

Inexact-aware architecture design for ultra-low power bio-signal analysis

ISSN 1751-8601

Received on 20th October 2015

Revised on 18th February 2016

Accepted on 22nd March 2016

doi: 10.1049/iet-cdt.2015.0194

www.ietdl.org

Soumya Basu¹ ✉, Pablo Garcia Del Valle¹, Georgios Karakonstantis², Giovanni Ansaloni³, Laura Pozzi³, David Atienza¹

¹Embedded Systems Laboratory, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland

²School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, Belfast, UK

³Faculty of Informatics, Università della Svizzera Italiana, Lugano, Switzerland

✉ E-mail: soumya.basu@epfl.ch

Abstract: This study introduces an inexact, but ultra-low power, computing architecture devoted to the embedded analysis of bio-signals. The platform operates at extremely low voltage supply levels to minimise energy consumption. In this scenario, the reliability of static RAM (SRAM) memories cannot be guaranteed when using conventional 6-transistor implementations. While error correction codes and dedicated SRAM implementations can ensure correct operations in this near-threshold regime, they incur in significant area and energy overheads, and should therefore be employed judiciously. Herein, the authors propose a novel scheme to design inexact computing architectures that selectively protects memory regions based on their significance, i.e. their impact on the end-to-end quality of service, as dictated by the bio-signal application characteristics. The authors illustrate their scheme on an industrial benchmark application performing the power spectrum analysis of electrocardiograms. Experimental evidence showcases that a significance-based memory protection approach leads to a small degradation in the output quality with respect to an exact implementation, while resulting in substantial energy gains, both in the memory and the processing subsystem.

1 Introduction

Busy and unhealthy lifestyles are becoming common, resulting in a rise in the number of people developing or living with cardiovascular conditions. Moreover, a significant part of the world population is ageing, and hence becoming in danger of contracting cardiac diseases. This scenario calls for increased levels of medical supervision and management, which are resulting in high costs, and traditional healthcare infrastructures are finding it increasingly difficult to cope with these demands [1].

Emerging wireless body sensor [2] network technologies can offer large-scale and cost-effective solutions to this problem. These wearable devices for bio-signal monitoring are bringing about a revolutionary change in healthcare systems by allowing long-term monitoring of chronic patients, while providing a low-cost and unobtrusive solution. Wireless body sensor nodes (WBSNs) [3], represented in Fig. 1, are the building blocks of such a network, since they can provide real-time and personalised monitoring of patients. They are designed to monitor different organs of the human body, including the heart. Such devices involve sensing of bio-signals and then transmitting them wirelessly to receiver devices, for further analysis. The analysis of the received sensed data generally consists of labour-intensive manual inspection or offline execution on a server infrastructure.

Recently, however, a new generation of smart WBSNs has emerged which are able to perform digital signal processing directly on-board to analyse the acquired bio-signals and extract clinically-relevant features, in addition to data acquisition and transmission [4]. These devices pave the way for truly autonomous and versatile health monitoring devices.

Energy efficiency is of paramount importance in these battery operated WBSNs, as they work under tight energy constraints, defined by the battery-based power supplies on the device. Performing on-chip signal analysis results in increased computation, thereby increasing energy consumption. For this reason, there is an urgent need for efficient energy management in

WBSNs, which has fuelled significant research interest. An efficient method to achieve large energy savings in the execution of these applications is voltage scaling. However, voltage scaling is limited by the low reliability of static RAM (SRAM) memories when operating in a near-threshold regime, which results in a high probability of random bit-flip errors.

In this work, we explore the ensuing inexact computation (or approximate computation) as a new method to achieve higher energy efficiency in WBSNs. It involves *trading off the accuracy of logic circuits in order to save energy*, by applying techniques like voltage scaling [5] and circuit pruning [6], among others. In particular, in this paper we present an architecture that embodies the paradigm shift from exact to inexact computation, targeting the WBSN scenario. Bio-signal processing applications normally acquire noisy input and produce qualitative outputs, thereby being error-resilient in nature. Also, they frequently involve storing data that is *sparse* in nature [7], with the memory components accounting for a key part of the energy consumed [8].

The noise-resilient properties of bio-signal processing applications motivate the use of an inexact computing paradigm. In this context, we introduce the concept of *data significance*, and we derive a significance-based heterogeneous protection of the memory, working at near-threshold supplies. We explore the effects of the proposed scheme on the energy savings obtained w.r.t protecting the whole memory subsystem, and the effect of errors introduced on the final output.

We applied the proposed methodology to a typical WBSN architecture (depicted in Fig. 2), which comprises a low power processor along with supporting memories. We assumed that the application content is stored in a non-volatile memory (NVM), and transferred into SRAMs during bootstrap. This initialisation phase only requires a tiny fraction of the total run-time of typical WBSN acquisitions, which can last from hours to days. Our work therefore focuses on the architectural components active at run-time: the microprocessor and the SRAM system (highlighted by the shaded rectangle in Fig. 2)

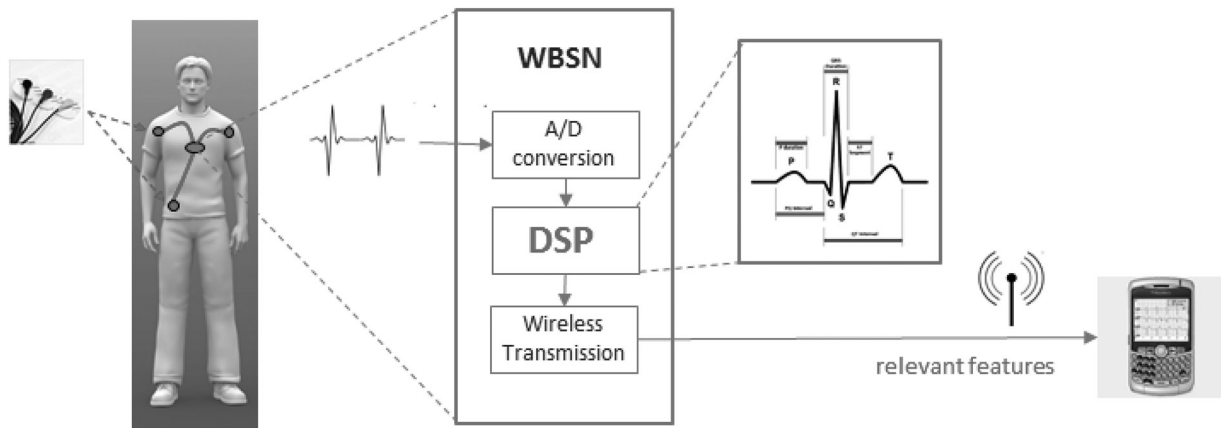


Fig. 1 Schematic representation of a WBSN node

To evaluate our proposed methodology, we selected the power spectral analysis (PSA) application of the heart rate variability (HRV). Many applications have been proposed to predict heart diseases, ranging from the automated detection of epileptic seizures [9], to the predictive risk assessment of atrial fibrillations [10]. The choice of PSA is based on the fact that it is one of the most widely used strategies for predicting cardiac failures, as it allows the monitoring of various health conditions associated with the heart, as well as other organs [11, 12].

