# A Self-Aware Architecture for PVT Compensation and Power Nap in Near-Threshold Processors

**Davide Rossi and Igor Loi**
University of Bologna

**Antonio Pullini**
ETHZ

**Christoph Müller and Andreas Burg**
EPFL

**Francesco Conti and Luca Benini**
University of Bologna
ETHZ

**Philippe Flatresse**
STMicroelectronics

*Editor's note:*

As technology edges closer to fundamental limits, variations of process parameters, operational voltage, and temperature (PVT variations) have to be accounted for. This paper proposes to make higher levels aware of these variations and bring them under system control by using software controlled body biasing. PVT variations can thus be exploited to reduce power and energy consumption.

—Axel Jantsch, TU Wien

■ **INTERNET OF THINGS** (IoT) end-nodes require extreme energy efficiency in active state and idle state, coupled with software-programmable architectures to cope with complex and highly dynamic near-sensor data analytics algorithms. Nanometer CMOS technologies, in combination with aggressive voltage scaling and parallel processing, have enabled major improvements in active energy efficiency [1]. However, efficiency gains are ulti-

mately limited in near-threshold (NT) [2] operation by

· conservative design margins for countering the increased sensitivity to PVT variations and
· the dominance of leakage currents, both

in active state and during frequent and short idle periods, which require full state retention.

Widely used techniques to deal with variations in NT involve the design of architectures able to probe the PVT and aging conditions of the circuit and to provide a hardware feedback to knobs exposed at system level. Automatically compensating the variations by adjusting the supply [3] or body bias voltage [4], these architectures can reduce the PVT margins applied at design time, significantly improving the efficiency of the system (by 32% in this work). Among the various techniques, the Razor concept [3] guarantees a reliable operation by lowering the supply voltage to the point of first failure using speed monitors

within the respective circuit block. This eliminates design margins required to cope with both global and local on-chip PVT variations. The main drawback of this approach lies in the high intrusiveness in the design and in the area overhead, as it requires to manually insert the shadow latches and the transition detectors, as well as the related control logic, into the microarchitecture of the systems on chip (SoC). This makes the general adoption of Razor very challenging and not feasible when the IP cores, which compose the system architecture, are released as hard macro or encrypted netlists by IP providers, which is nowadays the most common way to design complex SoCs.

Such intrusive approach can be avoided in tiny (i.e., a few millimeters) near-threshold designs, where the main source of variation comes from die-to-die variation and thermal inversion, while on-chip variations (OCV) and temperature gradients are small [5], and voltage droop is nearly homogenous within the whole die area, being mainly dominated by the off-chip power supply network [6]. In this context, intrachip variations can still be handled with margining and decoupling capacitance (DECAP) cells [6], while global PVT status can be tracked with a less intrusive approach based on critical path monitors (CPMs). A CPM for a digital logic circuit is a replica of the critical path with the addition of some delay elements, making the monitor super-critical [7]. However, CPMs are not generic reusable components, since they have to be designed for a specific SoC. Moreover, they cannot provide the direct information when the critical path is through SRAMs, since the mixed-signal propagation path within SRAMs cannot be easily replicated nor properly tracked by digital logic cascades. To address this issue, some programmability capabilities can be added to CPMs to match the SoC critical path for a given voltage and maximum frequency [7]. Process monitoring boxes (PMBs) take CPMs one step further in terms of generality [8]. They contain arrays of ring oscillators with different characteristics (e.g., PMOS only, and NMOS only). After proper calibration, the readings of these "on-chip sensors" can be fused to provide a reliable indication of the operating corner that can be used as an input for operating point compensation.

While Razor-like schemes rely usually on hardware feedback loops [3], the modularity of CPMs

and PMBs, allows using both hardware [7] and software [9] strategies to close the loop. Moreover, ultra-low power (ULP) devices are not subject to fast self-heating or instruction-dependent fast voltage droops [5], hence requiring to track only slow variations such as ambient temperature variations, degradation of the power supply network, and aging. In this case, a software approach where general-purpose processors periodically query the on-chip monitors provides several advantages over designs with hardware feedback loops:

- it allows to use more accurate and flexible software models to track the best energy point under PVT variations [6], [9],
- it allows to exploit application knowledge to implement advanced power management policies, and
- it leads to a more generic not intrusive architecture, and causes less hardware overhead.

