. Typically, neural signals contain the activity of multiple neurons that need to be *sorted out*, and this information remains in the compressed data at the output of the DWT block. We have shown elsewhere that sorting the multi source neuronal signals can be performed directly on the wavelet transformed data [10], [28], and this topic is outside the scope of this paper.

B. Multichannel/Level Implementations

Using (19) and Table V, the relative area for pipeline and sequential architectures as a function of levels and channels is shown in Fig. 9. These results demonstrate that the pipeline requires significantly more chip area than the sequential approach and its area needs grow faster with larger number of channels and levels. This is due primarily to the relatively large number of registers that must be stored per channel or level (11 for pipeline compared to 4 for sequential). Fig. 9 also shows the relative

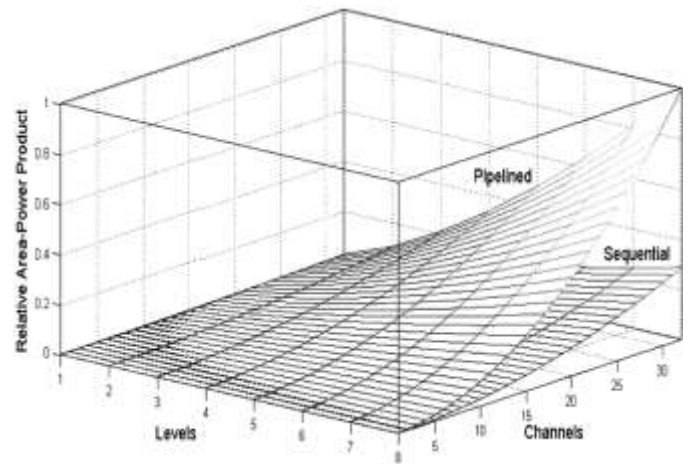


Fig. 10. Power-area product versus level and channel for pipeline and sequential approaches.

power consumption for the two approaches based on (20). The linear increase in power per channel is slightly higher with the sequential design than the pipeline. Although there is a sharp jump in power from to , further increases in levels require less and less additional power as the usage rate approaches one. The most important observation from Fig. 9 is that the power consumption of the two implementations is almost similar but the sequential design requires significantly less chip area.

Due to size and power constraints in implantable systems, an important figure of merit is the relative area-power product, which is plotted in Fig. 10 versus both level and channel. Fig. 10 illustrates that the sequential approach is increasingly preferable as the number of channels or the number of decomposition levels increases. The only significant benefits of the pipeline within the enforced design constraints are that it can be clocked at a higher rate and that it takes fewer clock cycles

to complete a computation. Both of these factors result in the pipeline having a *higher threshold* on the maximum number of channels that can be simultaneously processed. However, based on the parameters defined above, the sequential execution architecture has an estimated maximum of around 500 data channels (at $f_{clk} = 100$ MHz). Given the chip area limitations, the area-efficient sequential approach is best suited for this application. In an example implementation with 32 channels and 4 levels of decomposition, the models predict that the sequential approach will require 0.692 mm² and 50.1 nA in 0.13- μ m CMOS, indicating the feasibility of performing front-end signal processing within the constraints of an implanted device.

Another interesting result of this study is the comparison of the area required by the computational node circuitry versus the area required by the memory that holds register values required for multichannel/multilevel operation. Fig. 11 illustrates this result for both sequential and pipeline configurations as a function of channels at $f_{clk} = 100$ MHz. Notice with the pipeline that memory dominates the area when the number of channels is greater than four. For the sequential design, memory dominates when the number of channels is greater than ten. With 10-bit data resolution, at $f_{clk} = 100$ MHz and $V_{DD} = 1.0$ V, the pipeline requires over 14 000 bits

