

Alma Mater Studiorum Università di Bologna
Archivio istituzionale della ricerca

Performance-aware predictive-model-based on-chip body-bias regulation strategy for an ULP multi-core cluster in 28 nm UTBB FD-SOI

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

Di Mauro A., Rossi D., Pullini A., Flatresse P., Benini L. (2020). Performance-aware predictive-model-based on-chip body-bias regulation strategy for an ULP multi-core cluster in 28 nm UTBB FD-SOI. INTEGRATION, 72, 194-207 [10.1016/j.vlsi.2019.12.006].

Availability:

This version is available at: <https://hdl.handle.net/11585/766914> since: 2020-07-23

Published:

DOI: <http://doi.org/10.1016/j.vlsi.2019.12.006>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

Alfio Di Mauro, Davide Rossi, Antonio Pullini, Philippe Flatresse, Luca Benini (2020). Performance-aware predictive-model-based on-chip body-bias regulation strategy for an ULP multi-core cluster in 28 nm UTBB FD-SOI. In Integration, Volume 72, pag. 194-207. DOI: 10.1016/j.vlsi.2019.12.006

The final published version is available online at:
<https://doi.org/10.1016/j.vlsi.2019.12.006>

Rights / License:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website:

<https://www.elsevier.com/journals/integration/0167-9260/open-access-options>

<https://www.elsevier.com/about/policies/sharing>

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>)

When citing, please refer to the published version.

Performance-Aware Predictive-Model-Based On-Chip Body-Bias Regulation Strategy for an ULP Multi-Core Cluster in 28nm UTBB FD-SOI

Alfio Di Mauro[†], Davide Rossi^{*}, Antonio Pullini[†], Philippe Flatresse[‡], Luca Benini^{*†}

^{*}DEI, University of Bologna, Via Risorgimento 2, 40136 Bologna, Italy

[†]Integrated Systems Laboratory, ETHZ, Gloriastr. 35, 8092 Zurich, Switzerland

[‡]SOITEC, Crolles, France

Abstract—The performance and reliability of Ultra-Low-Power (ULP) computing platforms are adversely affected by environmental temperature and process variations. Mitigating the effect of these phenomena becomes crucial when these devices operate near-threshold, due to the magnification of process variations and to the strong temperature inversion effect that affects advanced technology nodes in low-voltage corners, which causes huge overhead due to margining for timing closure. Supporting an extended range of reverse and forward body-bias, UTBB FD-SOI technology provides a powerful knob to compensate for such variations. In this work we propose a methodology to maximize energy efficiency at run-time exploiting body biasing on a ULP platform operating near-threshold. The proposed method relies on on-line performance measurements by means of Process Monitoring Blocks (PMBs) coupled with an on-chip low-power body bias generator. We correlate the measurement performed by the PMBs to the maximum achievable frequency of the system, deriving a predictive model able to estimate it with an error of 9.7% at 0.7V. To minimize the effect of process variations we propose a calibration procedure that allows to use a PMB model affected by only the temperature-induced error, which reduces the frequency estimation error by 2.4x (from 9.7% to 4%). We finally propose a controller architecture relying on the derived models to automatically regulate at run-time the body bias voltage. We demonstrate that adjusting the body bias voltage against environmental temperature variations leads up to 2X reduction in the leakage power and a 15% improvement on the global energy consumption when the system operates at 0.7V and 170MHz.

I. INTRODUCTION

PERVASIVE use of embedded computing platforms in e-Health, wearables, smart sensors and Internet of Things (IoT) applications is pushing the research community to extensively explore the performance-energy tradeoff. IoT edge computing requires Ultra-Low-Power (ULP) devices consuming few mW and delivering GOPS. On the other hand, Moore's law is slowing down, and CMOS scaling does not lead huge energy efficiency benefits any longer [1]. This pushes for ever more aggressive voltage reductions, i.e. operating transistors very close to their threshold voltage, an approach known as Near-Threshold Computing (NTC) [2] [3]. The major challenge that near-threshold IoT devices have to face is operating in many different scenarios that are not always completely predictable at design time. More specifically a serious concern is validation on very wide range of operating temperatures.

This aspect is crucial, as devices implemented in most advanced technological nodes have a strong dependency between environmental temperature, operating frequency and leakage power [4] [5]. This behaviour is caused by a phenomenon called Temperature Effect Inversion (TEI) [6]. In deep sub-micron technologies the detrimental effect that the temperature has on the maximum frequency of a device is reverted [7]. Due to TEI [8], the performance and leakage current have a positive sensitivity to temperature, especially for those devices operating in low-voltage corners. Another effect which must be taken into account in ULP devices implemented in recent technological nodes is process variation¹. Accounting for temperature and process variation requires to take very large margins at design time, which causes huge overheads in power due to the needs of buffers for setup and hold time fixing in different corners, often making it impossible to achieve timing closure [9].

To achieve correct operation and acceptable power-performance ranges in near-threshold CMOS chips, post-fabrication compensation of process and temperature variations is critically needed. Body biasing (BB) is a well-known approach for post-fabrication compensation. It consists in polarizing with a voltage potential the N and P well of CMOS transistors, thereby changing threshold voltage. The performance of a device, in terms of maximum achievable frequency, is directly related with the threshold voltage of transistors, more specifically, when the threshold voltage decreases the maximum frequency of the device increases [10] [11]. On the other side, the leakage current increases when the threshold voltage decreases, thus increasing static power consumption. Body biasing represents an alternative to adaptation of supply voltage, the latter is less efficient and requires more complex support circuitry, like on chip DC/DC converters or voltage regulators resulting in higher power overheads and coarser granularity [12].