Body biasing (BB) is a well-known technique that consists in applying a voltage to the body contact of CMOS transistors to shift the effective transistor threshold voltage. It can be used either to boost speed [forward BB, (FBB)] or to reduce the leakage power consumption [reverse BB, (RBB)]. BB has been widely exploited in the past to compensate process and temperature induced variations [4], [7]. The principle has been widely applied in bulk technology [4], but its effectiveness decreases with CMOS scaling. Indeed, in traditional bulk technologies, body bias range is limited to ±300 mV due to gate induced drain leakage (RBB), and because of source–drain junction leakage and latch-up risk at higher voltage and temperature (FBB). In contrast, ultra thin body and box fully depleted silicon on insulator (UTBB FD-SOI), thanks to the buried oxide providing complete dielectric isolation between source and drain, allows to vary back-bias voltage over a wider voltage range. Past work on process and temperature compensation of UTBB FD-SOI processors focused on the high-performance flavor of the technology (flip-well) [7]. This technology flavor allows aggressive FBB (up to 3 V) leading to significant frequency boost, but RBB for leakage power minimization has very limited range (up to 300 mV). In this work, we leverage the conventional-well flavor of the technology [1], featuring much less leakage power and a wide-range

negative (up to –3 V) and positive [up to voltage drain drain (VDD)/2 + 300 mV] BB.

We propose a self-aware architecture for body biasing (BB) in a 28-nm UTBB FD-SOI technology to boost energy efficiency. We leverage the wide-range FBB and RBB to

· compensate for process and temperature variations, thereby reducing design-time margins and
· introduce a low-power mode and a state-retentive sleep mode [called power-nap, (PN)] with strong reverse RBB.

We designed an on-chip low-power body bias generator (BBGEN) to operate in a software-controlled closed loop with embedded Process Monitoring Boxes (PMBs) with minimal overhead during active, sleep, and PN modes, while preserving the capability to rapidly switch with low energy overhead between the different power modes. Leveraging the environmental- and application-awareness, the proposed approach enables an adaptive management of on-chip resources at both the circuit and architecture levels, opening the way for energy efficient and dependable processors targeting the end-nodes of the IoT.

## Low-cost self-aware architecture

The self-aware architecture is integrated in an NT multicore reference design shown in Figure 1. The system features three fully asynchronous power- and clock domains controlled by two frequency-locked loops (FLLs): the SoC domain, a small safe domain, and the cluster domain. The SoC and the cluster domains can be independently body-biased. The design was implemented with SoC encounter exploiting a multimode multicorner flow with Common Power Format for the description of the power domains. The default digital implementation flow for the cluster and the SoC requires design margins for worst-case process and temperature corners. Unfortunately, especially the low-voltage corners that reflect the strong impact of temperature inversion are highly pessimistic and impose significant overhead for timing closure. With postfabrication BB compensation, timing closure can be performed with almost zero-margins in the range from 20 MHz at 0.5 V to 200 MHz at 0.7 V for a typical die at 25 °C with native –0.4 V BB. As shown in Figure 2, the choice of a baseline –0.4 V BB, centered on the BB range of the BBGEN, allows applying strong FBB and RBB to operate at the best energy point across the target supply voltage/frequency range for a typical die at 25 °C.

The self-aware architecture leverages three hardware components that can be accessed by the processors through a memory mapped interface: the process monitors, the BBGENs, and the event unit. To distribute the body bias voltages to the power domains, the outputs of the BBGEN connect to a dedicated power grid within each body bias region. In contrast with traditional methodologies, the well-tap cells supporting triple-well disconnect the wells from the power rails, and separate body bias voltages can
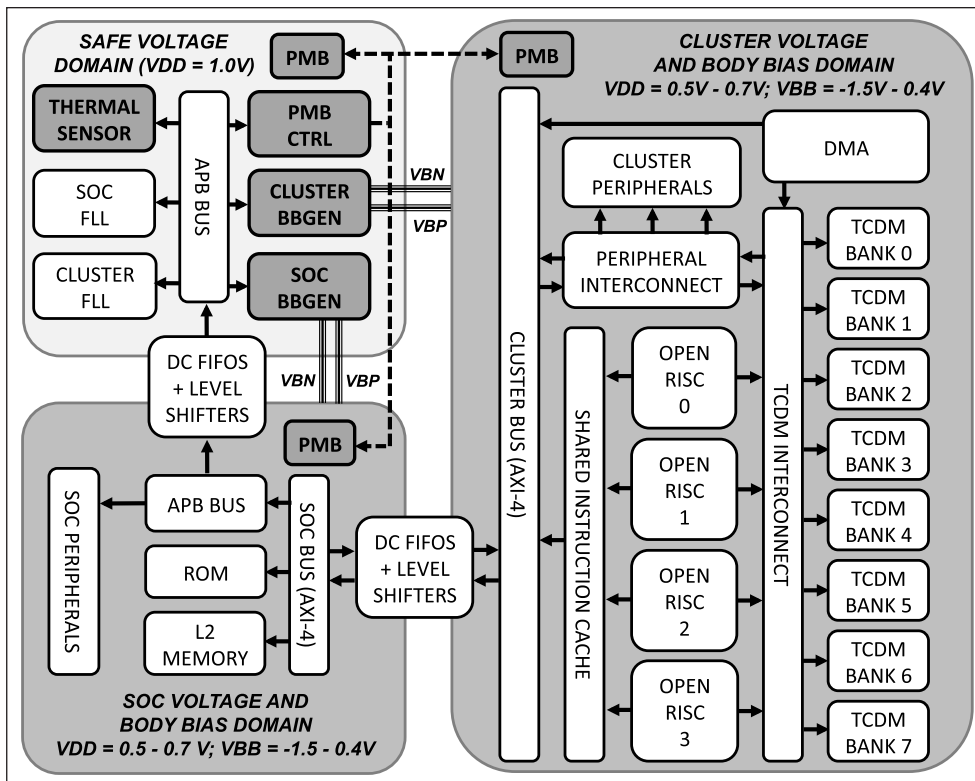