μW μ

Fig. 11

Fig. 11 - 32

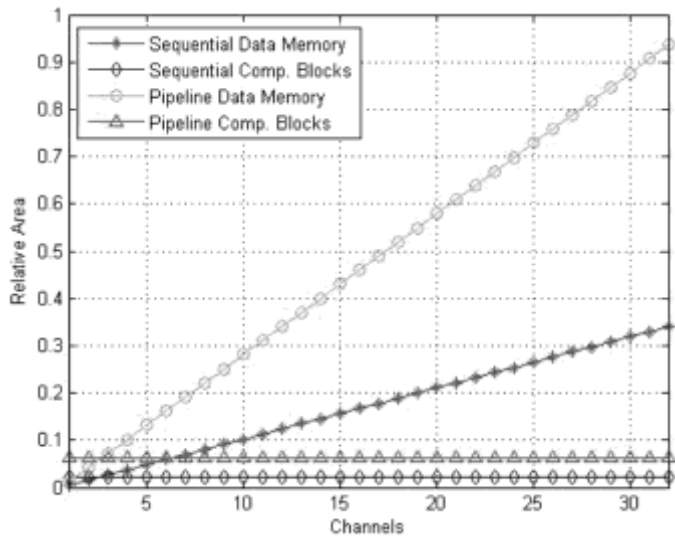


Fig. 11. Relative area versus channels of data memory compared to all other blocks for sequential and pipeline designs, at $L = 4$.

of SRAM, while the sequential circuit requires only about 5000 bits. Reducing memory requirements becomes increasingly important in multichannel applications, again highlighting the advantage of the sequential approach.

C. Lifting versus B-Spline

As illustrated in Fig. 11, the memory required to store intermediate calculation values will dominate circuit area in multichannel implementations. Careful analysis of an optimized sequential B-spline implementation [25] has shown that eight memory registers are required per channel/level, compared to four for sequential lifting and 11 for pipeline lifting. Based on this information and the comparisons above, B-spline has a slight advantage over pipeline lifting but incurs a significant penalty relative to sequential lifting in terms of area. Furthermore, the sequential lifting implementation requires only about 25% of the dynamic power of sequential B-spline, primarily because B-spline takes 18 cycles to execute sequentially compared to 5 cycles for lifting [25]. The advantage of sequential lifting becomes even more profound when static power is considered, especially in deep submicron technologies. Fig. 12 provides an additional comparison, where the number of required gates, synthesized from Verilog descriptions of lifting and B-spline circuits, are plotted. These results illustrate that lifting is increasingly preferable over B-spline as the number of channels and levels increase.

D. Multiplication-Free Lifting

The CC unit proposed in this paper uses one multiplier so that the calculations required per sample are 8 multiplications and 8 additions that can be completed in 5 cycles as listed in Table III. It is noteworthy that a general purpose lifting approach based on only shifts and additions

was proposed in [3]. For the sake of completeness, we compared the demands of a CC unit with a multiplier (proposed in this paper) to a CC unit without a multiplier, i.e., composed of only a shifter and an adder. The later

approach resulted in 12 shift operations and 21 add operations, and required 21 cycles per sample. This is because the equations required to compute multiplication-free lifting DWT did not show any regular structure such as the ones in (7). Therefore, substituting another adder and shifter in the data path did not help in reducing the number of cycles required to complete the computation. With respect to area demands, we found that for one sample pair, a CC unit without a multiplier requires 52% less area compared to a CC with multiplier. This obviously translates into large savings in chip area. However, these savings were not substantial when the system is scaled up. For example, a 32-channel/4-level DWT system using a CC with multiplier would occupy 6.5% of the total chip area as opposed to 3.3% using a CC without multiplier. So the overall savings in chip area are only 3.2%. In contrast, the CC without multiplier requires 13.3% more power than a CC with multiplier for this specification. We therefore concluded that the reduction in area using a shift and add strategy in the lifting approach is overshadowed by the increase in power dissipation when multichannel/multilevel decomposition is sought.

CONCLUSION

VLSI architectures to compute a 1-D DWT for real-time multichannel streaming data under stringent area and power constraints have been developed. The implementations are based on the lifting-scheme for wavelet computation and integer fixed-point precision arithmetic, which minimize computational load and memory requirements. A computational node has been custom designed for the quantized integer lifting DWT and characterized to estimate the maximum achievable computation frequency. Negligible degradation in the signal fidelity as a result of these approximations has been demonstrated.

