# Computation as Estimation: A General Framework for Robustness and Energy Efficiency in SoCs

Sriram Narayanan, Girish Vishnu Varatkar,
Douglas L. Jones, *Fellow, IEEE*, and
Naresh R. Shanbhag, *Fellow, IEEE*

*Abstract*—Traditional integrated circuit design achieves error-free operation by designing with margins (clock frequency and supply voltage) and/or including hardware replication and recomputation, which may counter the full energy and area benefits of aggressive technology scaling. It is thus desirable that modern systems-on-chip (SoCs) permit hardware errors while maintaining robust system-level performance. Treating hardware errors as *computational noise* and extending traditional estimation theory to include practical SoC design constraints yields a novel and general design optimization framework. This work demonstrates the breadth of applicability of the estimation-theoretic framework for system design by showcasing two different application classes that demonstrate 36% to 50% power reduction.

*Index Terms*—Applications of statistical signal processing techniques, HDW-LPWR, low-power signal processing techniques and architectures.

## I. INTRODUCTION

Increased transistor density, limited improvements in battery technology, and demand for high performance have made power reduction an important concern when designing SoCs. At the same time, technology scaling in the nanometer regimes has introduced numerous sources of nonidealities such as process variations and soft errors [1]. In order to reap the full benefits of aggressive technology scaling, it is desirable that modern SoCs be designed to tolerate the maximum number of hardware errors and incur the least amount of overhead while guaranteeing minimum acceptable performance specifications of the application (e.g., bit error rate), as identified by the authors in [2].

Reliability and power reduction are often related problems because the ill effects of some CMOS nonidealities may be overcome by overprovisioning power or operating at slower clock frequencies. However, traditional design methods address either power reduction or error-tolerance in isolation. Conventional dynamic voltage/frequency scaling (DVFS) [3] methods have been shown to offer power benefits. However, with increasing IC density and thinning design guardbands, they will increase the risk of timing violations, requiring methods that simultaneously address robustness. Error-tolerant approaches often mitigate hardware errors at the cost of very significant power, area, or design overheads. Redundant residue number systems provide

circuit-level error tolerance [4], and approaches such as $N$-modular redundancy [shown in Fig. 1(a)] provide system-level tolerance. But the computational power, area, and cost overheads of such techniques are often substantial. Furthermore, $N$-modular redundancy methods may not be applicable for overcoming timing errors that depend on the input. Even with sufficient architectural diversity to decorrelate timing errors across modules, the $N$-times area and power overhead of this approach will be prohibitive for many applications. In contrast, the Class I systems presented in this work are able to address timing errors with an area overhead of around 11% [5]. Other system-level techniques such as algorithmic noise tolerance (e.g., [6] and [7]), are designed to address particular error models or input statistics, and usually lack a notion of system-level optimality. For the motion-estimation kernel in Class II systems presented here, area overheads of an ANT-based redesign is around 26% [8]. In contrast, the methods presented here incur no overhead associated with recomputation, and do not require redesign of the video processing block. The better-than-worst-case design methods [9], canary replica feedback [10], crystal ball flip-flops [11], and BISER [12] employ gate-level redundancy to overcome every instance of timing errors, and ignore any available system-level error masking. Correcting errors at the latch level imposes significant redesign, and large overheads of power and area. Applying Razor to the Class 1 systems will incur up to 27% area overhead.[1] Additionally, this overhead does not include the cost of recomputation that is associated with such techniques.

The design philosophy adopted here is to use existing computational blocks and develop system-level methods to achieve robustness. These methods do not introduce architectural modifications to the computational blocks of the conventional system; instead, they either modify the postprocessing blocks that already exist in the conventional design, or expose the application to hardware errors. This is in contrast to other system-level techniques such as ANT where a reliable estimator block is constructed to operate in parallel with the main block, thereby representing a significant architectural redesign.

Our system design is also indifferent to particular hardware-error models, which are often hard to obtain. We develop a formalization that defines a notion of optimal robustness and guides the system designer to exploit all available statistical information. This framework treats hardware errors as computational noise that is analogous to system/measurement noise, and applies results from estimation theory to design robust SoCs [see Fig. 1(c)].

## II. ESTIMATION THEORY FOR COMPUTATION

Traditional estimation theory deals with the problem of optimally determining an underlying parameter or signal based on a set of noisy measurements by minimizing the average risk of misestimating it. When viewing computation as a special case of estimation, soft errors and hardware errors due to process variations, voltage or frequency scaling represent a new source of noise in addition to measurement or system noise. This analogy enables us to leverage many well-known algorithms in estimation theory. Traditional estimation theory is unconcerned with engineering constraints of arriving at an estimate. However, in the present context, we are interested in the complexity of the estimator, and therefore require a nontrivial extension of estimation theory to robust SoC design [13].