**Figure 1. SoC architecture and power domains.**

be supplied to the wells. Figure 3 shows a die micrograph and summarizes its main features. It is interesting to note the extremely low area impact of the components adopted to implement the self-aware architecture (less than 1% over just 3 mm² of silicon).

## PVT compensation

During active periods, the body bias voltage is constantly adjusted by a closed loop that comprises a temperature monitor and PMBs, a software-based controller, and a BBGEN.

The BB voltages are generated on-chip, with a generator that is similar to the design in [10] with several improvements (Figure 4). Conventional-well technology flavor [1] is employed to widen the BB range and lower leakage power for low-duty cycle operation. Consequently, a push–pull approach is used for positive regulation of both wells, while a dual phase charge pump is used to provide negative regulation on the p-well. A resistive digital to analog converter generates the reference for comparators, providing the feedback for the two control loops. To reduce power consumption to 4.5 µW, we employ a sleep mode that monitors only the impact of leakage of the wells.

As shown in Figure 5, the software controller periodically probes the status of the PMBs, consequently adjusting the BB value for the given target frequency. The PMBs include two sensors that can independently probe the process corner and static operating point of PMOS and NMOS transistors exploiting a ring oscillator (Figure 5a). Two register fields (*CLOCK DIVIDER* and *REFERENCE COUNTER*) are used to tune the sensor readout time according to the desired precision. Each sensor provides indirect information of the maximum operating frequency (*speedout*) computed with the following formula:

$$SPEEDOUT = \frac{CLOCK\ FREQUENCY * SENSOR\ STATUS * DIVIDER}{REFERENCE\ CONUTER}.$$

Then, the *speedout* has to be translated through the linear model shown in Figure 5. A margin of 5% on the predicted operating frequency, extracted from the characterization of the samples, has to be taken into
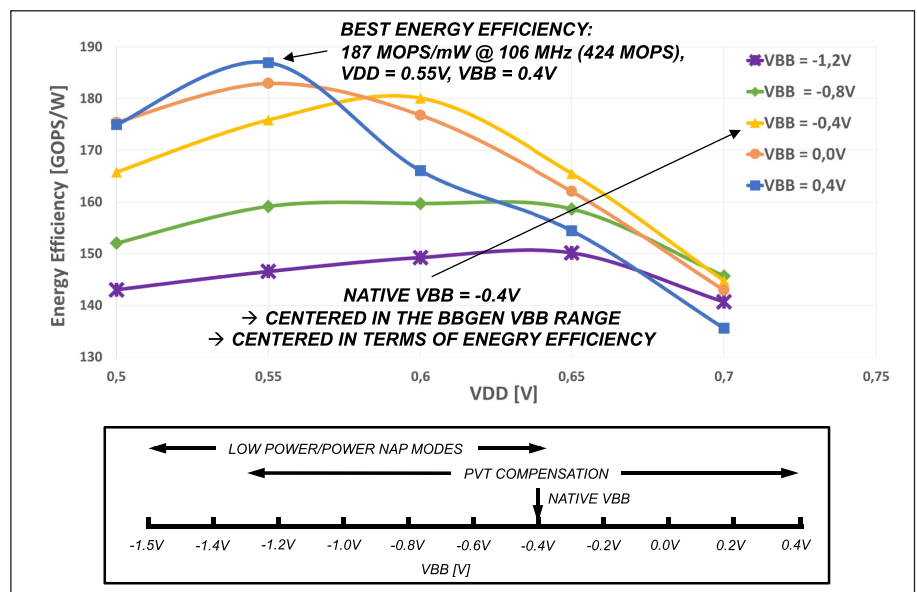


**Figure 2. Energy efficiency of the cluster vs. VDD at max frequency of a typical chip at 25 °C in the target voltage range (0.5–0.7 V): slight FBB provides best energy at 0.5 V, slight RBB provides best energy at 0.7V, and full RBB is useful for power nap mode. The native –0.4 V BB is well centered in both terms of energy efficiency and BB range of the BBGEN. Software control of body biasing (BB) allows tracking the best energy point from 0.5 to 0.7 V.**

account to deal with local variations that may occur at a smaller spatial granularity than the PMBs can probe. A control strategy based on a 64-entry lookup-table [4] translates the operating frequency into a BB value.