The use of body biasing is twofold: *i*) Forward Body Biasing (FBB) allows to increase the operating frequency of the chip with a minimum power overhead *ii*) Reverse Body Biasing (RBB) could significantly reduce the leakage current when both process and environmental temperature would allow the

¹We neglect the effect of aging since self-heating is minimum in ULP devices

device to run faster than necessary, given a target application [13].

The aim of this work is to propose a performance-aware body biasing strategy to compensate for process and temperature variations while guaranteeing with minimum power overhead a user-specified target frequency. The methodology exploits a calibration procedure which can be executed at boot time, to tune all the involved software models and compensate static performance variations (e.g. process and aging effects). Then, dynamic variations (e.g. temperature) are compensated at run-time by exploiting on-chip frequency measurement obtained from Process Monitoring Boxes (PMB) properly calibrated to minimize the frequency measurement error. Energy efficiency improvement is obtained by exploiting the capabilities of UTBB FD-SOI (Ultra Thin Body and Box Fully Depleted Silicon On Insulator) technology to apply a wide-range forward and reverse body bias voltage. Thanks to on-chip ULP Body Bias Generators (BBG) we modulate the V_{bb} from -1 V corresponding to full RBB to $V_{dd}/2 + 300\text{ mV}$ corresponding to full FBB. The proposed methodology is suitable to be implemented as a software control strategy on any processor featuring both on-chip PMBs and BB generators. The software overhead can be considered negligible, since the controller can be activated with periods in the order of seconds to track temperature variations, and the Body bias regulation requires approximately Rossi2017. Additionally, exploiting flexibility of software models, higher accuracy can be achieved with respect to pure HW-based controllers [14].

As a preliminary step we operated a complete characterization of the on-chip frequency measurement modules (PMB) available to us. We performed this operation on an advanced chip testing equipment: Advantest SoCV93000 tester system. In this phase we studied the accuracy of the PMBs versus the environmental temperature and we derived a mathematical model correlating their frequency measurement with the actual maximum frequency of the chip. To characterize the PMB in the widest possible temperature and voltage operating range we performed our measurements at three different temperature ($T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$) and supply voltages ($V_{dd} = \{0.5\text{ V}, 0.7\text{ V}, 0.9\text{ V}\}$). At every operating point identified by V_{DD} and T we computed the correlation between the PMB frequency measurement and the real maximum frequency of the device on the entire body bias range (-1 V to $V_{dd}/2 + 300\text{ mV}$). We repeated this operation on chips belonging to different process corners.

Note that the scope of the PMB characterization is to establish a relationship between the output of the PMBs and the maximum frequency of the chip. The maximum frequency could deviate from the trend predicted by the PMBs at different supply and body bias voltages. The PMBs itself cannot provide an exhaustive representation of all the reasons that limit the frequency of the device (e.g. critical path going to the memories). However, every mismatch between the output of the PMBs and actual maximum frequency is embedded in the software model constructed during the PMB characterization.

In a successive phase we developed a strategy exploiting both the PMB and BB models to dynamically modulate the

body bias voltage and achieve the desired target frequency against temperature variations. This task is executed by a software body bias controller which probes the performance of the chip and as a consequence programs the body bias generators to apply the required V_{BB} . To compensate the model inaccuracy we adopt forward body bias margins that are added to the V_{BB} regulation operated by the controller. These forward body bias margins allow to prevent chip failures when the PMBs are over-estimating the chip performance and the V_{BB} controller is applying less forward body bias than necessary. The model error reduction comes with a cost, adding body bias margins to the V_{BB} regulation increases the leakage current. We demonstrate that with a calibration procedure it is possible to reduce the frequency estimation error by a factor of 2.4X, delivering accurate predictions on chip performance with a minimum leakage cost. We propose a controller architecture based on such models, dynamically adjusting the V_{BB} either to achieve the target frequency or to reduce the leakage current as much as possible. Although the system features two different power domains with body biasing capabilities, to simplify the implementation, the body bias regulation proposed in this paper has been tested and validated on the most computationally capable and power-hungry power domain, that is the cluster. In this scenario only the body bias generator related to the cluster domain is controlled by the regulation loop; the second body bias generator is set to generate a constant $V_{BB} = 0\text{ V}$. In this work we demonstrate that the proposed body biasing controlling methodology can achieve up to 2X in leakage reduction and 15% in energy efficiency improvement at 170MHz and $V_{DD} = 0.7\text{ V}$ on the related power domain.

The remainder of this work is organized as follows. Section II presents an overview of the different methodologies to monitor on-chip the performance and the strategies to dynamically modulate the body bias voltage already presented in the past works. Section III describes the chip used as test vehicle and the most relevant IPs enabling the adaptive body biasing methodology. In section IV we report all the required steps to derive a model for the PMBs and how to reduce its error. Section V describes a PMB calibration procedure and the body bias controller; additionally, it also shows two working examples of the dynamic body biasing controller. Section VI provides the results in terms of leakage current reduction and energy efficiency improvement when the controller is active. Finally, section VII contains concluding remarks.

II. RELATED WORK

With the diffusion of Ultra-Low-Power (ULP) devices operating in low voltage corners, robustness to variations introduced by process and temperature has become a critical issue. The problem of variability in this field has been analyzed in the past by several works. The solutions proposed to address this challenge leverage design-time and run-time techniques employed at the circuit, micro-architectural and architectural levels. These techniques lead to mitigate, improve resiliency, or optimize energy or performance in presence of variations. M. Alioto [9] highlighted the inefficiency to apply conventional design paradigms to ULP-oriented devices, analyzing

the impact of process and temperature variation on low power designs operating in the near- and sub-threshold region, and providing guidelines for design of standard cells, memories and microarchitectures. Other circuitual and architectural solutions to reduce the impact of variability are presented by M. Seok et. al. [15], their work mostly focuses on design-time techniques addressing the problem of variations in logic, memories and clock tree.