The implementation of PSA on ultra-low power embedded devices requires a carefully tailored digital architecture. Lowering supply voltages might result in quadratic energy savings, but memory subsystems using conventional 6-transistor (6T) SRAM cells become unreliable in nanometre technologies [13]. Soft errors, like bit-flips, start occurring at near-threshold voltages, and their probability of occurrence increases as the supply voltage is further lowered [14]. These issues regarding reliability of memories become more prominent with modern technology scaling, as transistors with smaller lengths are conceived. In this context, larger memory cells have been proposed in previous research work, which consist of 8-transistor (8T) or 10-transistor (10T), because they ensure reliable operation of the memories at much lower supply voltages compared with 6-transistor (6T) cells [15].

However, 8T and 10T cells present a high area footprint, limiting the amount of memory that can be included in area-constrained WBSNs. The majority of the silicon real estate of typical WBSN processors is devoted to the memory subsystem, a characteristic that will become even more preeminent in the foreseeable future, as forecasts predict that as much as 95% of the entire chip area will be devoted to memories [16]. Thus, area-hungry 8T or 10T cells

should be employed judiciously. Another proposed approach in the literature, to deal with errors in memories, is to use error correction codes (ECCs) [17]. Nonetheless, even ECCs require significant area and energy when all or large parts of the memories need to be protected. These characteristics highlight the need of an effective memory protection strategy that, combining different memory implementation and error detection/correction mechanisms, allows reliable operation at ultra-low supply levels, while at the same time presenting small overheads.

To sum up, in this article we introduce a novel inexact architecture featuring a heterogeneous memory protection scheme which takes advantage of the *sparsity* of data in the targeted application. We advocate the application of the proposed scheme based on significance of data, rather than based just on significant bits, in order to achieve high energy savings, while binding the error introduced in the output of the application within permissible limits.

The main contributions of this paper are the following:

- (i) We analyse the PSA application's software code in order to explore its statistical properties. This includes analysis of the data elements in the intermediate steps of the application and the classification of the data into more significant and less significant, depending on their contribution to the output quality of the system.
- (ii) We introduce a novel hybrid memory protection scheme, which involves a significance-based protection in hardware of the data memory for the PSA system, using ECC bits. Extending our previous work [18], herein we explore the system-wide impact of ultra-low voltage scaling, reporting the effects on the entire data memory (including sections outside the one dedicated to storing intermediate buffers), on the instruction memory and on the processor system.
- (iii) We estimate the effects of using voltage scaling with the proposed memory protection scheme on the performance of the PSA application and the energy savings that the scheme results in.

The rest of the paper proceeds as follows. In Section 2, we introduce the proposed scheme of memory protection driven by data significance. Section 3 presents a case study where the PSA application is analysed. In Section 4 we explain our experimental setup, followed by the results obtained, before finally concluding the paper.

2 Data-significance driven criticality approach in WBSNs

The impact of adopting ultra-low voltage supplies is not homogeneous across architectural blocks. Combinational logic circuits such as arithmetic logic units (ALUs) are the least affected, because they are stateless and do not present internal feedback connections. Indeed, combinational circuits have been

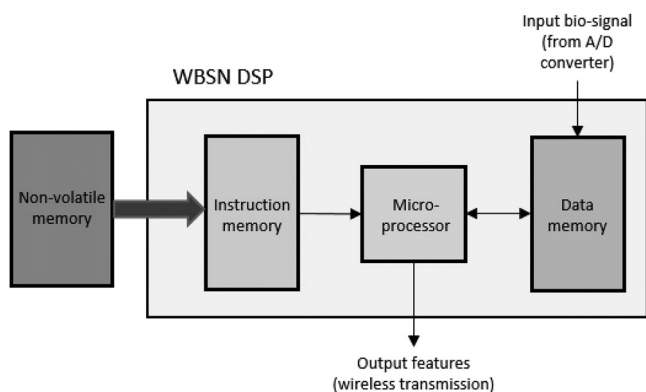


Fig. 2 Block scheme of a typical WBSN's processing unit. Only the microprocessor and the data/instruction memories are active after the bootstrap phase

proposed operating at supply voltages (V_{dd}) in the range of few hundreds of millivolts [19]. Closely coupled with ALUs is the register file, which embodies the first level of the memory hierarchy. In a load/store architecture such as ARM, target of the present work, all instructions except explicit loads and stores operate with data residing in registers. The implementation of the register file (the Cortex M3 has 16 architecturally-visible registers) is usually based on standard cell memories, which can reliably operate at extremely low voltage levels [20]. From these two observations, we conclude that the computing core, integrating the ALU and the register file, is not the resiliency bottleneck of the system.

Conversely, the data and instruction memories (DM and IM, respectively) of WBSNs are commonly implemented as 6T SRAMs, which are more prone to random bit-flips when operated in near-threshold regime. In [14], a non-negligible probability of error of 1.3×10^{-5} is reported for a V_{dd} of 0.75 V, which rapidly escalates as the supply level drops. A platform adopting a low voltage level for the processor and its associated registers, and a higher one for the instruction and data memories, could concurrently allow high efficiency and reliability. However, such an approach requires multiple voltage regulators, as well as voltage converters on each signal crossing the boundaries of the voltage islands. This solution is thereby not efficient for ultra-low power platforms [8]. More complex cell structures employing dedicated read and write paths (8T- or 10T-SRAMs) can reliably operate at a lower voltage supply, but at the cost of important area and energy overheads. In [8] and from the CACTI memory modelling tool [21], it is found that 8T cells occupy 30% more real-estate, and consume on an average 25% more dynamic energy, than comparable 6T implementations at 40 nm technology. An alternative path to ensure reliability is to detect and correct errors using redundant representations of the memory content, adding dedicated ECCs to transparently recover from bit-flips. Even in this case, area and power overheads have to be accounted for, due to the dedicated memory cells required to store the redundant information, and the logic required to recover from errors.

Herein, we propose to minimise the above-mentioned overheads by imposing different reliability guarantees, depending on the criticality of the data. The IM content is highly susceptible to errors, as a single bit-flip can lead to an unpredictable execution flow, or even to unrecoverable states (e.g.: if a jump to a random location is made). The impact of bit-flips in DM can instead be less pronounced. While control variables and address manipulations have indeed high criticality levels, a large portion of the DM in bio-signal analysis applications is employed to store windows of data containing inputs and outputs of the various stages of digital processing. Errors in these buffers do not lead to catastrophic failures, but can cause an unacceptable degradation of the output quality.