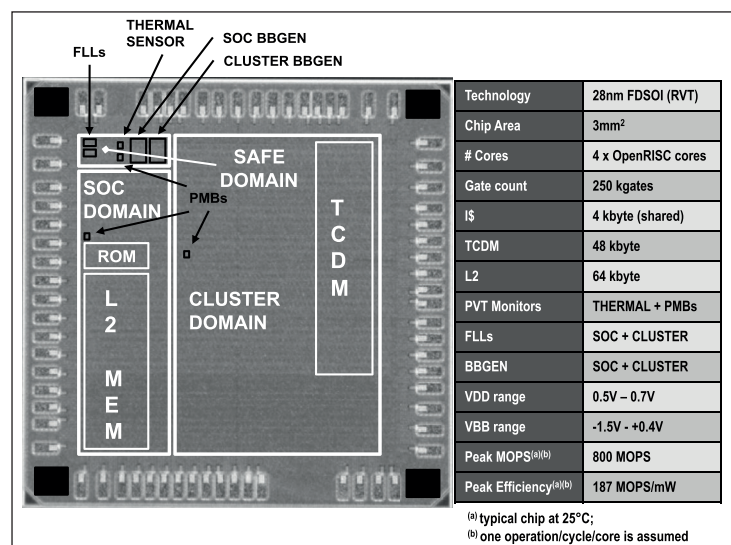The main potential drawback of this approach is the reduced speed and the energy cost implied



| Technology | 28nm FDSOI (RVT) |
|---|---|
| Chip Area | 3mm² |
| # Cores | 4 x OpenRISC cores |
| Gate count | 250 kgates |
| I$ | 4 kbyte (shared) |
| TCDM | 48 kbyte |
| L2 | 64 kbyte |
| PVT Monitors | THERMAL + PMBs |
| FLLs | SOC + CLUSTER |
| BBGEN | SOC + CLUSTER |
| VDD range | 0.5V – 0.7V |
| VBB range | -1.5V - +0.4V |
| Peak MOPS[a][b] | 800 MOPS |
| Peak Efficiency[a][b] | 187 MOPS/mW |

[a] typical chip at 25°C;
[b] one operation/cycle/core is assumed

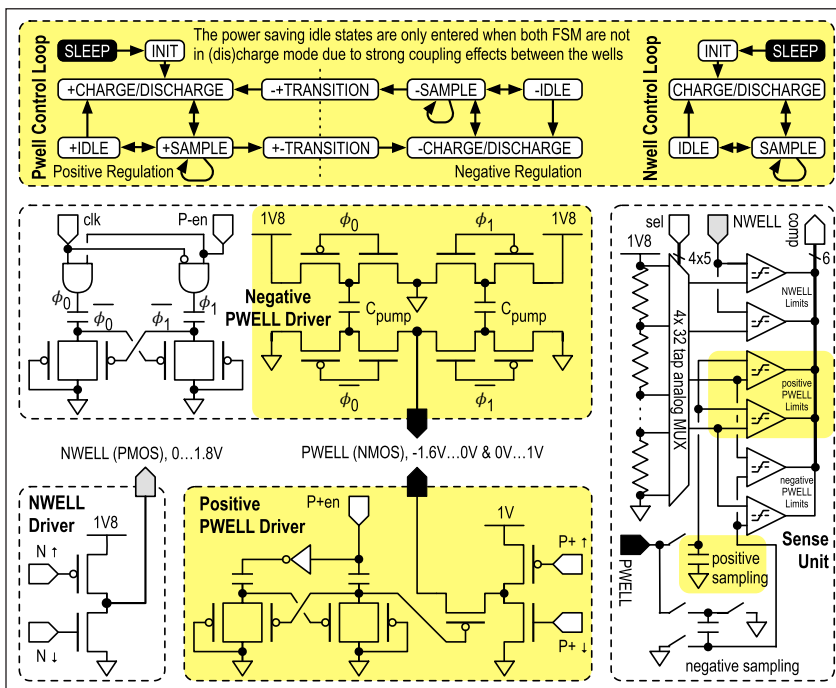**Figure 3. Chip micrograph and main features.**