In the proposed view of computation, we treat the subsystems of a complex SoC as providing noisy estimates of the overall computation. The problem is to optimally find the final result based on these

---

[1]This conservatively assumes that the computational overhead of the filter banks is 90% and the overhead of the flip-flops in each filter bank that need Razor protection is 30%; see gate-complexity details in [5].
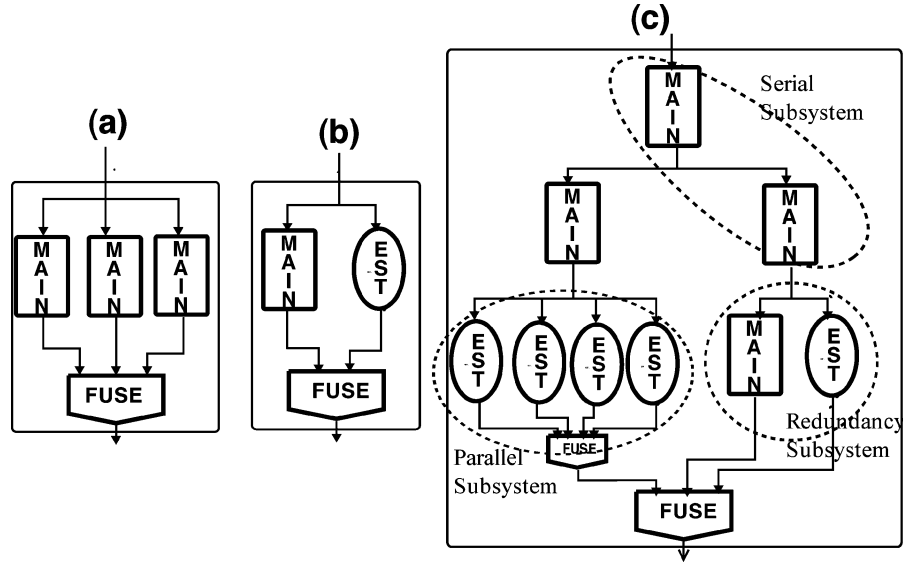
Fig. 1. (a) Traditional NMR systems replicate computation a number of times and discard erroneous outputs using a majority voter. (b) ANT systems employ a lower-complexity estimator in place of replicated computation. (c) Our novel view of complex SoCs identifies serial subsystems, parallel subsystems (main computation decomposed into a set of estimators), and subsystems with explicitly built-in redundancy blocks (such as ANT). This approach generalizes previous error-tolerance mechanisms to fully exploit statistical correlations present within the SoC. (a) NMR. (b) ANT. (c) Complex SoC.

estimates. In the context of SoC design, this problem translates to minimizing the computational risk of misestimating the final result from a vector of error-prone subsystem outputs, $\vec{Y} = \{Y_1, \ldots, Y_N\}$. Our design philosophy is to gain maximum amount of robustness by using only these already available subcomputation results.

### A. Computation Cost Function

Estimation-theoretic system cost functions are used to define computational risk. These cost functions usually depend on the specific application. For a given computation, let $\theta$ denote the desired result of computation, and $\hat{\theta}(\vec{Y})$ denote the error-prone output of the SoC. The application-specific system cost is determined by $\theta$ and $\hat{\theta}(\vec{Y})$. The risk function is the expectation of this cost function over the random variable $\theta$ and the observed variables $\vec{Y}$.

Different computation problems may call for different cost functions. Squared-error cost functions may be appropriate for the many signal-processing applications in which it is common to minimize the mean squared error (MSE), given by $C[\hat{\theta}(\vec{Y}), \theta] = (\hat{\theta}(\vec{Y}) - \theta)^2$. For general-purpose computing systems that tolerate errors below some limit, $\Delta$, we suggest the following 0–1 cost function:

$$C[\hat{\theta}(\vec{Y}), \theta] = \begin{cases} 1, & |\hat{\theta}(\vec{Y}) - \theta| \geq \Delta \\ 0, & |\hat{\theta}(\vec{Y}) - \theta| < \Delta \end{cases}. \qquad (1)$$

Another common cost function is the absolute error, $C[\hat{\theta}(\vec{Y}), \theta] = |\hat{\theta}(\vec{Y}) - \theta|$.

### B. Architectural Design Space and Constraints

Estimation-theoretic design of SoCs presents unique issues that need to be addressed. Cost and technology considerations may limit the system designer to a finite number of architectures or technology choices. For a given architecture, the choices for various operating parameters such as supply voltage, clock frequency, and register word-lengths may also be constrained. Because these parameters have a direct impact on both the system power consumption and hardware error rate, it is important to optimally choose them. Therefore, we need a general estimation-theoretic framework that optimizes system

performance or power consumption while accounting for design constraints.

### C. Canonical Problems of the Estimation-Theoretic Framework

The estimation-theoretic system design optimization may be stated in two canonical problems. The first problem seeks to minimize power consumed in arriving at a computational result (analogous to an estimate), while constraining the average system accuracy (analogous to average risk). Let $\theta$ be the result being computed that belongs to some set $\Lambda$, and let $\hat{\theta}$ be the estimator that operates on the input, $\vec{Y}$. We can state this problem as follows.