Crucially, the loss in end-to-end quality of service, or in other words the net performance of the system, deriving from random errors, is dependent on the statistical properties of the stored data. This characteristic is leveraged here to guide the design of heterogeneous protection schemes, which provide correction,

detection or only a best-effort guarantee on different memory sections, maximising the quality of service for a target energy budget. In our approach, we distinguish between two important cases, addressing *sparse* and *non-sparse* buffers. In the latter case, the magnitude of each entry in a buffer array is randomly distributed. For these arrays, each entry equally contributes to the overall correctness of the computation, so each stored word must expose the same reliability level. Protection of non-sparse buffer must be therefore performed at the bit-level, ensuring the correctness of high-order bits, while possibly allowing a degree of inexactness for low-order ones.

This strategy, however, becomes sub-optimal when the memory content is mostly centred on an expected value, with only few words significantly deviating from it, which is often the case in WBSN applications [22]. This *sparsity* property allows the adoption of a word-based, instead of bit-based, protection scheme, in which error correction is employed for the small subset of data, which is not close to the expected value (and we term this subset *significant*), while a much simpler error detection mechanism is used for the rest. In this way, correctness is ensured for significant words, while bit-flips in the non-significant parts are only partially countered, by adopting the expected value of the data upon the detection of an error.

We apply the above-mentioned considerations to the design of the target inexact architecture, whose block scheme is provided in Fig. 3. To increase reliability at low supply levels, the instruction memory is realised with 8T-SRAM cells. As for the data memory, energy- and area-efficient 6T-SRAMs are employed, coupled with heterogeneous error detection/protection features. This arrangement results in a simple implementation of the IM (which is protected in its entirety), while allowing a fine-grained tuning of DM reliability, dependent on the criticality of the stored values. Scalar and control variables (non-buffer data), as well as the highly significant portions of the buffer data, are fully protected with multi-bit error detection codes. Conversely, only error detection (implemented as a 1 bit parity code), but not correction, is employed for non-significant buffer data in sparse arrays. Finally, for the non-sparse arrays, errors in the most significant bits are corrected, while no detection or correction is performed for the least significant parts of each word.

3 Data-significance analysis of the power spectral analysis system

3.1 Functionality of the PSA system

In this paper, we use as a case study the PSA of the HRV, which is a powerful tool for evaluating the autonomic control of the heart rate and for identifying various health conditions [23]. Fig. 4 shows the block scheme of the application.

The PSA system is composed of four essential steps:

- (i) In the first step, the time differences between consecutive heartbeats (known as RR intervals) are extracted from

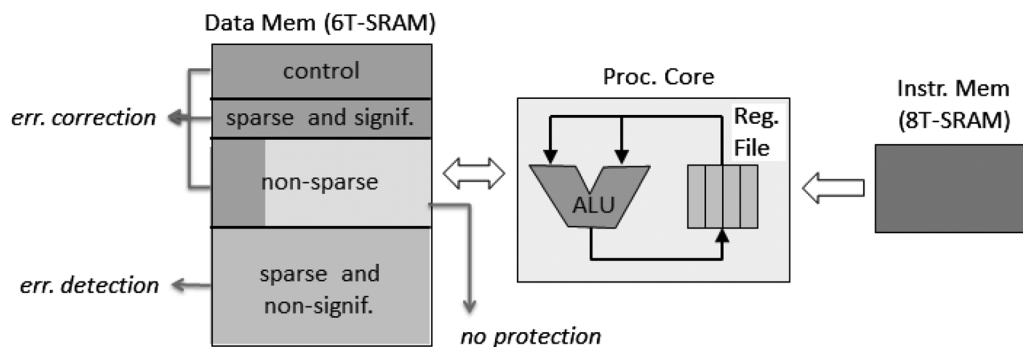


Fig. 3 The proposed heterogeneous protection scheme applies varying exactness guarantees depending on the data criticality

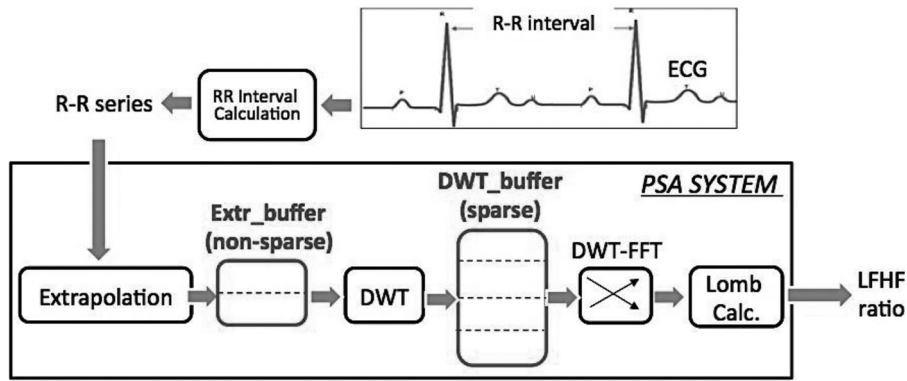


Fig. 4 Block scheme of the PSA application

electrocardiogram (ECG) recordings of patients. The RR intervals are non-periodic signals, and thus require processing by dedicated algorithms, such as the fast-Lomb periodogram method [24].

(ii) In the next step, according to the fast-Lomb method, the extracted RR intervals are extrapolated to a fixed size window (i.e. 512 samples). This procedure essentially converts the non-periodic signals into uniformly sampled ones.

(iii) Then, the uniformly-sampled data are processed to estimate the specific trigonometric functions required by fast-Lomb. Traditionally, such an estimation is performed by applying fast-Fourier transform (FFT). Instead, and similarly to [24], in this paper we use a wavelet-based FFT (WFFT), which reduces substantially (up to 28% w.r.t the state-of-the-art) the complexity of the fast-Lomb method [25], while also introducing *sparsity* in the data being processed [26]. In particular, the wavelet transform involved in the WFFT helps in revealing the sparse nature of bio-signals in the wavelet domain, exposing eventually the terms that are zero (or close to zero). Such close-to-zero terms and the following butterfly operations applied in the second stage of the WFFT can then be pruned, eventually reducing the computational complexity.

(iv) Finally, the Lomb calculator combines the output data obtained from WFFT, estimating the real-time power spectrum information. In clinical practice [24], the most used metric derived from PSA is the ratio between the power in low frequencies (LFP, defined as 0.04–0.15 Hz) and high frequencies (HFP, 0.15–0.4 Hz), with $LFHF\ ratio = LFP/HFP$. A deviation of the $LFHF\ ratio$ above or below normal values is indicative of various health issues [27].

3.2 Analysis of the PSA system

The application of our scheme requires the identification of the statistical characteristics of the target application, for identifying blocks of data where it can be applied. To this end, we have analysed the PSA system and estimated the data distribution in the various stages, by performing several experiments with the ECG recordings. Fig. 5 focuses on the distribution of the data in the two memory buffers used in the system: the `Extr_buffer` is used to store the output of the extrapolation on the input data and the `DWT_buffer` is used to store the DWT outputs, as indicated in Fig. 4.