**Figure 4. Schematic of the BBGEN highlighting differences to [10].**

by the software handler that reads the PMBs and closes the self-adaptation loop on the BBGEN, which requires ~100 clock cycles. This implies a latency of up to 5 μs (at 20 MHz and 0.5 V) for
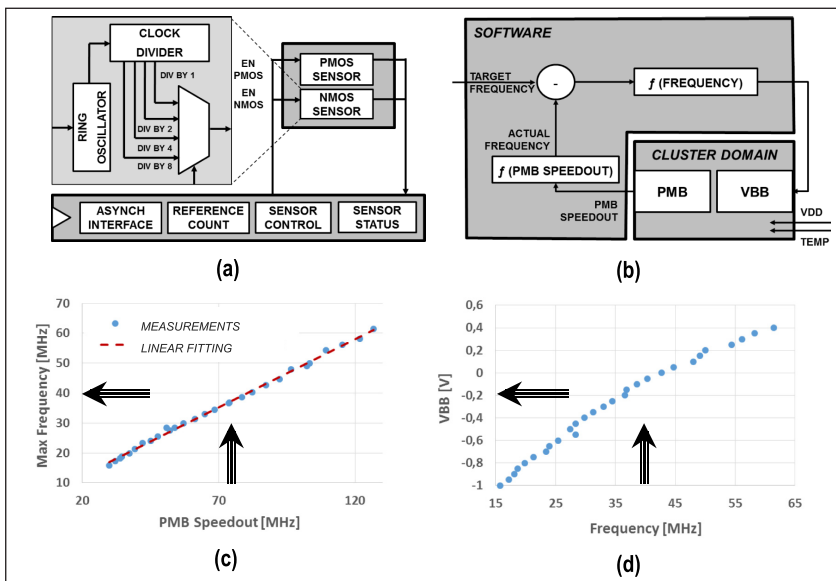


**Figure 5. (a) Block diagram of the PMBs, with zoom of the speedometer circuit architecture. (b) Block diagram of the process and temperature compensation loop. (c) Relation between PMBs *speedout* and cluster operating frequency. (d) Relation between the target operating frequency and body bias voltage.**

the software control plus ~7.5 ns/mV for BBGEN settling (in active mode), which depends on the applied settings and leads to an energy overhead of less than 2nJ with respect to a hardware solution. Considering the ambient temperature variations as the main target of this work, and assuming to probe the PMBs every tenth of a second, the proposed architecture has an overhead of the order of 0.01% both in terms of runtime and energy. Hence, it is clear that even more complex control approaches (e.g., model-predictive control [6], [9]) can be adopted, without significant impact on the performance and energy overheads required for the additional computations.

## Power nap

In addition to process compensation during active periods, we exploit the ability of the conventional-well FD-SOI process to implement a low-power low-frequency operation (a few megahertz, moderate RBB) mode and a PN mode (subkilohertz, strong RBB) for the SoC and cluster, leveraging the application knowledge enabled by the software programmability of the self-aware architecture. The low-power mode is mainly intended to maintain energy efficient operation during data transfers from peripherals to L2 memory, while the PN mode of the cluster saves significant leakage power for short, but frequent sleep periods for which the overhead of state-lossy power gating would be too high.

An event unit automatically manages transitions of the cores between the active and idle states. Processors go into the low-power and PN modes by software with a dedicated instruction after programming the desired frequency on the FLL and RBB mode on the BBGEN. After entering the PN mode, the cores remain idle (clock gated) until a configurable event is triggered by one of the peripherals. The transition time target from no BB to deep RBB and vice versa determines the

time required to enter and leave the modes but also the energy overhead for the transition, since higher operating frequencies are required in the BBGEN for its charge pump. We chose a relatively slow transition time of ~10 μs that can be accomplished already with 10-MHz clock, which is slow compared to other BBGENs (i.e., 1–2 GHz; see Table 1 [10]–[12]). Yet, the cluster can enter its PN mode in only 7.1 μs and leave in 11.5 μs, with an energy overhead of just 25 nJ. This architecture provides an energy advantage already for sleep periods beyond 203 μs where leakage is reduced from 145 to 22 μW, and a save and restore of the cluster state at 20 MHz would be too expensive.