**Problem 1: Performance-constrained system.**

$$\hat{\theta}(\vec{Y}) = \arg\left\{ \min_{\theta \in \Lambda, A \in \mathbb{A}} P(\hat{\theta}(\vec{Y})) \right\}$$

subject to

$$E_\theta\{C[\hat{\theta}(\vec{Y}), \theta]\} \leq C_{\text{Target}}$$
$$\mathbb{A} = \{A_1(\vec{\lambda}), A_2(\vec{\lambda}), \ldots\}$$

where $\mathbb{A}$ is the set of architectural choices, the vector $\vec{\lambda}$ defines the tunable parameters of an architecture (e.g., supply voltage, clock frequency), and the function, $P(\cdot)$, computes the power consumed in arriving at an estimate.

Problem 2 seeks to minimize the average risk incurred in misestimating the result while constraining the power consumed to be within budget. For a battery-operated system, for example, the amount of stored energy may be used to arrive at a power budget. Cooling and packaging costs may define this budget for tethered systems.

**Problem 2: Power-constrained system.**

$$\hat{\theta}(\vec{Y}) = \arg\left\{ \min_{\theta \in \Lambda, A \in \mathbb{A}} E_\theta\{C[\hat{\theta}(\vec{Y}), \theta]\} \right\}$$

subject to

$$P(\hat{\theta}(\vec{Y})) \leq P_{Budget}$$
$$\mathbb{A} = \{A_1(\vec{\lambda}), A_2(\vec{\lambda}), \ldots\}.$$

The canonical problems of the estimation-theoretic framework are very general, and different applications may lead to very different es-

timation problems. *The value of the abstraction presented through the canonical problems is that it allows the system designer to quickly recognize relevant results from estimation theory and apply them to many important system-design problems.* We present two classes of systems in which this framework offers marked improvements in robustness. The chosen classes were deliberately made general. The following sections also detail specific applications from each of the classes. The specific solution methodology that results from the estimation-theoretic problem statement is identified, and relevant results that highlight power reduction and system robustness are then presented.

## III. CLASS I: LINEAR SYSTEMS COMPRISING PARALLEL SUBSYSTEMS

In the first class of systems, the input may be divided into subsets that are processed by a set of parallel processing elements (PEs), $\{\mathrm{PE}_1, \ldots, \mathrm{PE}_N\}$. The results of the PEs, $\{y_1, \ldots, y_N\}$, comprise the noisy measurements, $\vec{Y}$, of the final output, $\theta$. A fusion block uses these noisy subcomputations to produce an estimate of the final computation. The PEs may be designed such that their total complexity equals that of the overall computation; i.e., the new system does not use any explicit redundant computation. Alternatively, the complexity of the PEs may exceed that of the overall system; i.e., we introduce redundant computation in the parallelization process. The unifying characteristic of this class of systems is that robustness is achieved by means of postprocessing the outputs of the PEs.

Redundancy-free parallel systems commonly occur in many signal-processing applications such as polyphase FIR filtering, and applications in computational geometry in a variety of fields such as pattern recognition, image processing, and computer graphics [14]. Examples of redundancy-aided systems range from traditional $N$-modular redundancy systems in which the parallel subcomputations are identical and each performs the overall computation, to algorithmic noise tolerance (ANT) systems where one lower-complexity parallel subcomputation checks the main computation.

We adopt a mean squared-error cost for this class. Thus, for each $\mathrm{PE}_i$, $C(\theta_i, y_i) = E_{\theta_i}\{y_i - \theta_i\}^2$. For applications in this class, the system cost function, $C_{\mathrm{sys}}(\theta, \hat{\theta})$, is given by

$$C_{\mathrm{sys}}(\theta, \hat{\theta}(\vec{Y})) = E\{C(\theta_i, y_i)\} \qquad (2)$$

where $C(\theta_i, y_i)$ is the cost function of the $i$th processing element. The expectation in (2) is taken over the random variable representing the event that a subcomputation is declared to be in error. If the subcomputations are identical, then it simplifies to a simple average. The system power consumption, $P_{\mathrm{sys}}$, is given by

$$P_{\mathrm{sys}} = \sum_{i=1}^{m} P_{\mathrm{PE}_i} + P_{\mathrm{fusion}} \qquad (3)$$

where $P_{\mathrm{PE}_i}$ is the power consumed by $\mathrm{PE}_i$ and $P_{\mathrm{fusion}}$ is the power consumed by the fusion center.