We can observe in these two figures that the elements of the `DWT_buffer` (Fig. 5b) are mostly centred on zero, justified by their sparse nature, while the elements of the `Extr_buffer` (Fig. 5a) have a non-sparse distribution.

The different data distribution patterns indicate that different protection approaches against memory faults can be applied for limiting the overhead. Intuitively, for taking advantage of the error resiliency of such an application we apply a scheme in which the most significant bits (MSBs) of every word in the `Extr_buffer` are protected with a state-of-the-art mechanism such as ECC, whereas the least significant bits (LSBs) are not protected against memory faults by any specific mechanism, as depicted in Fig. 6.

On the other hand, in case of the `DWT_buffer`, the distribution of the stored data allows application of a more elaborate protection scheme. In particular, rather than protecting groups of bits, here we can protect complete words, distinguishing between significant and less significant ones. In fact, in case of the less significant

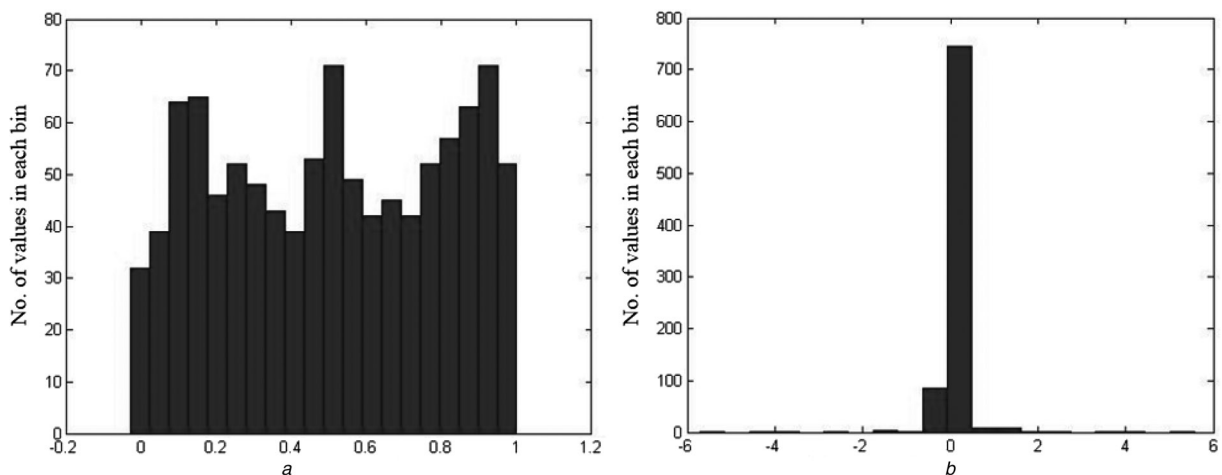


Fig. 5 Histogram of data values in `Extr_buffer` and `DWT_buffer` (normalised), distributed in 20 bins

- a `Extr_buffer` presents a non-sparse distribution
- b `DWT_buffer` presents a sparse distribution

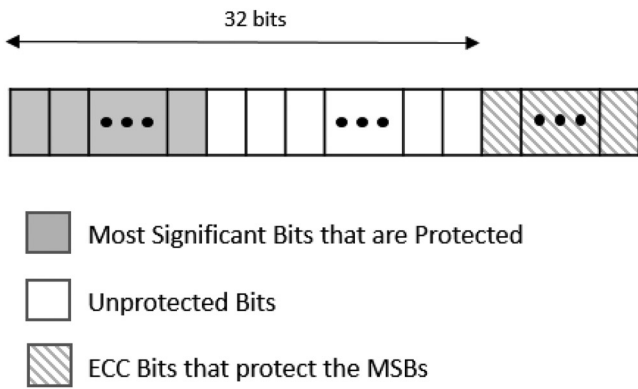


Fig. 6 Memory word in the *Extr_buffer*

data, as most of the values are close to zero, it is possible to replace them with their expected value (zero) if an error occurs in a word. This ensures that the impact of such an error will not drastically affect the expected data, since a flipped bit within each of the close-to-zero data can alter completely the magnitude of the stored value. Error detection is supported by a single parity bit per word, resulting in a small overhead with respect to error correction.

For the significant data, which has magnitudes much larger than zero, a more expensive error correction scheme must be applied for ensuring their correct storage. In the PSA application, such elements reside in the low-frequency outputs of the DWT. In this case, we considered single error correction, double error detection (SECDED) ECC (Fig. 7). Six ECC bits are required to protect a 32 bit word, as used in the data memories.

Note that the partition between significant and less significant words can be selected statically, i.e. independently from the

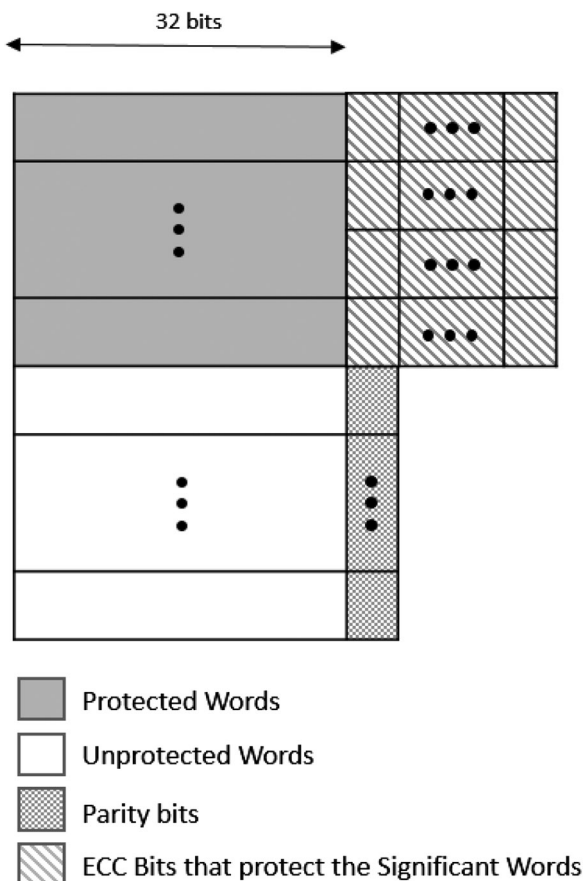


Fig. 7 Block diagram of the *DWT_buffer*

particular window of inputs being processed. Identification of sparse and non-sparse data can be performed with an off-line analysis of the application, and thereby does not imply any run-time overhead. In case of the fast-Lomb, the data are distinguished based on the inherent properties of the DWT, resulting in the separation of the processed data into high and low frequencies. High frequencies (close-to-zero data) are termed as non-significant, whereas low frequencies are termed as significant.