**Table 1. Comparison with other BBGEN architectures.**

| | This Work | | [10] | | [11] | [12] |
|---|---|---|---|---|---|---|
| Technology | FD-SOI 28nm conv. well | | FD-SOI 28nm flip well | | 65nm Bulk | 90nm Bulk |
| BBGEN Area | 0.00913mm² | | 0.012mm² | | 0.0052mm² | 0.03mm² |
| Bias Area | 0.9mm² | | 1mm² | | 0.22mm² | 1mm² |
| BBGen Supply | 1.8V $V_{DDIO}$, 1V$V_{DD}$ | | 1.8V $V_{DDIO}$, 1V$V_{VDD}$ | | 0.5V-1V | 1.2V |
| Range N/P-WELL | 0V…1.8V | 0.4V…-1.5V | 0V…1.8V | 0V…-1.4V | VDD/GND±250mV | 1.2V…0.7V 0V…0.5V |
| $F_{sampling}$/$F_{Cpump}$ | 10MHz/100MHz | | 1-2GHz/1GHz | | - | - |
| Power | 4.15μW | | 10μW | | 600μW | 177μW |
| Standby $t_{sampling}$ | 1.31ms | | 2μs | | - | - |
| Max. Drive Current | 1.07mA@1.8V | | 40-200mA@1.8V | | - | - |
| $t_{Transition}$ N/P-WELL | 2.3μs | 11.5μs | 27ns | 160ns/90ns* | 2μs | 4μs |
| Transition Energy | <25nJ | | - | | - | - |

*with one / two drivers

## Chip measurements

Figure 6 shows frequency and power consumption measurements of three prototypes featuring different process corners [i.e. fast-fast (FF), slow-slow (SS), and typical-typical (TT)] selected among 60 samples, at the operating temperature of –20, 25, and 85 °C. The characterization was performed at the operating voltage of 0.5 V at 20 MHz. Due to the thermal inversion effect typical of NT devices implemented in deep submicron technologies, chips are faster at high temperature and slower at low temperature [1]. Variations across process corners and a temperature range between –20 and +85 °C can be fully compensated for a supply voltage ranging from 0.5 to 0.7 V, where FBB restores the target operating frequency (20 MHz at 0.5 V in Figure 6) of the chips when operating at low temperature, while RBB reduces leakage power consumption by up to 3.5× in the high-temperature and leakage-dominated regime, improving energy efficiency by up to 2×.

The impact on power of the low-margin design optimization and run-time power management techniques is shown in Figure 7a. During the active state, the relaxed implementation constraints coupled with BB compensation loop improve energy efficiency by 32% compared to a design that targets the same frequency with standard multicorner design. The power consumption of the cluster can be further reduced by 6.6× and 25× when entering the low-power and PN modes, respectively.

Figure 7b shows a comparison with a recent design implementing the PVT compensation loop with BB. As opposed to [7], where a hardware-based control loop with embedded CPM and a silicon-expensive
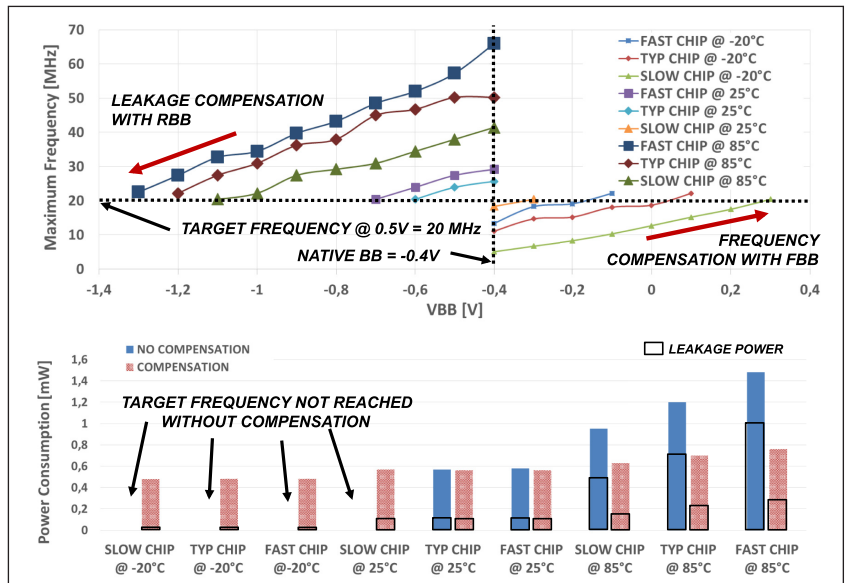


**Figure 6. Compensation of cluster's frequency and leakage power at the lowest operating voltage is 0.5 V, where the process and temperature variations are significant. Target frequency is 20 MHz.**