Detailed comparison between the lifting and the B-spline schemes was presented. It was shown that the lifting approach is more suited when floating point operations are eliminated, thereby superseding the gain achieved by the B-spline approach where adders replace multipliers. Two power and size efficient hardware alternatives for computing the single-level, single-channel wavelet transform have been described and analyzed. The memory management efficiency of the pipeline design results in slightly less power dissipation, while the sequential execution design requires significantly less chip area. Design considerations for scaling these architectures to multichannel and multilevel processing have been discussed. Area and power consumption models with detailed transistor count and switching frequency parameters have been described and used to compare the performance of the

two design alternatives in multichannel and multilevel implementations. The results show many interesting characteristics of each design when it scales to an arbitrary number of levels and channels. When the number of channels is two or more, the sequential execution architecture was shown to be more efficient than the pipeline approach in terms of both power and chip area. Furthermore, results indicate that, using this architecture, multilevel processing of many channels simultaneously is

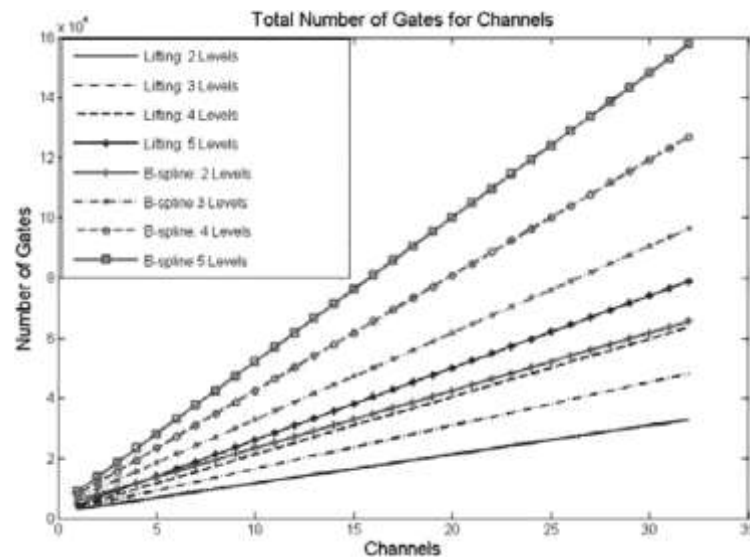


Fig. 12. Total number of gates as a function of the number of channels and the number of levels for the lifting and B-spline implementation.

feasible within the constraints of a high-density intracortical im-plant. This work demonstrates that on-chip real-time wavelet computation is feasible prior to data transmission, permitting large savings in bandwidth requirements and communication costs. This can substantially improve the overall performance of next generation implantable neuroprosthetic devices and brain-machine interfaces.

ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their helpful suggestions and constructive comments.

REFERENCES

K. K. Parhi and T. Nishitani, "VLSI architectures for discrete wavelet transforms," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 1, no. 6, pp. 191–202, Jun. 1993.

C. Chakrabarti, M. Vishwanath, and R. M. Owens, "Architectures for wavelet transforms: A survey," *J. VLSI Signal Process.*, vol. 14, pp. 171–192, 1996.

H. Olkkonen, J. T. Olkkonen, and P. Pesola, "Efficient lifting wavelet transform for microprocessor and VLSI applications," *IEEE Signal Process. Lett.*, vol. 12, pp. 120–122, 2005.

P. K. Campbell, K. E. Jones, R. J. Huber, K. W. Horch, and R. A. Normann, "A silicon-based, three-dimensional neural interface: Manufacturing processes for an intracortical electrode array," *IEEE Trans. Biomed. Eng.*, vol. 38, no. 8, pp. 758–768, Aug. 1991.

K. D. Wise, D. J. Anderson, J. F. Hetke, D. R. Kipke, and K. Na-jafi, "Wireless implantable microsystems: High-density electronic interfaces to the nervous system," *Proc. IEEE*, vol. 92, no. 1, pp. 76–97, Jan. 2004.

D. M. Taylor, S. I. Tillery, and A. B. Schwartz, "Direct control of 3-D neuroprosthetic devices," *Science*, vol. 296, pp. 1829–1832, 2002.

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. 361–365, 2000.

M. D. Serruya, N. G. Hatsopoulos, L. Paninski, M. R. Fellows, and J. P. Donoghue, "Instant neural control of a movement signal," *Nature*, vol. 416, pp. 141–142, 2002.

H. Liao, M. K. Mandal, and B. F. Cockburn, "Efficient architectures for 1-D and 2-D lifting-based wavelet transforms," *IEEE Trans. Signal Process.*, vol. 52, no. 5, pp. 1315–1326, May 2004.