An orthogonal approach to address the problem of variation is that of compensation. While the aforementioned methodologies rely on decisions that have to be taken at design time (e.g. very conservative supply voltage margins), run-time variation compensation allow to reduce margining, leaving rooms for energy efficiency improvements. Most common approaches to compensate variations at run-time exploit Dynamic Voltage Scaling (DVS) or Dynamic Body Biasing (DBB) as compensation knobs. In many architectures V_{DD} and V_{BB} are controlled in closed loops that rely on circuitual parameters tracking (e.g. frequency or temperature). Therefore, probing the system status in terms of maximum achievable frequency and temperature is needed to properly regulate the compensation knob.

A. Probing approaches

A powerful approach to estimate system maximum frequency in the context of DVS methodologies is "Razor" [16] [17] [18]. It exploits in-situ error detection on the processor pipeline stages, the functionality of the registers is augmented with "Razor" shadow-latches capable to detect a timing failures and recover from it. D. Bull et. al. [19] applied Razor to a 32bits microprocessor implemented in UMC 65nm process, reporting 3% of frequency estimation error. Fojtik et al. propose an improved Razor-based solution called "Bubble Razor" [20], targeting designs operating in low-voltage corners and exploiting flip-flops based datapath to two-phase latches datapath conversion; this allows the use of two-phases clocking methodology, enabling larger timing speculation windows with respect to [16] [17] [18]. The advantage in using these strategies is that they allow to significantly reduce design margins, mitigating the effect of both global and local process variations when used in combination with compensation approaches. However, if on one hand they enable promising energy savings, on the other hand their application is very intrusive, and it is limited to all those cases where deep knowledge of the architecture is possible, and modification to the RTL are allowed. Very often IPs composing a System on Chip (SoC) are provided as encrypted macros, or simply come as hardened macros from different design teams, hence no access to RTL to modify pipeline stages is possible.

Other approaches for maximum frequency on-chip probing are based on Critical Path Monitor (CPM). CPM are critical path replicas to which extra delay elements are added to make the path super-critical. Drake et al. [21] propose a CPM device to track intra-pipeline stage critical path (CP) where 5 different critical path replicas featuring different types of gate (e.g. NAND4, NOR3) are synthesized to emulate different sensitivities of datapath logic gates to supply voltage. Tschanz

et al. [22] propose a more general approach. Instead of pre-determined logic gate types, configurable buffer delay chains are tuned to emulate largest delay path between pipeline stages, providing a feedback on the circuit maximum frequency. A similar solution is proposed by Clerc et al. [12], to take into account process variability, different CPMs types can be tuned to match the critical path of the device at a given voltage and temperature. Due to the low level of intrusiveness, and since ULP devices are only marginally affected by intra-chip variations (devices in few mm^2 range) CPM-based approaches can be very effective to estimate the maximum frequency. However, when critical path involves RAM, CP can be rarely emulated by cascading logic gates or with delay buffer chains because of the mixed-signal nature of internal signals, making these methodologies hardly usable in complex designs. Clearly identify a critical path is a difficult task in devices fabricated in deep-sub micrometer technology nodes and operating in low-voltage corners. To improve the accuracy, Beigné et al. [23] propose a more general approach combining the use of Critical Path Replica (CPR) and timing fault detection, coupling two different performance estimation methodologies. Such frequency tracking approach is implemented in 28nm FDSOI technology node and covers multiple voltage operating points (0.4V to 1.3V).

As Shown by Zandrahimi et al. [24], a big improvement in generality of probing approaches is introduced by the adoption of Process Monitoring Blocks coupled with a software model (This solution is adopted in our work). Thanks to hardware sensors based on several types of ring oscillators (e.g. NMOS only and PMOS only), frequency estimation can be performed with an accuracy of 7.6%. As we demonstrate in this work, by properly calibrating PMB models this error can be further reduced to 4%. Indeed, models implemented in software can be more flexible than LUT implemented in hardware and better fit the behaviour of PMBs. Moreover, adopting a PMB-based solution allows also to overcome limitations coming from architecture or CP awareness. Finally, area overhead of this probing methodology is very limited, as visually evident from the example of Fig. 1.

B. Actuation approaches

Actuation can be seen as the next step of probing, once the performance has been estimated, some action is performed to align the current behaviour of the system with the desired one. Some approaches dynamically change the operating frequency depending on the specific process/temperature/aging conditions. The solution proposed by Constantin et al. in [25] tackles the problem of margins reduction from a different perspective with respect to the adoption of specific circuitual solutions [9] [15]. Timing constraints are relaxed at design-time, simplifying timing closure. Clock frequency is then modulated at run-time according to the specific critical paths triggered by the executed instructions. Even though this methodology can allow significant energy improvements, its application requires to know both critical paths and delay paths triggered by a certain processor instruction. Moreover, the application of this methodology cannot ensure a target frequency matching.

Compensation of temperature and process variation has traditionally been performed by supply voltage adaptation. Various strategies exploiting Dynamic Voltage Scaling (DVS) in bulk technologies are presented in [16] [18] [19] [20] [26], where design-time margins are reduced thanks to the in-situ timing violations detection strategies previously mentioned. Specifically, the maximum performance of the chip is probed and used as feedback to modulate the supply voltage right before the Point-of-First-Failure (PoFF). While DVS can be very effective to reduce design-time margins, compensating run-time performance degradation caused by process and temperature variations, as demonstrated by [12], is a task that can be executed much more efficiently leveraging body biasing.