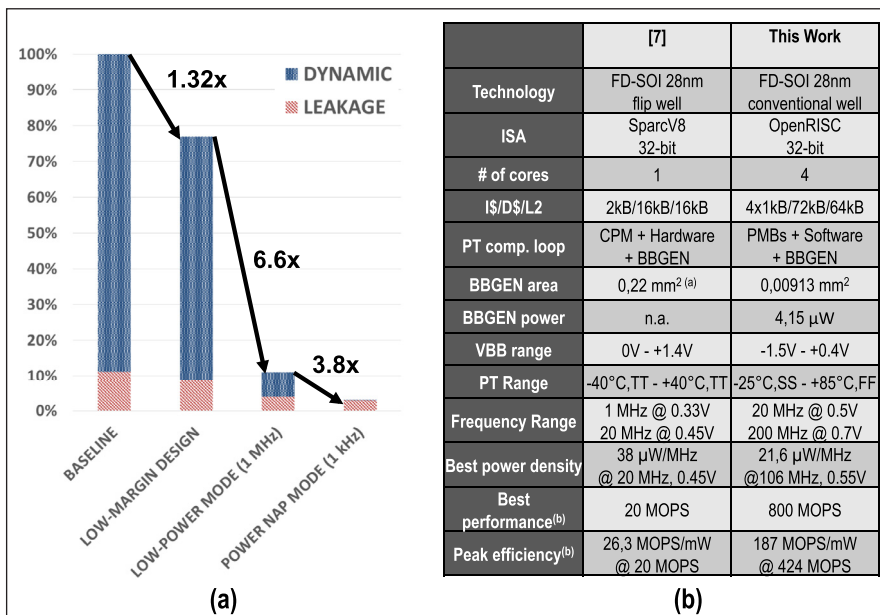| | [7] | This Work |
|---|---|---|
| Technology | FD-SOI 28nm flip well | FD-SOI 28nm conventional well |
| ISA | SparcV8 32-bit | OpenRISC 32-bit |
| # of cores | 1 | 4 |
| I$/D$/L2 | 2kB/16kB/16kB | 4x1kB/72kB/64kB |
| PT comp. loop | CPM + Hardware + BBGEN | PMBs + Software + BBGEN |
| BBGEN area | 0,22 mm$^{2\ (a)}$ | 0,00913 mm$^2$ |
| BBGEN power | n.a. | 4,15 µW |
| VBB range | 0V - +1.4V | -1.5V - +0.4V |
| PT Range | -40°C,TT - +40°C,TT | -25°C,SS - +85°C,FF |
| Frequency Range | 1 MHz @ 0.33V 20 MHz @ 0.45V | 20 MHz @ 0.5V 200 MHz @ 0.7V |
| Best power density | 38 µW/MHz @ 20 MHz, 0.45V | 21,6 µW/MHz @106 MHz, 0.55V |
| Best performance[b] | 20 MOPS | 800 MOPS |
| Peak efficiency[b] | 26,3 MOPS/mW @ 20 MOPS | 187 MOPS/mW @ 424 MOPS |

(a)               (b)

Figure 7. (a) Impact of design-time optimizations and power management modes on the operating power of a typical chip at 0.5 V and 25 °C. An operating frequency of 20 MHz is considered during active operation. (b) Comparison with other IoT architectures integrating PVT compensation loop with BB.

BBGEN are employed to track fast variations, this is not necessary in ULP designs [5]. The proposed software approach allows adopting pro-active techniques mixing hardware monitoring application knowledge [6], [9], while relaxing the design of the circuit components implementing the closed loop, significantly improving energy efficiency. With respect to [7], our architecture compensates for a wider PT range (SS, –25 °C to FF, 85 °C) with 24× smaller area overhead and only 4.15 µW of power consumption. Moreover, thanks to the low-margin design methodology enabled by the proposed software-controlled self-aware architecture, our design surpasses [7] by 21.2× in energy efficiency and 7.1× in performance.

**THIS WORK PRESENTS** a software-controllable self-aware architecture exploiting BB for compensation of PVT variations and for implementation of low-power modes in NT processors, implemented in 28-nm UTBB FD-SOI technology. The proposed architecture reduces design margins enormously, improving energy efficiency of the processor by 32% with a hardware cost of less than 1% and a runtime cost for software control of less than 0.01%. We demonstrate 24× area reduction for the compensation loop and 21.2× better efficiency

with respect to the previous design implementing the loop in hardware. Our future work will focus on the development of methodologies to automatically correlate and calibrate PMBs with the maximum frequency at different operating voltages and to test them on a large amount of samples. We also plan to explore advanced control strategies such as a model-predictive control with application awareness. ∎

## Acknowledgments

## ■ References

[1] D. Rossi et al., "A 60 GOPS/W,-1.8V to 0.9V body bias ULP cluster in 28nm UTBB FD-SOI technology," *Solid State Electron.*, vol. 117, pp. 170–184, 2016.