The estimation-theoretic design problem is to minimize the maximum asymptotic variance of the computation result while constraining the power consumption to a given value. This can be stated as

$$\hat{\theta}(\vec{Y}) = \arg \min_{\theta \in \Lambda} C_{\mathrm{sys}}(\theta, \hat{\theta}(\vec{Y}))$$
$$\text{subject to}$$
$$P(\hat{\theta}) \leq P_{\max} \quad \& \quad \mathbb{A} = \mathbb{V} \qquad (4)$$

where $P_{\max}$ is the maximum allowable power consumption and $\mathbb{V}$ is the set of supply voltages at which the system may be operated. The problem in this application corresponds to Problem 2 of the estimation-theoretic framework. Solving the performance-constrained Problem 1 for this class of problems involves finding the architectural parameters

from an ordered set that minimizes power while meeting performance constraints.

The particular form of the optimal estimate of the computational result will depend on assumed probabilistic models of the hardware errors and system noise present in $\vec{Y}$. Some researchers have tried to model errors caused by timing violations (for example, see [15]), but hardware errors may be time varying (for example, dependent on chip temperature) or very different for different types of failure mechanisms. To accommodate this, we instead design for the worst case and develop techniques that are robust to deviations of the actual errors from assumed models.

In the presence of hardware errors that occur with probability $\epsilon$, the output can be modeled as random variables drawn from a class of distributions that is Gaussian with probability $(1 - \epsilon)$ and some unknown distribution with probability $\epsilon$ for some $0 < \epsilon < 1$, as follows:

$$P_\epsilon = \{F | F = (1 - \epsilon)\Phi + \epsilon H, H \in \mathcal{R}\} \qquad (5)$$

where $\Phi$ is the class of standard normal distributions (a good model for input/measurement noise) and $H$ is the class of arbitrary densities with zero mean and finite but unbounded variance (used to model hardware errors such as those caused by timing violations), and $\mathcal{R}$ is the set of all probability measures on the real line. An element of the set $F$ is the probability distribution of the computational noise (i.e., a combination of system/input noise and hardware errors caused by power-reduction schemes). It is important to note that because the exact probability model of the hardware errors may be unknown and time-varying, this mixture model allows us to design for the worst-case error model. In this way, this design methodology is indifferent to any specific error model.

The theory of robust statistics offers a solution to the estimation problem in this class. The author in [16] derives the maximum-likelihood estimate for the least informative distribution from the class described in (5). We find this SoC design problem to be a novel application of this well-known statistical signal processing technique. The fusion block used in this class of applications computes this optimal estimate of the final computation result.

As a specific example as in [17] we use voltage overscaling (VOS) as an error-prone power-reduction technique. VOS results in increased errors due to the slower computations not meeting timing constraints, and the output is contaminated by a mixture of the noise that may be already present in the input and the large-magnitude VOS errors. Studies performed in [6] link the voltage overscaling factor, $k_{\mathrm{VOS}}$ (i.e., the factor by which the supply voltage is reduced with respect to the critical supply voltage), and the probability of hardware error $\epsilon$. Because such hardware errors are an artifact of voltage overscaling and the input may be modeled as an independent random process, we assume that they are uncorrelated with the input.

The preceding description of this class has been fairly abstract. The following application case-study of a redundancy-free system illustrates a concrete instance of estimation-theoretic design; see [13] for a redundancy-aided application. Spread-spectrum communication systems use finite impulse response (FIR) filters as matched filters for pseudonoise (PN) codes when identifying a user in a multiple access channel. The peaks in the output of the matched filters are used for detection and synchronization of PN sequences. This code acquisition is a computationally critical block in a spread-spectrum communication receiver [18]. Fig. 2(a) shows an implementation of a parallel matched filter. In the event of a successful match, the signal components of the outputs of the filter banks of a polyphase matched-filter are completely correlated. Such a statistical relationship between the subcomputations of the filtering operation can be exploited to gain robustness to errors and power reduction without the need for redundant computation. A
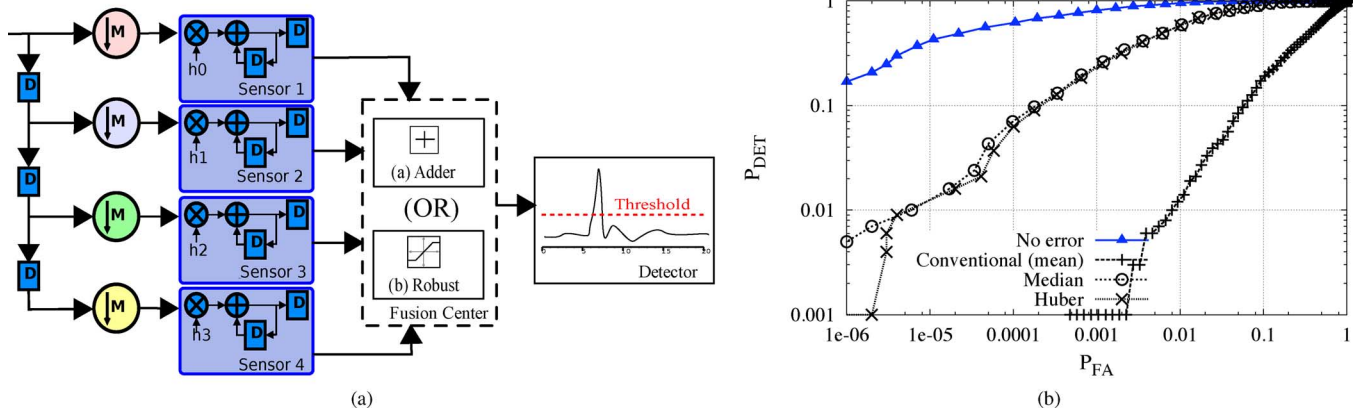
Fig. 2. The robust detector offers a better probability of detection for a fixed probability of false alarm. (a) Robust PN Matched Filter. (b) Receiver Operating Characteristics ($k_{\mathrm{VOS}} = 0.75$).