4 Experimental setup

To evaluate the application of the heterogeneous scheme in the PSA application, we retrieved the input ECG signals from real-world recordings, available in the Physionet PAF prediction database [28]. This database includes 300 recordings, each of 30 min. Each recording consists of data acquired from two simultaneously operating ECG sensors, and the root mean squared value of the data from the two sensors. We have considered time windows of 6 min, with an overlap of 5 min, for processing the input data. As a result we obtained 25 time windows for each of the 30 min long recordings. The data from each window in each recording were independently processed by the application to retrieve their LFHF ratio.

To estimate the power consumption of the proposed system, we developed high-level models of its different components, including the processor, the data and the instruction memory, whose implementation is detailed as follows.

We considered a technology node of 40 nm as a realistic example for next-generation WBSN nodes (which are currently 65 and 90 nm), and an ambient temperature of 300 K (27°C). To obtain the static energy from static power figures, we assumed a working frequency of 168 MHz and measured the run-time of the PSA application on all input windows on the target Cortex-M3 processor, which corresponds to 2.23 s. The dynamic energy consumption of the processor was calculated from the values reported in its datasheet [29]. Its leakage energy was derived by considering the leakage power values of the processor while in standby mode, with the clock inactive. Both static and leakage energy figures were scaled according to the target supply voltage.

The IM, which always needs to operate reliably, is implemented by using 8T SRAM cells. The static power and dynamic energy figures of the IM were obtained from [21], adapting them to the target 40 nm technology. To calculate the IM dynamic read energy, we considered a worst-case scenario in which an instruction is fetched every clock cycle. The IM has a size that is able to store the whole application. While we employed full shadowing in our model, our framework can also be employed in conjunction with partial allocation policies [30]. This combined approach would result in further energy and area gains by minimising the amount of required 8T SRAM, at the cost of frequent page transfers from the NVM [31], the latter possibly employing recently proposed NVM structures such as in [30].

The data memory, realised with 6T SRAMs, is itself divided into two sections. The first one comprises the part outside the *Extr_buffer* and *DWT_buffer*, which must also operate without errors irrespective of the supply voltage. It is therefore entirely protected by SECDED codes. The non-buffer data memory (*DM_Rest*) was modelled using CACTI [21] to retrieve the dynamic energy (read and write) per access, while the total number of accesses was estimated using software counters. The leakage power reported by CACTI was adopted to compute the leakage energy, considering the application run time.

The second data-memory section is composed of the data buffers (*Extr_buffer* and *DWT_buffer*), abbreviated as *DM_Buff*, and target of our approach of data-driven inexact scheme. For our experiments, we have considered a maximum of one error occurring in a memory word. In the case of the sparse *DWT_buffer*, we have employed six ECC bits for the protection of the most significant words, and one parity bit for detecting errors in the less significant words. To simulate this heterogeneous memory structure, two separate CACTI models were used as a

starting point, either employing 6 bit ECC or 1 bit parity for the whole memory content. To derive the dynamic energy per read access of intermediate configurations, corresponding to the partial protection schemes, we employed the formula as described in the following equation

$$E_t = p \times E_p + (1 - p) \times E_u \quad (1)$$

where E_p is the read energy per access in the protected memory, and E_u is the read energy per access in the unprotected memory. Also, p is the percentage of considered significant words. E_t is the net read energy per access in the heterogeneous memory. The write energy per access and the leakage power of the hybrid memory were also calculated in the same manner. For the memories having ECC protection, CACTI reports the total static and dynamic energies, including the overhead due to the additional check bits and the data bus. However, it neglects the energy consumed by the encoding and decoding logic of the ECC. Therefore, the CACTI model underestimates the benefits of partial protection with respect to an alternative approach in which the data memory is fully protected. In fact, since only some parts of the memory require SECDED, smaller dedicated circuits can be activated, resulting in reduced dynamic energy consumption. In addition, some memory locations only require a simpler parity calculation, or no protection at all.

To evaluate the impact of errors in unreliable memories, binary error masks were randomly generated for each buffer. We considered single bit-flip errors with probabilities of 0.07 and 0.22%, relative to the behaviour of a 6T SRAM cell working with supply voltages of 0.65 and 0.6 V, respectively [14]. A '1' in a mask induces a bit-flip in the corresponding position of the target buffer, but only if that value belongs to the part of the memory that is unprotected. While we considered a maximum of one error per memory word in the experimental evaluation of Section 5, more complex error detection and correction methods can be adopted, allowing recovery from multiple bit-flips. In those cases, the ECC overhead would increase significantly, highlighting even more the benefit of only protecting significant words or bits.

The impact of these errors on the quality of the output of the PSA application, under the different protection schemes, was then measured by comparing the obtained LFHF ratio with respect to an error-free execution.

We compared our inexact architecture against two different baselines:

(i) *High V_{dd}* : In the first case we considered a high supply voltage (1.1 V), which does not impact the reliability of the system. All memories in this case were implemented as 6T SRAMs, whose energy values were computed by modelling them in CACTI.

(ii) *Low V_{dd} and total ECC protection*: In the second case, we have considered exact operations at low supply voltage levels (0.65 V). This requires the implementation of the instruction memory with 8T SRAMs, while all of the data memory (buffer and non-buffer) is completely protected by SECDED ECC codes.

5 Experiment results

We evaluated our system in three parts. First, we analysed the performance degradation of the PSA system in calculating the LFHF ratio, under the different configurations of the heterogeneous memory scheme. Next, we studied the energy savings achieved by using the proposed configurations. Finally, we reported the energy-performance trade-offs for the different protection schemes.

5.1 Analysis of the introduced error

The results of the error simulations have been achieved by averaging individual results obtained by processing of data in each time window for each ECG recording. Figs. 8a and b show the percentage of error in the computation of the LFHF ratio by the PSA application, when compared with an error-free version of the same, under the different test-points of the proposed heterogeneous memory scheme at supply voltages of 0.65 and 0.6 V, respectively.

Our obtained results show that selective protection of a small fraction of words (the significant ones) in the sparse buffers can still guarantee high-quality performance of the system, with respect to an error-free version. This shows the error-tolerance capabilities of WBSN applications. As an example, 1.3% relative error is incurred in the LFHF ratio by protecting 11 MSBs in the *Extr_buffer* and 15% significant words in the *DWT_buffer* (Fig. 8a). In the clinical practice, this error magnitude can be well tolerated. Indeed, usually the medical literature reports results on LFHF analysis with only 2–3 significant digits [27]. Moreover, results must be aggregated and weighted with respect to multiple factors including, but not limited to, the age, race and gender of the patient [32]. However, the percentage of error in the LFHF ratio obtained by protecting only 4 MSBs in the *Extr_buffer* exceeded 20%, even in the case of full protection of the