M. Miyazaki et.al. [27] demonstrated the effectiveness of process compensation by exploiting -1.5 RBB to + 0.5 FBB body bias range in a 200 nm CMOS technology. Similarly, Tschanz et. al. [28] implemented in 150nm CMOS technology node an adaptive body biasing scheme (-0.5V to + 0.5V) for process compensation. More recently, DBB has been also used to compensate not only for process and aging variation but also to compensate short-term performance variation caused by temperature changes. Kumar et.al. [29] developed an algorithm for temperature compensation in the range 35°C to 65°C, exploiting body biasing in devices fabricated in 65nm and 45nm technology nodes. Similarly, Tschanz et. al. [30] used body biasing for temperature compensation of a TCP/IP processor in 90nm technology operating at 1V in a range from 60°C to 80°C. More recently, Kang et al. [31] propose a PLL-based performance feedback circuit solution exploiting body biasing, implemented in IBM 130nm to compensate process, temperature and aging induced variations. In successive works, Kumar et al. [32] and Ono et al. [33] presented two hybrid methodologies exploiting adaptive forward body biasing combined with DVS to maintain optimal performance of a device against variations introduced by aging. More recently, Gammie et. al. [11] exploited a -0.5V to 0.5V body biasing for process variation compensation and for high performance and low-power states enabling.

All the listed approaches refer to performance compensation for devices implemented in bulk technologies. Unfortunately, such technologies only achieve good results either in limited temperature ranges or leveraging too small body bias ranges. In deep-submicron bulk technologies, maximum body bias range is limited by p-n junction leakage and potential latch up. In FinFET technologies [13], instead, the lack of an easy way to access the back gates represents the main obstacle to body bias voltage application [13]. In near-threshold, the impact of temperature variations is huge and cannot be fully compensated with the limited body biasing capabilities provided by the bulk CMOS technologies. On the contrary, 28nm UTBB-FDSOI provides a very powerful knob for process and temperature variations compensation because of its wide-range body biasing capability (theoretically from -3V to 3V). A good demonstration of the body biasing capabilities in FDSOI technology is provided by Clerck et al. [12], however the proposed design is fabricated with LVT cells, allowing only forward body biasing of the transistors in the range 0 V to 3 V. Contrarily, as anticipated by the analysis reported by Rossi

et al. in [14], and as implemented in this work, a processor implemented in UTBB FDSOI with RVT can exploit wide-range body biasing to implement advanced power management and compensation strategies.

In this work we fully characterize the PMB sensors, studying the correlation with the maximum frequency of the device in presence of process variations (i.e. among multiple chips belonging to different process corners), and demonstrating that with PMB sensors and properly calibrated software model it is possible to obtain a performance estimation accuracy comparable with the probing approaches previously described (4% at 0.7V, including temperature variations), with a negligible overhead on the device area, and with an extreme low level of intrusiveness. We then present a body bias controller capable to boost the performance through FBB and enter low leakage full-state-retention modes through RBB. Thanks to a finely calibrated control strategy implemented in software, we demonstrate how it is possible to compensate dynamic performance degradation and reduce the leakage current on an a real embedded platform, exploiting the higher efficiency of compensating temperature and process variations with body biasing [14] [12].

III. PULP SYSTEM

A. UTBB FD-SOI technology

PULPv3 has been implemented in 28nm Ultra-Thin Body and Box Fully Depleted SOI technology (UTBB FD-SOI) from STMicroelectronics. This technology features an improved channel electrostatic control thanks to thin-film technology, reducing leakage currents and Short Channel Effects (SCE). On this technology, for the same leakage current target the threshold voltage can be strongly scaled, thanks to the ultra-thin buried oxide. This ensures low variability when circuits operate close to the threshold voltage of transistors. Ultra-thin buried oxide also enables the use of very wide body biasing range from -3V to 3V using conventional and flip well transistor [13] [34]. In UTBB FD-SOI technology channel length modulation (i.e. poly biasing) is used at design time to statically optimize circuit critical path. In the case of PULPv3 implementation, we used the conventional-well flavor of the technology [13] and the entire system is implemented with the same type of Poly-Biasing 0 (PB0) Regular Voltage Threshold (RVT) cells.

B. Architecture

Parallel Ultra-Low-Power platform [35] version 3 (PULPv3) [14] is a multi-core SoC for ULP applications operating in near-threshold to achieve extreme energy efficiency on a wide range of operating points. The SoC is built around a cluster featuring four cores and 64kbyte of L2 memory. The cores are based on a highly power optimized micro-architecture implementing the OpenRISC-32bit ISA featuring 4kB of shared instruction cache. The cores do not have private data caches, avoiding memory coherency overhead and increasing area efficiency, while they share a L1 multi-banked Tightly Coupled Data Memory (TCDM) acting as a shared data scratchpad memory. The TCDM features 8 4kB

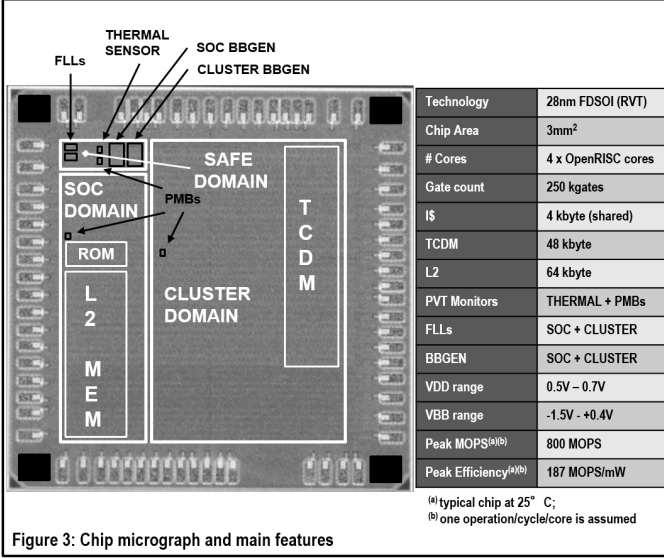


Figure 3: Chip micrograph and main features

Fig. 1: This figure shows a micrograph of the chip where it is possible to identify: the SoC domain, the cores cluster and the safe domain. Additionally, it is possible to note the PMB sensors in every power domain and PMB controller. Table on the right shows the main architectural features of the PULPv3 chip.