robust fusion block that implements the optimal estimator defined in [16] is used to process outputs of the filter banks.

Prior work in [5] used circuit-level simulations to show that the robust PN acquisition system can reduce power consumption by up to 36% when the median is used to approximate the optimal estimator. The architecture in [5] consisted of $N$ identical processors and was chosen only because of the low design cost and commercial availability of such multiprocessors, and not due to any notion of optimal design. By invoking our estimation-theoretic framework, we can show that the architecture chosen in [5] is indeed optimal in the sense of minimizing system cost. The following theorem defines the optimal design of the subsystems for this robust system.

**Theorem 1:** *A robust $k$-phase matched filter implementation produces an estimate with minimum variance if the filter banks are of equal length.*

*Proof:* Follows from convexity. [19] ∎

Therefore, the architectural constraint set, $\mathbb{A}$, is the discrete set of voltages, $\mathbb{V}$, at which the $N$ identical processing elements may be operated.

Fig. 2(b) compares the performance of the median, the optimal maximum-likelihood estimate, and a brute-force solution. Each plot also shows the performance of an error-free (nonoverscaled) matched filter as a comparison. The performance is measured by the best achievable probability of detection, $P_{\mathrm{DET}}$, for a given probability of false alarm, $P_{\mathrm{FA}}$.

## IV. CLASS II: SYSTEMS COMPRISING SERIAL SUBSYSTEMS

Class II systems exploit the fact that many applications are inherently robust to some number of hardware errors in computation. The estimation-theoretic framework allows the designer to trade off hardware errors and system noise in a joint manner.

In this class, the overall computation can be divided into smaller subcomputations, each of which is processed by one or more processing elements chosen from the set $\{\mathrm{PE}_1, \ldots, \mathrm{PE}_N\}$. The PEs are serially linked, and produce outputs $\vec{Y} = \{y_1, \ldots, y_N\}$ of varying degrees of error tolerance. The subcomputation results may be subject to different cost functions $\{C_1(\theta_1, y_1), \ldots, C_m(\theta_N, y_N)\}$. Here $\{\theta_1, \ldots, \theta_N\}$ denote the desired outputs of the PEs. The system cost function, $f(\cdot)$, is a known function of the individual cost functions, i.e.,

$$C(\theta, \hat{\theta}(\vec{Y})) = f\left(C_1(\theta_1, y_1), \ldots, C_N(\theta_m, y_N)\right) \qquad (6)$$

where $\theta$ is the desired system output and $\hat{\theta}(\vec{Y})$ is the estimate produced by this serial system. The performance of the overall system is

sensitive to the performance of the individual components by varying amounts, and hence the performance of some PEs can be worsened without significant degradation in the overall system performance. The system power consumption, $P_{\mathrm{sys}}$, is given by

$$P_{\mathrm{sys}} = \sum_{i=1}^{m} P_{\mathrm{PE}_i} \qquad (7)$$

where $P_{\mathrm{PE}_i}$ is the power consumed by the $i$th PE.

This class of applications naturally lends itself to the performance-constrained Problem 1. Gradient-descent approaches often work well for this class of problems. Optimizing this class of applications using the power-constrained Problem 2 involves optimally allocating power amongst the PEs to minimize the overall system cost.

The remainder of this section describes an application in this class in which the estimation-theoretic framework offers up to 50% power reduction.

In wireless video transmission systems, the raw video input is fed into a video encoder, the video encoder is linked to a channel encoder, and the channel encoder is linked to a wireless transmitter. Modern video encoding standards such as the H.264/AVC have made substantial advances in compression, but only at the expense of greatly increased computational complexity. Joint source coding and power-management techniques, such as those in [20], have been shown to optimize only the communication power consumption in such applications. By adopting aggressive power-reduction techniques at the video encoder and compensating for any resulting computational errors, mobile video systems can gain additional power reduction. But timing errors can cause the encoder to choose suboptimal motion vectors, thereby lessening its compression efficiency [21]. The estimation-theoretic framework allows optimally trading-off power reduction and compression efficiency in scenarios with low or moderate communication demand to lower overall system power consumption.