A. M. Kambh, A. Mason, and K. G. Oweiss, "Comparison of lifting and B-spline DWT implementations for implantable neuroprosthetics," *J. VLSI Signal Process. Syst.*, to be published.

K. G. Oweiss, "A systems approach for data compression and latency reduction in cortically controlled brain machine interfaces," *IEEE Trans. Biomed. Eng.*, vol. 53, no. 7, pp. 1364–1377, Jul. 2006.

K. G. Oweiss, "Multiresolution analysis of multichannel neural recordings in the context of signal detection, estimation, classification and noise suppression," Ph.D. dissertation, Univ. Michigan, Ann Arbor, 2002.

K. G. Oweiss, D. J. Anderson, and M. M. Papaefthymiou, "Optimizing signal coding in neural interface system-on-a-chip modules," in *Proc. 25th IEEE Int. Conf. Eng. Med. Biol.*, Sep. 2003, pp. 2016–2019.

I. Daubechies and W. Sweldens, "Factoring wavelet transforms into lifting steps," *J. Fourier Anal. Appl.*, vol. 4, no. 3, pp. 245–267, 1998.

K. A. Kotteri, S. Barua, A. E. Bell, and J. E. Carletta, "A comparison of hardware implementations of the biorthogonal 9/7 DWT: Convolution versus lifting," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 52, no. 5, pp. 256–260, May 2005.

C.-T. Huang, P.-C. Tseng, and L.-G. Chen, "Analysis and VLSI architecture for 1-D and 2-D discrete wavelet transform," *IEEE Trans. Signal Process.*, vol. 53, no. 4, pp. 1575–1586, Apr. 2005.

C.-T. Huang, P.-C. Tseng, and L.-G. Chen, "Flipping structure: An efficient VLSI architecture for lifting based discrete wavelet transform," *IEEE Trans. Signal Process.*, vol. 52, no. 4, pp. 1080–1089, Apr. 2004.

M. Unser and T. Blu, "Wavelet theory demystified," *IEEE Trans. Signal Process.*, vol. 51, no. 2, pp. 470–483, Feb. 2003.

C.-T. Huang, P.-C. Tseng, and L.-G. Chen, "VLSI architecture for forward discrete wavelet transform based on B-spline factorization," *J. VLSI Signal Process.*, vol. 40, pp. 343–353, 2005.

S. Mallat, *A Wavelet Tour of Signal Processing*, 2nd ed. New York: Academic, 1999.

D. Donoho, "Denosing by soft thresholding," *IEEE Trans. Inf. Theory*, vol. 41, no. 5, pp. 613–627, May 1995.

K. Andra, C. Chakrabarti, and T. Acharya, "A VLSI architecture for lifting-based forward and inverse wavelet transform," *IEEE Trans. Signal Process.*, vol. 50, no. 4, pp. 966–977, Apr. 2002.

Y. Suhail and K. G. Oweiss, "A reduced complexity integer lifting wavelet based module for real-time processing in implantable neural interface devices," in *Proc. 26th IEEE Int. Conf. Eng. Med. Biol.*, Sep. 2004, pp. 4552–4555.

R. Calderbank, I. Daubechies, W. Sweldens, and B.-L. Yeo, "Wavelet transforms that map integers to integers," *Appl. Comput. Harmon. Anal.*, vol. 5, no. 3, pp. 332–369, 1998.

B. F. Wu and C. F. Lin, "A rescheduling and fast pipeline VLSI architecture for lifting-based discrete wavelet transforms," in *Proc. IEEE Int. Symp. Circuits Syst.*, May 2003, vol. 2, pp. 732–735.

P. Y. Chen, "VLSI implementation for one-dimensional multilevel lifting-based wavelet transform," *IEEE Trans. Comput.*, vol. 53, no. 4, pp. 386–398, Apr. 2004.

F. Rieke, D. Warland, R. R. van Steveninck, and W. Bialek, *Spikes: Exploring the neural code*. Cambridge, MA: MIT press, 1997.

K. Oweiss, "Compressed sensing of large-scale ensemble neural activity with resource-constrained cortical implants," *Soc. Neurosci. Abstr.*, vol. 13.11, Oct. 2006.

[View publication stats](#)