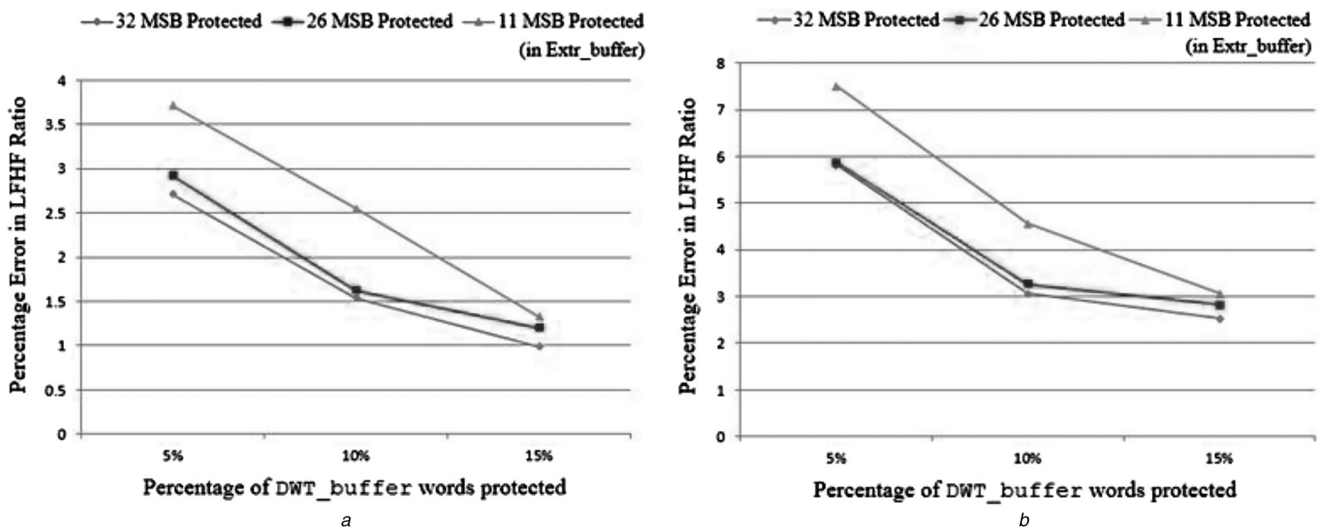
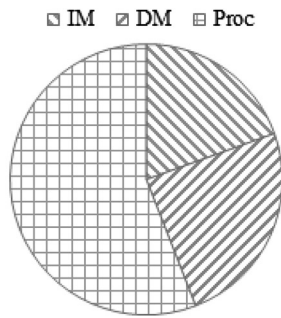


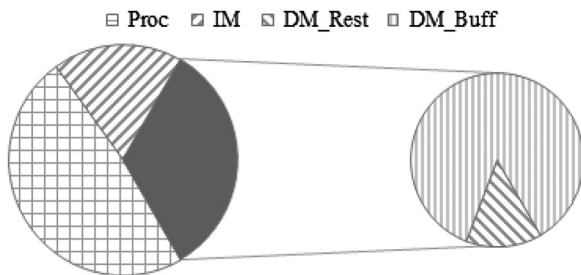
Fig. 8 Percentage of error in the calculation of the LFHF ratio under the different protection schemes

a At 0.65 V supply
b At 0.6 V supply



Energy consumed by different parts of the system at 1.1V
Total energy consumed = 0.58J

Fig. 9 Energy consumption at baseline 1 (high V_{dd})



Energy consumed by different parts of the system at 0.65V
Total energy consumed = 0.26J

Fig. 10 Energy consumption at baseline 2 (low V_{dd})

DWT_buffer. This high error-rate is not acceptable in bio-signal processing and thus we have excluded the condition of protecting only 4 MSBs in the Extr_buffer from the following sections of this paper.

In the case of 32 MSBs protected in the Extr_buffer and 15% of significant words protected in the DWT_buffer, the relative error is $<1\%$. This figure is bound below 4% even when we consider the worst-case protection from our experimental setup (11 MSBs protected in the Extr_buffer and 5% of significant words protected in the DWT_buffer).

The percentage error in the LFHF ratio for a supply voltage of 0.6 V is shown in Fig. 8b. The same trends as in Fig. 8a are noticed also in this case, but with higher relative error with respect to operation at 0.65 V supply voltage. This is due to the much higher number of bit-flip errors in the memories at 0.6 V supply, when compared with 0.65 V. Interestingly, even in this case, the error in the LFHF ratio, with respect to a fault-free execution, can be limited to 5% by allowing errors in the 21 LSBs of Extr_buffer and only checking (but not correcting) errors in 90% of DWT_buffer.

5.2 Analysis of energy consumption

Figs. 9 and 10 compare the energy consumption of the different system components, at the first and second baselines considered (high V_{dd} and low V_{dd} with complete protection, respectively). It can be seen that at low supply voltage, the system already shows substantial energy savings when compared with operation at high supply voltage. Moreover, at low V_{dd} , it can be observed that data memory accounts for a major part of the energy budget and that the targeted buffers account for most of the energy consumed by the data memory (Fig. 10). This justifies the application of the proposed memory protection scheme to save even more energy in these buffers, thereby further enabling energy benefits at low-voltage operating points.

Fig. 11 shows the total energy consumption of the targeted memory buffers under the different protection schemes. We can observe from it that by using our proposed scheme with the condition where we protect 11 MSBs in the Extr_buffer and 10% significant words in the DWT_buffer (corresponding to a relatively low error of 2.6% in the LFHF ratio as in Fig. 8a), we were able to save about 18% of the energy in the buffers compared with the second baseline, which results in an overall gain of 5.2% for the whole system.