[2] R. G. Dreslinski, M. Wieckowski, D. Blaauw, D. Sylvester, and T. Mudge, "Near-threshold computing: Reclaiming Moore's law through energy efficient integrated circuits," in *Proc. IEEE*, vol. 98, no. 2, pp. 253–266, 2010.

[3] D. Blaauw et al., "Razor II: In situ error detection and correction for PVT and SER tolerance," in *Proc. Int. Solid-State Circuits Conf.*, 2008, pp. 400–622.

[4] J. Tschanz et al., "Adaptive frequency and biasing techniques for tolerance to dynamic temperature-voltage variations and aging," in *Proc. Int. Solid-State Circuits Conf.*, 2007, pp. 292–604.

[5] G. Paci et al., "Exploring "temperature-aware" design in low-power MPSoCs," in *Proc. Des. Autom. Test Europe*, 2006, p. 180.

[6] F. Ye et al., "On-chip droop-induced circuit delay prediction based on support-vector machines," *IEEE Trans. Computer-Aided Des. Integr. Circuits Syst.*, vol. 35, no. 4, pp. 665–678, 2016.

[7] S. Clerc et al., "A 0.33V/-40°C process/temperature closed-loop compensation SoC embedding all-digital

clock multiplier and DC-DC converter exploiting FDSOI 28nm back-gate biasing," in *Proc. Int. Solid-State Circuits Conf.*, 2015, pp. 1–3.

[8] M. Zandrahimi et al., "Challenges of using on-chip performance monitors for process and environmental variation compensation," in *Proc. Des. Autom. Test Europe*, 2016, pp. 1018–1019.

[9] S. V. Kumar, C. H. Kim, and S. S. Sapatnekar, "Body bias voltage computations for process and temperature compensation," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 16, no. 3, pp. 249–262, 2008.

[10] M. Blagojević et al., "A fast, flexible, positive and negative adaptive body-bias generator in 28nm FDSOI," in *Proc. IEEE Symp. VLSI Circuits*, 2016, pp. 60–61.

[11] N. Kamae et al., "A body bias generator with wide supply-range down to threshold voltage for within-die variability compensation," in *Proc. 2014 IEEE Asian-Solid-State Circuits Conf.*, 2014, pp. 53–56.

[12] M. Meijer et al., "A forward body bias generator for digital CMOS circuits with supply voltage scaling," in *Proc. IEEE Int. Symp. Circuits and Syst.*, 2010, pp. 2482–2485.

**Davide Rossi** is an Assistant Professor at the Energy Efficient Embedded Systems Laboratory at the University of Bologna. His current research interests include ultra-low power multicore SoC design and applications. He received a PhD in electronics engineering from the University of Bologna.

**Igor Loi** is an Assistant Professor at the Energy Efficient Embedded Systems Laboratory at the University of Bologna. His current research interests include ultra-low power multicore systems and memory systems evolution. He received a PhD in electronics engineering from the University of Bologna.

**Antonio Pullini** is a PhD student at the Swiss Federal Institute of Technology, Zurich. His current research interests include ultra-low power SoC design with a special focus on the peripheral subsystems. He received an MSc in electronics engineering from the University of Bologna.

**Christoph Müller** is pursuing a PhD degree with EPFL, Lausanne, Switzerland, with a research focus on digital implementation and an emphasis on low power, variation mitigation, and design methodology. He received a master's degree in system-on-chip from Lund University, Lund, Sweden, in 2013.

**Andreas Burg** is a Professor at the École Polytechnique Fédérale de Lausanne, where he leads the Telecommunications Circuits Laboratory. His current research interests include low-power very large scale integration signal processing. He received a PhD in electronics engineering from ETHZ.

**Francesco Conti** is a Post-Doc at the University of Bologna and the Swiss Federal Institute of Technology, Zurich. His current research interests include ultra-low power SoC design with a special focus on embedded artificial intelligence. He received a PhD in electronics engineering from the University of Bologna.

**Luca Benini** is a Full Professor of electronics at the University of Bologna and at the Swiss Federal Institute of Technology, Zurich. His current research interests include energy-efficient system design and multicore SoC design. He received a PhD degree in electrical engineering from Stanford University.

**Philippe Flatresse** is a Design Architect at STMicroelectronics Central R&D. His research interests include low-power and high-performance digital design techniques in both bulk and SOI technologies. He received a PhD in microelectronics from the Grenoble Institute of Technology.

■ Direct questions and comments about this article to Davide Rossi, Energy Efficient Embedded Systems Laboratory, 40123, Bologna, Italy; e-mail: davide.rossi@unibo.it.