SRAM banks and 8 1kB Standard Cell Memory (SCM) banks connected to the processors through a single clock latency non-blocking interconnect, implementing a word-level interleaved scheme to minimize banking conflict probability. Off-cluster (L2) memory latency is managed by a tightly coupled DMA featuring private per-core programming channels, ultra-low programming latency and lightweight architecture optimized for low-power and high transfer efficiency. Fig.2 shows the architecture of the system while Fig. 1 shows a micrograph of the chip.

Three isolated power domains enable advanced power management: *i)* The "Safe Voltage Domain" hosting the Frequency Locked Loop generators, the two Body-Bias Generators for the SoC and Cluster regions, the PMB Controller and additional infrastructural control logic *ii)* The "SoC Body-Bias Domain" *iii)* The "Cluster Body-Bias Domain". Each domain is monitored by a PMB, which will be described in the next section. In our tests, we will focus only on the Cluster Domain, applying $V_{bb} = 0V$ to the other body-bias domains.

1) *Process Monitor Boxes*: A Process Monitor Box (PMB) is as an on-chip sensor based on ring oscillators, connected to the system through a memory mapped interface. It provides a measurement of the maximum frequency achievable by the device. These sensors are designed and optimized to behave consistently with the other library logic gates, emulating also performance variations induced by temperature and process variations. Since PULPv3 has been implemented with PB0 cells, the ring oscillators of the PMB sensors instantiated in each isolated power domain are implemented with the same type of PB0 cells.

2) *Body-Bias generators*: The body bias voltage is modulated thanks to a fully integrated body-bias generator. Such body bias generator is capable to supply an area of 1 mm². It applies a fine-grain body-bias voltage with a minimum step

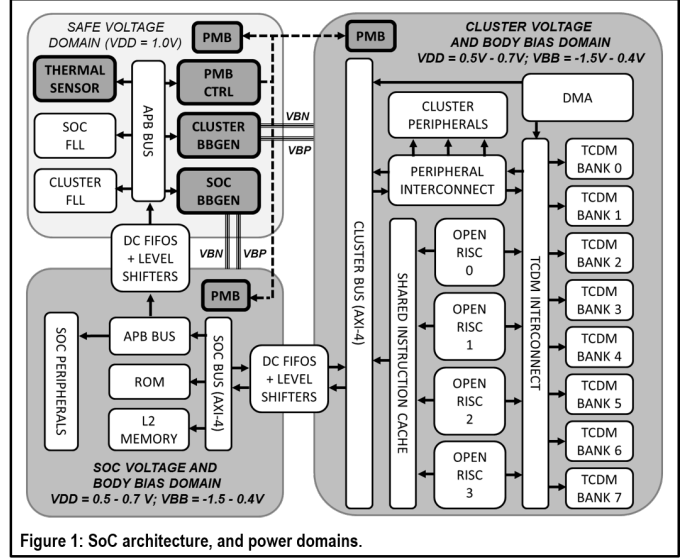


Figure 1: SoC architecture, and power domains.

Fig. 2: Architecture of PULPv3 system. The picture shows the three different power domains: the SoC domain, the cores cluster and the safe domain

of 50 mV and can cover a voltage range from $-1.5V$ to $V_{DD}/2 + 300mV$ in a $4.15\mu W$ power budget. This device has been designed to target ULP system on chips [36] and further optimized to maximize the energy efficiency [14]. For positive regulation of both wells it uses a push-pull approach. Negative regulation on the p-well is obtained by means of a dual phase charge pump. Two control loops allow to monitor the impact of the leakage and compensate for wells discharge. The comparators of the feedback use resistive Digital to Analog Converters (DAC) as references. To reduce the power consumption of the BB gen to $4.5\mu W$, its operation is duty-cycled, entering a sleep mode where only the leakage is monitored. The generator is programmed through a series of memory mapped registers and can independently bias with different voltages the Nwell and the Pwell of the transistor. Table I reports the main features of the body bias generator.

BBGEN Area	0.00913 mm ²
BBGEN Supply	1.8V V_{DDIO} , 1V V_{DD}
Power	4.15 μW
Transition time N-WELL	23 μs
Transition time P-WELL	11.5 μs
Transition energy	$\leq 25nJ$

TABLE I: Main features of the body bias generator.