The estimation-theoretic cost function for the video encoder is the output bit-rate for a fixed video quality (in PSNR). The IBM 130 nm process model simulations were used to characterize the increase in gate delay of full-adders as a function of voltage overscaling factor, $k_{\mathrm{VOS}}$; [22] provided models for computing the bit-rate. The cost function for the channel encoder is the packet error rate suffered by the receiver. The 180 nm, 1.8 V CMOS reconfigurable Reed–Solomon (RS) encoder presented in [23] offers a quantifiable means of balancing redundancy and power consumption. The cost of the transmitter is the signal-to-noise ratio seen by the receiver. The power-amplifier design presented in [23] provides a model for the power consumption of a wireless transmitter. The system-level cost function is the average
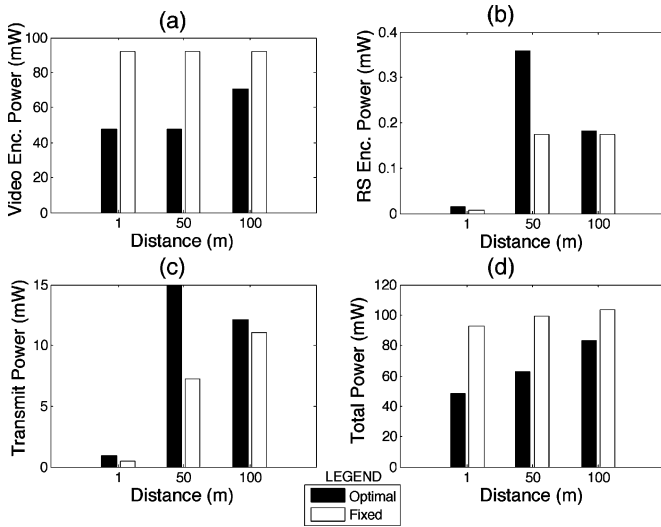
Fig. 3. The power consumed by the video encoder, the channel encoder, wireless transmitter and the total system power consumption are compared in Figures (a), (b), (c), and (d), respectively. Relative power allocation for the different subsystems varies depending on the range of communication. (a) Video Encoding. (b) Channel Encoding. (c) Transmission. (d) Total System Power.

end-to-end distortion, and the exact relationship of this quantity on the source encoder, channel, and transmitter depends on any error-concealment strategies used at the decoder; see [20] for a detailed discussion.

The elements of the architectural constraint, $\mathbb{A}$, each represent a different choice for the supply voltage for motion estimation and number of correctable symbols, $t$, of the channel coder. Each choice in this set can be configured using a continuously variable transmit power, $P_t$. Therefore,

$$\mathbb{A} = \mathbb{K}_{\mathrm{VOS}} \times \mathbb{T} \qquad (8)$$

where $\mathbb{K}_{\mathrm{VOS}}$ is the set of choices for the supply-voltage overscaling factor of the motion estimation and $\mathbb{T} = \{1, \ldots, \text{number of symbols in a packet}\}$. The transmit power, $P_t$, is a continuous parameter of the architectural constraint set, $\lambda = P_t$.

The problem of optimally designing the mobile video communication system corresponds to Problem 1 (minimize power subject to a cost constraint) that finds the element $A \in \mathbb{A}$ defined by optimal choice of $k_{\mathrm{VOS}} \in \mathbb{K}_{\mathrm{VOS}}$ and $t \in \mathbb{T}$. The objective function comprises the video-encoder power consumption, $P_{\mathrm{venc}}(k_{\mathrm{VOS}})$, the power consumed by the channel encoder, $P_{\mathrm{chenc}}(t, f_{bit})$, and the power consumed by the transmitter, $P_{\mathrm{xmit}}(P_t, L)$. The parameter $f_{bit}$ is the bit rate, and $L$ is the range of communication. In our model, the distortion constraint, $\mathcal{D}_{\mathrm{Target}}$, is automatically satisfied by fixing the quantization parameter of the encoder at a satisfactory value and adopting a retransmission scheme for communication. A gradient-descent algorithm was used to optimize over transmit power, $P_t$. Since the simulation results for the gate delay were obtained at a discrete set of values for $k_{\mathrm{VOS}}$ and the number of correctable symbols, $t$, we searched over these parameters to find the optimum.[2] Fig. 3 shows a comparison of the optimized system (left bar plots) relative to a system that does not allow voltage overscaling at the video encoder, but optimally chooses only the transmit power and RS encoder redundancy (right bar plots). Over short communication links, the system power is dominated by the video encoder, and optimal choice of $k_{\mathrm{VOS}}$ yields around 35%–50% power savings. Over longer communication links, where the communication power dominates, system power savings are around 18%.

[2]Many practical systems only allow operation at a discrete set of supply voltages.