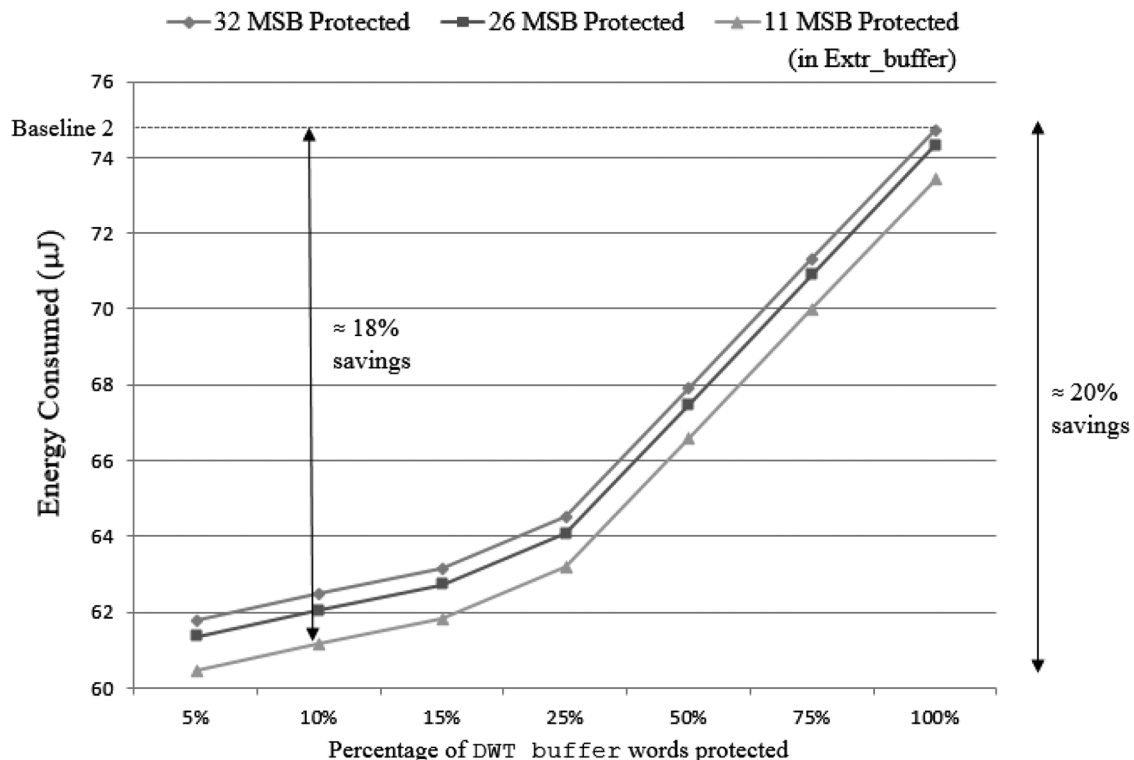


Fig. 11 Total energy consumption by the targeted memory buffers at 0.65 V supply under different memory protection schemes for an execution time of 2.23 s

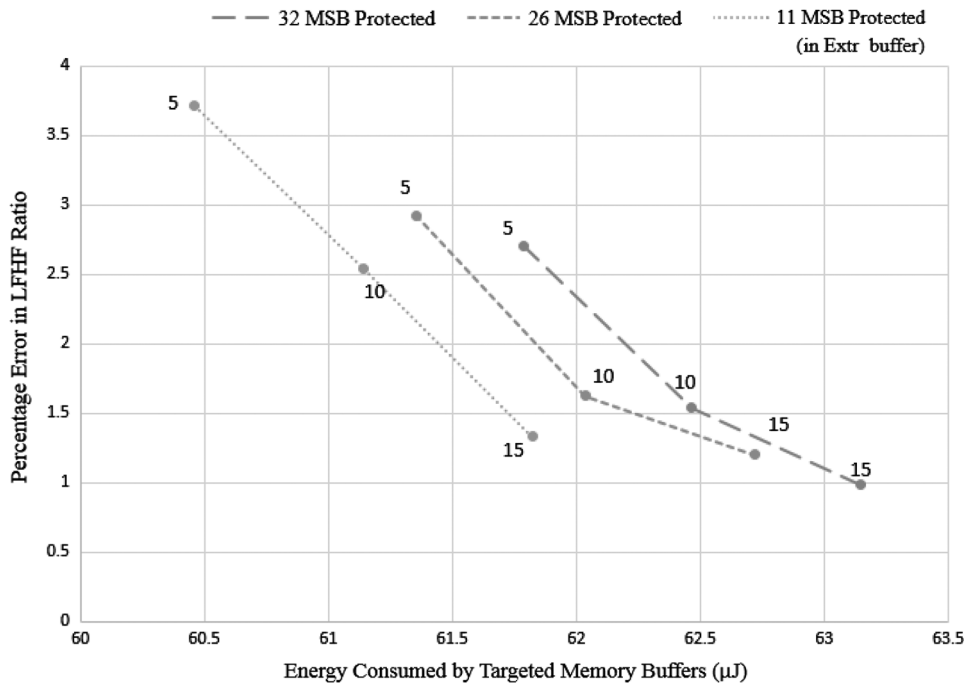


Fig. 12 Percentage error in LFHF ratio and the corresponding energy consumption under different memory protection schemes at 0.65 V. Numbers beside the points represent the % of DWT_buffer protected

It can be further observed from Fig. 11 that by using the proposed heterogeneous memory protection scheme we could achieve almost 20% of savings in energy in the targeted buffers in the most energy-efficient case of protection considered (11 MSBs in Extr_buffer and 5% of significant words in DWT_buffer), over a scheme which involves protecting the buffers completely with SECDED codes.

5.3 Energy/performance trade-off analysis

Combining the results obtained in the previous sections, we can investigate the trade-offs between relative error in LFHF ratio and the energy consumed. This enables selection of the optimal memory architecture considering both energy consumption and performance degradation. The total energy consumed by the targeted buffers under the different protection schemes are plotted against the percentage error in the computation of the LFHF ratio at 0.65 V, as shown in Fig. 12.

It shows that protection schemes with less energy consumption result in higher performance degradation. As an example for selecting the optimal protection scheme, for a maximum tolerable error of 3%, a solution with 11 MSBs protected in the Extr_buffer and just 10% of the most significant words in the DWT_buffer is the best one in terms of energy efficiency. On the other hand, for an energy budget of 61.5 µJ, the smallest percentage error can be achieved by protecting just 11 MSBs in the Extr_buffer and 10% of the DWT_buffer.

6 Conclusions

In this study we have introduced a novel heterogeneous memory architecture to increase the power efficiency of WBSNs by selectively protecting data with high criticality. Our experiments show that, by guaranteeing different amount of reliability in the bits and words of varying *significance*, the energy required by the considered PSA bio-signal processing application can be reduced beyond the levels attainable by voltage/frequency scaling alone, with a minimal degradation in the quality of service.

The results of our experiments have shown that by applying the resulting heterogeneous protection scheme, we were able to reduce

~20% of the energy budget of the data memory used in the intermediate data buffers in prospective real-life wearable ECG monitoring systems. Moreover, our approach is able to tolerate the high error rates incurred at ultra-low voltage supply levels. This approach also supports scaling to ultra-low operating voltages, which provides substantial energy benefits compared with high voltage operation.

Finally, our complete approach and system design framework is applicable to various health-monitoring applications beyond PSA, as they share similar sparse data distribution characteristics. These common features refer to dealing with noisy signal inputs, providing a statistical or qualitative output and consisting of intermediate buffers, which are very common requirements nowadays.

7 Acknowledgments

This work has been partially supported by the ONR-G grant no. N62909-14-1-N072 and the E4Bio RTD project (no. 200021_159853) evaluated by the Swiss NSF.