C. PULPv3 Embedded Test System

Since the final goal of this work is to develop a fully automated body-bias control strategy, the analysis of the PMBs has been performed on the chip testing equipment while the proposed methodology has been validated on the PULPv3 evaluation board; this is an embedded platform that can be used to develop and test software applications like commercial microcontroller evaluation boards. The body bias voltage and the leakage current have been measured from the PULPv3 board jumpers by means of a power analyzer. The external

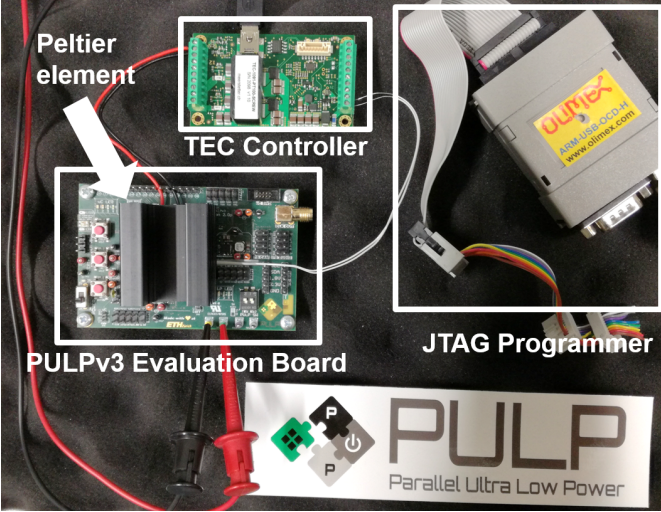


Fig. 3: PULPv3 body-bias controller testing setup. In the picture it can be observed the PULPv3 evaluation board hosting the PULPv3 chip and the JTAG programmer. Picture also shows the Peltier's element cell and the controller used to enforce a known temperature on the chip package

temperature variation have been simulated by enforcing a temperature on the chip package by means of a Peltier's element cell regulated by an industrial TEC controller. The Peltier's cell can modulate the chip package temperature in a range 10°C to 80°C ; the complete testing setup is reported in Fig. 3.

IV. MODEL

A. Experimental Setup

The PMB models derivation, as well as the body bias model described in the following sections are performed with an Advantest SoCV93000 tester system, in connection with a Thermonics 2500E temperature forcing system, able to force an environment temperature ranging from -80°C to 220°C . Since we are interested in using the body-biasing voltage as independent variable both for temperature and process variations compensation, we structured the measurements as follows: *i)* We defined the operating point (OP) in terms of {Supply voltage V_{dd} , Temperature}. More specifically, the voltage corners are: $V_{dd} = \{0.5\text{ V}, 0.7\text{ V}, 0.9\text{ V}\}$. Temperature corners are: $T = \{-20^{\circ}\text{C}, 25^{\circ}\text{C}, 80^{\circ}\text{C}\}$; *ii)* For each OP we swept the body-biasing voltage (V_{bb}) in the range -1 V to $V_{dd}/2 + 300\text{ mV}$ using the body-bias generator. At every operating point we measured: the leakage power (P_{LKG}), the active dynamic power (P_{DYN}), the total power (P_{TOT}) and maximum frequency (F_{MAX}) achievable by the device. We measured the active dynamic power (P_{DYN}) as the difference between the total power (P_{TOT}) and the leakage power (P_{LKG}). The total power (P_{TOT}) has been measured as the power consumption when the device is executing a benchmark application and the leakage power (P_{LKG}) as the power consumption when the device is not clocked.

We extracted the maximum operating frequency (F_{MAX}) by means of a carefully crafted benchmark (i.e. a sequence of arithmetic operations and memory stores) tested in post-

layout simulation, able to trigger the most critical paths² of the circuit. On the chip, we verified that the result of the benchmark was returned with the correct timing, and the End-Of-Computation³ signal was properly asserted. As cross-check, we verified that the result of the benchmark returned a valid check-sum.

B. PMB characterization

First part of the analysis focuses on the characterization of the performance monitoring blocks placed in the power domain of interest, that is the one hosting the core cluster. The aim of this step is to obtain a model allowing to correlate the response of the frequency probes (PMBs) to the maximum frequency of the device. To obtain the model, we simultaneously measured both the response of the PMBs (F_{PMB}) and the maximum chip frequency (F_{MAX}) versus the full body bias range. This operation has been repeated per each operating point (V_{DD} , T). The supply voltages reported in this section have been generated by the Advantest hp93000 testing equipment while the temperatures have been externally enforced on the chip package by the Thermonics 2500E temperature forcing system, body bias voltages are generated on-chip.

PULP V_{DD}	V_{BB}	T_{op}	Process corner
0.7 V	-1 V to 650 mV	25°C	Typical

TABLE II: Measurement operating conditions under which the PMB analysis has been performed

The aim of this step is to obtain an evaluation of the PMB accuracy. To perform the measurement in a well-defined environment, the temperature is externally enforced and kept constant by the temperature forcing system. The measurements that we performed in this phase are referred to a single chip, this allows to exclude "a priori" the process variations as possible source of error. Since the entire set of PULPv3 chips have been pre-characterized in terms of maximum performance, we know if the device in exam belongs to a *Fast*, *Typical* or *Slow* process corner. Fig. 4 reports the measurement of the F_{PMB} and F_{MAX} versus V_{BB} for a typical device supplied at 0.7 V , TABLE II reports more details regarding the measurement operating conditions.

In the second phase of the PMB characterization we evaluated the correlation between F_{PMB} and F_{MAX} . As shown in Fig. 5, the model we found is a linear function, and it is reported by (1).

$$F_{MAX} = C_{corr} F_{PMB} + F_0 \quad (1)$$

We repeated the same measurements in different operating corners proving the robustness of the model: TABLE III reports the parameters of the model.

²The critical path has been identified during the post layout timing analysis in the communication between the cores and the scm memory

³Physical output pin of the device which certifies that the system completed all the operations and properly entered a known final state.