## V. CONCLUSION

Aggressive power-reduction schemes that potentially deplete design guardbands have recently received significant attention as a means of extending Moore's law [2], [6], [9]. The estimation-theoretic framework presented in this correspondence defines a notion of optimal design for these aggressive power-reduction techniques by exploiting redundancy present in the application and the underlying system. The mathematical formalization leverages powerful algorithms from estimation theory by unifying the problem statement for diverse classes of SoCs. The application classes chosen here are meant to highlight the diversity of applications and the generality of the proposed estimation-theoretic design framework. There are no doubt other new classes and mixed classes that can similarly benefit from this framework. We envision future systems that are capable of monitoring and learning the probability of error occurrence and are able to use this in an online fashion to balance power reduction and error resiliency.

## REFERENCES

[1] C. Constantinescu, "Trends and challenges in VLSI circuit reliability," *IEEE Micro.*, vol. 23, no. 4, pp. 14–19, Jul. 2003.

[2] M. Breuer *et al.*, "Defect and error tolerance in the presence of massive numbers of defects," *IEEE Des. Test. Comput.*, vol. 21, no. 3, pp. 216–227, May–Jun. 2004.

[3] T. Pering, T. Burd, and R. Brodersen, "The simulation and evaluation of dynamic voltage scaling algorithms," in *Proc. IEEE Int. Symp. Low Power Electronics Design (ISLPED)*, Aug. 1998, pp. 76–81.

[4] W. K. Jenkins, "Design of error checkers for self-checking residue number arithmetic," *IEEE Trans. Comput.*, vol. C-32, pp. 388–396, Apr. 1983.

[5] G. V. Varatkar, S. Narayanan, N. Shanbhag, and D. L. Jones, "Stochastic networked computation," *IEEE Trans. VLSI Syst.*, vol. PP, Oct. 2009.

[6] B. Shim, S. R. Sridhara, and N. R. Shanbhag, "Reliable low-power digital signal processing via reduced precision redundancy," *IEEE Trans. VLSI Syst.*, vol. 12, pp. 497–510, May 2004.

[7] J. W. Choi, B. Shim, A. C. Singer, and N. I. Cho, "Low-power filtering via minimum power soft error cancellation," *IEEE Trans. Signal Process.*, vol. 55, no. 10, pp. 5084–5096, Oct. 2007.

[8] G. V. Varatkar and N. R. Shanbhag, "Error-resilient motion estimation architecture," *IEEE Trans. VLSI Syst.*, vol. 16, pp. 1399–1412, Oct. 2008.

[9] T. Austin and V. Bertacco, "Deployment of better than worst-case design: Solutions and needs," in *Proc. IEEE Int. Conf. Computer Design (ICCD)*, Washington, DC, 2005, pp. 550–558.

[10] J. Wang and B. Calhoun, "Canary replica feedback for near-DRV standby VDD scaling in a 90 nm SRAM," in *Proc. IEEE Custom Integrated Circuits Conf. (CICC)*, Sep. 2007, pp. 29–32.

[11] L. Anghel and M. Nicolaidis, "Cost reduction and evaluation of a temporary faults-detecting technique," in *Proc. Design, Automation, Test in Europe Conf. Exhib.*, 2000, pp. 591–598.

[12] S. Mitra *et al.*, "Logic soft errors: A major barrier to robust platform design," in *Proc. IEEE Int. Test Conf. (ITC)*, Nov. 2005, pp. 687–696.

[13] S. Narayanan *et al.*, "Computation as estimation: Estimation-theoretic IC design improves robustness and reduces power consumption," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing (ICASSP)*, Apr. 2008, pp. 1421–1424.

[14] S.-J. Oh and M. Suk, "Parallel algorithms for geometric searching problems," in *Proc. ACM/IEEE Conf. Supercomput.*, Reno, NV, 1989, pp. 344–350.

[15] S. Sarangi *et al.*, "VARIUS: A model of process variation and resulting timing errors for microarchitects," *IEEE Trans. Semicond. Manuf.*, vol. 21, no. 1, pp. 3–13, Feb. 2008.

[16] P. Huber, *Robust Statistics*. New York: Wiley, 1981.

[17] R. Hegde and N. R. Shanbhag, "Soft digital signal processing," *IEEE Trans. VLSI Syst.*, vol. 9, pp. 813–823, 2001.

[18] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. New York: Cambridge Univ. Press, 2005.

[19] S. Narayanan, "Estimation-theoretic framework for robust and energy-efficient system design," Ph.D. dissertation, Univ. of Illinois, Urbana-Champaign, IL, 2010.

[20] Y. Eisenberg *et al.*, "Joint source coding and transmission power management for energy efficient wireless video communications," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 12, pp. 411–424, Jun. 2002.

[21] I. S. Chong and A. Ortega, "Dynamic voltage scaling algorithms for power constrained motion estimation," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing (ICASSP)*, Apr. 2007, vol. 2, pp. 101–104.

[22] G. V. Varatkar, S. Narayanan, N. R. Shanbhag, and D. L. Jones, "Trends in energy-efficiency and robustness using stochastic sensor network-on-a-chip," in *Proc. 18th ACM Great Lakes Symp. VLSI (GLSVLSI)*, 2008, pp. 351–354.