8 References

- MEP Heart Group: 'Cardiovascular diseases facts and figures'. Available at <http://www.mepheartgroup.eu/index.php/facts-a-figures>
- Hao, Y., Foster, R.: 'Wireless body sensor networks for health-monitoring applications', *Physiol. Meas.*, 2008, **29**, (11), pp. R27–R56
- Braojos, R., Dogan, A., Beretta, I., *et al.*: 'Hardware/software approach for code synchronization in low-power multi-core sensor nodes', Design, Automation and Test in Europe Conf. and Exhibition (DATE), 2014
- Braojos, R., Giovanni, A., Atienza, D.: 'A methodology for embedded classification of heartbeats using random projections'. Design, Automation and Test in Europe Conf. and Exhibition (DATE), IEEE, 2013, 2013, pp. 899–904
- Ganapathy, S., Karakonstantis, G., Teman, A., *et al.*: 'Mitigating the impact of faults in unreliable memories for error-resilient applications'. Proc. Design Automation Conf., 2015
- Du, Z., Lingamneni, A., Chen, Y., *et al.*: 'Leveraging the error resilience of machine-learning applications for designing highly energy efficient accelerators'. 19th Asia and South Pacific Design Automation Conf. (ASP-DAC), 2014, pp. 201–206
- Mamaghanian, H., Khaled, N., Atienza, D., *et al.*: 'Compressed sensing for real-time energy-efficient ECG compression on wireless body sensor nodes'. *IEEE Trans. Biomed. Eng.*, 2011, vol. 58, no. 9, pp. 2456–2466

- 8 Bortolotti, D., Bartolini, A., Weis, C., *et al.*: 'Hybrid memory architecture for voltage scaling in ultra-low power multi-core biomedical processors'. Design, Automation and Test in Europe Conf. and Exhibition (DATE), 2014, 2014, pp. 1–6
- 9 Massé, F., Van Bussel, M., Serteyn, A., *et al.*: 'Miniaturized wireless ECG monitor for real-time detection of epileptic seizures'. *ACM Trans. Embedded Comput. Syst. (TECS)*, 2013, **12**, (4), p. 102
- 10 Milosevic, J., Dittrich, A., Ferrante, A., *et al.*: 'Risk assessment of atrial fibrillation: a failure prediction approach'. Computing in Cardiology Conf.(CinC), 2014, 2014, pp. 801–804
- 11 Sörmmo, L., Laguna, P.: 'Bioelectrical signal processing in cardiac and neurological applications' (Academic Press, Burlington, USA, 2005)
- 12 Chou, C.C., Tseng, S.Y., Chua, E., *et al.*: 'Advanced ECG processor with HRV analysis for real-time portable health monitoring'. Consumer Electronics- Berlin (ICCE-Berlin), September 2011, pp. 172–175
- 13 Weckx, P., Kaczer, B., Toledano-Luque, M., *et al.*: 'Implications of BTI- induced time-dependent statistics on yield estimation of digital circuits'. *IEEE Trans. Electron. Devices*, 2014, **61**, (3), pp. 666–673
- 14 Bortolotti, D., Mamaghanian, H., Bartolini, A., *et al.*: 'Approximate compressed sensing: ultra-low power biosignal processing via aggressive voltage scaling on a hybrid memory multi-core processor'. Proc. of 2014 IEEE Int. Symp. on Low Power Electronics and Design (ISLPED 2014), EPFL-CONF-200128, IEEE/ACM Press, 2014, vol. 1, no. pp. 40–45
- 15 Verma, N., Chandrakasan, A.P.: 'A 256 kb 65 nm 8 T subthreshold SRAM employing sense-amplifier redundancy'. *IEEE J. Solid-State Circuits*, 2008, **43.1**, pp. 141–149
- 16 Di Carlo, S., Savino, A., Scionti, A., *et al.*: 'Influence of parasitic capacitance variations on 65 nm and 32 nm predictive technology model SRAM core-cells'. IEEE 17th Asian Test Symp. (ATS), November, 2008
- 17 Sanchez-Macian, A., Reviriego, P., Maestro, J.A.: 'Hamming SEC-DAED and extended hamming SEC-DED-TAED codes through selective shortening and bit placement'. *IEEE Trans. Device Mater. Reliab.*, 2013, **14**, (1), pp. 574–576
- 18 Basu, S.S., Garcia del Valle, P., Ansaloni, G., *et al.*: 'Heterogeneous error-resilient scheme for spectral analysis in ultra-low power wearable electrocardiogram devices'. IEEE Annual Symp. on VLSI, 2015
- 19 Wang, A., Chandrakasan, A.: 'A 180 mV FFT processor using subthreshold circuit techniques'. Solid-State Circuits Conf., 2004, vol. 1, pp. 292–299
- 20 Ashouei, M., Hulzink, J., Konijnenburg, M., *et al.*: 'A voltage-scalable biomedical signal processor running ECG using 13pJ/cycle at 1 MHz and 0.4 V'. 2011 IEEE Int. Solid-State Circuits Conf. Digest of Technical Papers (ISSCC), 2011, pp. 332–334
- 21 Muralimanohar, N., Balasubramonian, R., Jouppi, N.P.: 'CACTI 6.0: A tool to model large caches'. (HP Laboratories, Chicago, USA, 2009), pp. 22–31
- 22 Rajoub, B.: 'An efficient coding algorithm for the compression of ECG signals using the wavelet transform'. *IEEE Trans. Biomed. Eng.*, 2002, **49.4**, pp. 355–362
- 23 Akselrod, S., Gordon, D., Ubel, F., *et al.*: 'Power spectrum analysis of heart rate fluctuation: a quantitative probe of beat-to-beat cardiovascular control'. *Science*, 1981, **213**, (4504), pp. 220–222
- 24 Karakonstantis, G., Sankaranarayanan, A., Sabry, M.M., *et al.*: 'A quality-scalable and energy-efficient approach for spectral analysis of heart rate variability'. Design, Automation and Test in Europe Conf. and Exhibition (DATE), 2014, 2014, pp. 1–6
- 25 Karakonstantis, G., Sankaranarayanan, A., Burg, A.: 'Low complexity spectral analysis of heart-rate-variability through a wavelet based FFT'. Computing in Cardiology Conf. (CinC), 2012, September, 2012, pp. 285–288
- 26 Boichat, N., Atienza, D., Khaled, N.: 'Wavelet-based ECG delineation on a wearable embedded sensor platform' (BSN, Washington DC, USA, 2009)
- 27 Winchell, R.J., Hoyt, D.B.: 'Spectral analysis of heart rate variability in the ICU: a measure of autonomic function'. *J. Surg. Res.*, 1996, **63**, (1), pp. 11–16
- 28 PhysioBank Database. Available at <http://www.physionet.org/physiobank/>
- 29 The Cortex M3 Processor. Available at <http://www.arm.com/products/processors/cortex-m/cortex-m3.php>
- 30 Zuolo, L., Morandi, G., Zambelli, C., *et al.*: 'System interconnect extensions for fully transparent demand paging in low-cost MMU-less embedded systems'. Int. Symp. in System on Chip, 2013
- 31 Dong, X., Xu, C., Xie, Y., *et al.*: 'NVSIM: a circuit-level performance, energy, and area model for emerging nonvolatile memory'. *IEEE CAD*, 2012, **31**, (7), pp. 994–1007
- 32 Liao, D., Barnes, R.W., Chambless, L.E., *et al.*: 'Age, race, and sex differences in autonomic cardiac function measured by spectral analysis of heart rate variability – the ARIC study'. *Am. J. Cardiol.*, 1995, **76**, (12), pp. 906–12