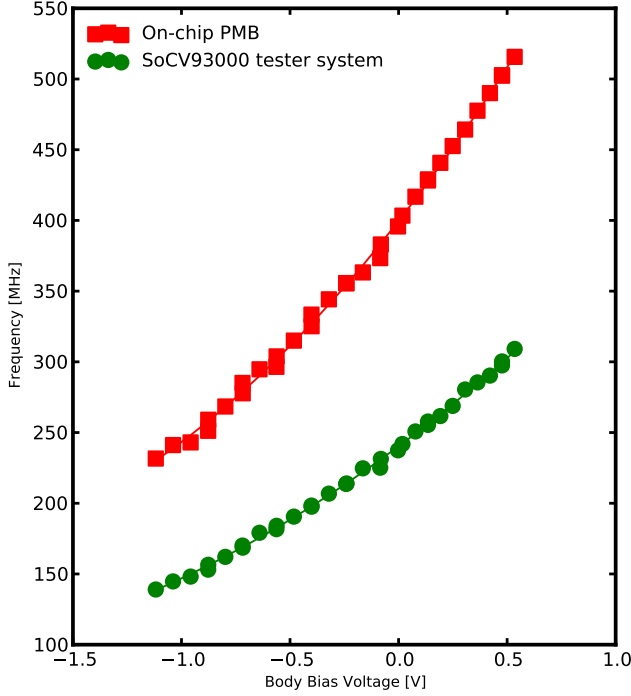


Fig. 4: This plot compares the maximum achievable frequency and the PMB frequency estimation for a "Typical" PULPv3 chip at 0.7 V and 25 °C.

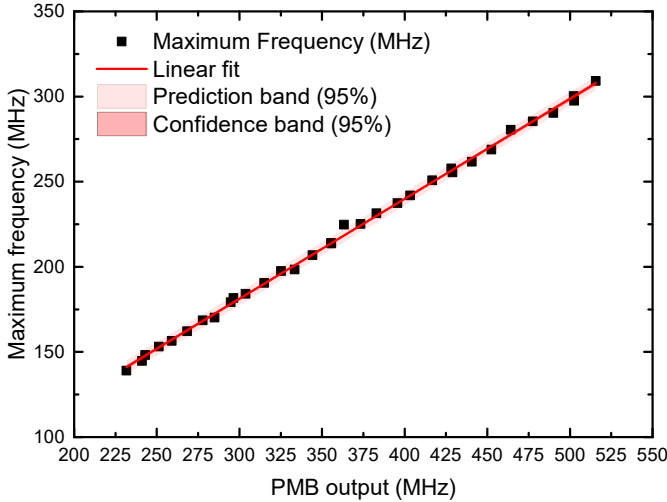


Fig. 5: This plot shows the correlation between the maximum frequency and the PMB frequency estimation for a typical chip at 0.7 V and 25 °C. The data have been fitted with a linear model, the parameters of the fit are: $y=mx+q$ where $m=0.59$ and $q=5.19$, R-Square (COD) = 0.998.

As a final step for this preliminary PMB characterization we estimated the error of the model, which is represented by the residuals of the measurements with respect to the fitting curve. Table VII shows the maximum errors we report for the correlation model (*Process-aware Temperature-aware*), for a single device. Comparable errors are reported also in [37], where a similar study performed on a DSP architecture implemented with the same FD-SOI technology, exploiting a performance monitor system based on Timing Fault Sensors (TMFLT).

V_{dd}	0.9 V	0.7 V	0.5 V
C_{corr}	0.6	0.59	0.47
F_0	8.72	5.19	3.21
R-Square	0.995	0.998	0.998

TABLE III: Parameters of the model for a typical chip at three different supply voltages and 25 °C.

C. Temperature Variations

Once the robustness of the PMB model has been proved in a given operating condition, it is possible to generalize the analysis to cover a wide temperature operating range. The measurements described in this section are performed on a single chip, supplied with a given voltage (V_{DD}), at multiple temperatures. The approach is the same we followed in the previous case, the only difference is that the fitted data are now related to three different temperatures, specifically $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$.

Fig. 6 shows the global linear fit. The red solid line, which is described by the equation 1, represents the general model fitting the data. Table IV reports the parameters of the model in this operating condition, as well as the R-Square. As in the previous case, we evaluated the error. Fig. 7 shows the residuals with respect to the fitting curve. As expected, a model which fits data belonging to measurements at various temperatures is affected by a larger error. This phenomenon is caused by the fact that we are assuming that the maximum frequency is limited only by the PMB related to the PMOS transistors, which is usually slower than PMB related to the NMOS ones. Correlation between the maximum frequency of the device and a linear combination of both PMB sensors could improve the accuracy.

D. Process variations

Operating temperature is not the only factor affecting device performance operating in near threshold. In this section we will present a short discussion on the performance variance among different chips fabricated in the UTBB FD-SOI technology node caused by different process corners. Then, we will show how the performance monitoring methodology described so far can be further generalized to cover also multiple chips belonging to various process corners. Particular attention should be given to the nature of the performance variation related to the technological process, the variations introduced by the process can be classified into two types: *i)* The Inter-chip variations, that can be observed in terms of performance gaps between different devices, as shown in Fig. 8 *ii)* Intra-chip variations, as demonstrated by [24], resulting in different first N critical path, that can affect the consistency between the

V_{dd}	0.7 V
C_{corr}	0.59
F_0	5.42
R-Square	0.996

TABLE IV: Parameters of the model fitting the data for a single chip, at a given supply voltage, at three different temperatures, $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$.