[23] S. Appadwedula, M. Goel, N. R. Shanbhag, D. L. Jones, and K. Ramchandran, "Total system energy minimization for wireless image transmission," *J. VLSI Signal Process.*, vol. 27, no. 1-2, pp. 99–117, Feb. 2001.

# Unbiased Model Combinations for Adaptive Filtering

Suleyman S. Kozat, Andrew C. Singer, Alper Tunga Erdogan, and Ali H. Sayed

*Abstract*—In this paper, we consider model combination methods for adaptive filtering that perform unbiased estimation. In this widely studied framework, two adaptive filters are run in parallel, each producing unbiased estimates of an underlying linear model. The outputs of these two filters are combined using another adaptive algorithm to yield the final output of the system. Overall, we require that the final algorithm produce an unbiased estimate of the underlying model. We later specialize this framework where we combine one filter using the least-mean squares (LMS) update and the other filter using the least-mean fourth (LMF) update to decrease cross correlation in between the outputs and improve the overall performance. We study the steady-state performance of previously introduced methods as well as novel combination algorithms for stationary and nonstationary data. These algorithms use stochastic gradient updates instead of the variable transformations used in previous approaches. We explicitly provide steady-state analysis for both stationary and nonstationary environments. We also demonstrate close agreement with the introduced results and the simulations, and show for this specific combination, more than 2 dB gains in terms of excess mean square error with respect to the best constituent filter in the simulations.

*Index Terms*—Adaptive filtering, gradient projection, least-mean fourth, least-mean square, mixture methods.

## I. INTRODUCTION

We investigate unbiased mixture methods to combine outputs of two adaptive filtering algorithms operating in stationary and nonstationary environments. The objective is to achieve a steady-state mean-square error (MSE) better than, or at least as good as, each individual adaptive branch by exploiting the cross correlation structure between them

through an adaptive combining scheme. We may achieve an unbiased output through the use of the convex or affine combination constraints on the combination weights. We focus on steady-state results for stationary and certain nonstationary data models, however, the transient analysis of the algorithms can be derived using similar methods. Furthermore, although we only use stochastic gradient updates to train the combination weights, one can extend these algorithms to other methods, such as those based on Newton or quasi-Newton updates.

The structure we consider consists of two stages [1], [2]. In the first stage, we have two adaptive filters, working in parallel, to model a desired signal. These adaptive filters have the same length, however, each may use a different adaptation algorithm. We also require that these constituent filters produce unbiased estimates of the underlying model. The desired signal has a random walk formulation to represent both stationary and nonstationary environments [3]. A more precise problem formulation is given in Section II. The second stage of the model is the combination stage. Here, the outputs of the adaptive filters in the first stage are linearly combined to yield the final output. We only consider combination methods that produce unbiased final estimates of the underlying model. A sufficient condition to satisfy this requirement is to assume that the second stage coefficients sum up to one at all times, i.e., affine combinations. In addition to unbiasedness, the combination coefficients can be further constrained to be nonnegative, which corresponds to the case of convex combination. We consider both of these cases.

The framework where multiple adaptive algorithms are combined using an unbiased linear combination with the goal of improving the overall performance has recently attracted wide interest [1], [4], and [5], following the result in [1] that the convex combinations can improve the resulting MSE performance. The requirement on unbiasedness may be motivated from some problem-specific constraints as well as implementation related issues. The combination weights are usually trained using stochastic gradient updates, either after a sigmoid nonlinearity transformation to satisfy convex constraints [1], [4] or after a variable transformation to satisfy affine constraints [5]. There are also Bayesian inspired methods that have extensive roots in machine learning literature [2]. The methods in [1], [2], [4], and [5] combine filters using least-mean squares (LMS) or recursive least squares (RLS) updates (or unsupervised updates). As demonstrated in [1] and [4], mixtures of two filters using the LMS or RLS updates (or a combination of the two) with the convex methods yield combination structures that converge to the best algorithm among the two for stationary data. As demonstrated in [1], the cross correlation between *a priori* errors of the two LMS filters (or LMS and RLS filters in [4]) remains sufficiently high that it limits the combination performance and the optimal convex combination solution converges to only selecting one of the two outputs.

In this paper, we first quantify the achievable gains using convex or affine constrained combination weights in steady-state for stationary and nonstationary data. We also provide the optimal combination weights to yield these gains. We next demonstrate that the update given in [5, Eq. (45)] (which tries to simulate the unrealizable optimal affine combiner) is a stochastic gradient update with a single tap input regressor and derive its steady-state MSE for both stationary and nonstationary environments. Here, we refrain from making variable transformations and directly adapt the combination weights using stochastic gradient updates. However, to preserve convexity or affinity, after each update, we project each updated mixture weight vector back to the convex or affine space. These methods update the weights directly instead of using variable transformations [1], [4]. As a by product of our analysis, we demonstrate that the update in [5, Eq. (45)] is also a stochastic gradient projection update. As a specific example,