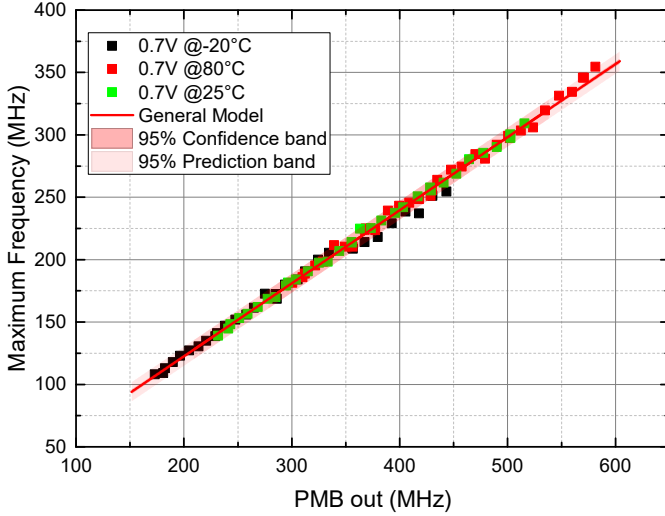


Fig. 6: In this plot are shown all the measurements performed at a single supply voltage $V_{dd} = 0.7\text{ V}$ and three different temperatures $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$. The red solid line represents the curve fitting the data measured at the three different temperature, $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$.

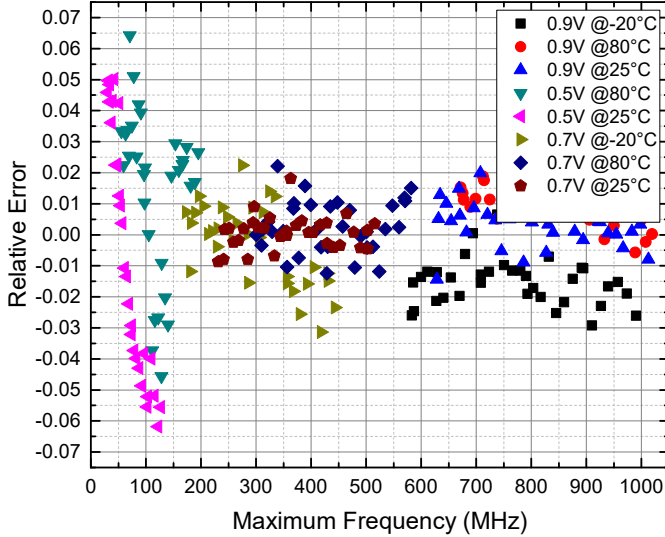


Fig. 7: This plot shows the distribution of the relative error versus the maximum frequency when the *Process-aware Temperature-unaware* model is used to convert the PMB output value in F_{max} .

behaviour of the circuit and an on-chip performance monitor, confirmed also by our analysis.

As in the case of one chip at multiple temperatures, the model has been obtained by performing a global fit on the data related to chips belonging to different process corners, finding the best function describing them with minimum error. Fig. 9 shows the PMB analysis related the chips in exams and the fitting curve.

It can be noted that also in this case the model is represented by the same linear function 1 shown previously. However, in this case the fitted data have a larger variability (Fig. 9), and as a consequence, we observed a larger error and lower R-square, Fig. 10 reports the error. Table V shows the values of the parameters and the R-Square. When the chip was fabricated the technology node was at a very early stage, according to

[38], the variance introduced by the manufacturing process is expected to decrease with the maturity of the process.

V_{dd}	0.7 V
C_{corr}	0.614
F_0	6.86
R-Square	0.88

TABLE V: Parameters of the model fitting the data related to chips belonging to different process corners, at a given supply voltage, at three different temperatures.

E. Body-Bias Model

The final goal of this study is to develop a performance-aware body bias regulation system. To this end, it is necessary to derive a model which links a variation of the V_{BB} to the F_{max} variation. Once the relationship between these two parameters has been obtained, given a performance gap to be compensated, it is possible to determine the necessary amount of body bias voltage to do it.

It is very important to note that the use of body-bias for this application is twofold: *i)* in an ideal case, that is assuming no errors on the F_{max} estimation, V_{BB} represents the knob to simply change the performance of the device *ii)* in a context where F_{max} is affected by uncertainty because of the reasons previously described, V_{BB} margins can be used to compensate F_{max} negative residuals; in other words, all those conditions where the F_{max} is overestimated by the model. Fig.11 shows the relationship between F_{max} and V_{BB} at different temperatures.

Despite the curve shifts with temperature, if we consider the relative frequency variation, we can simplify the body bias model approximating the curve with a linear function having 5%/100 mV (at 0.7V) as slope (Fig. 12). Table VI shows the same analysis at different voltage operating points.

This model allows to compute the additional body bias margin required to compensate the measurement errors described

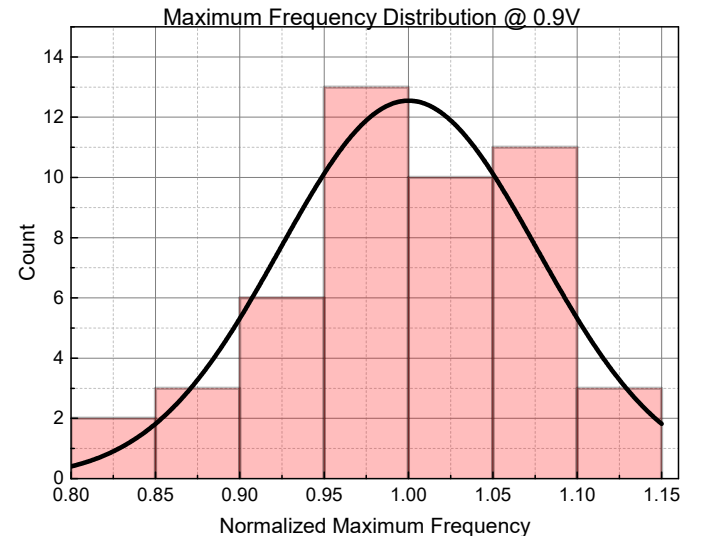


Fig. 8: Normalized distribution of the maximum frequency for the entire population of 48 chips at 0.9 V, 25 °C.

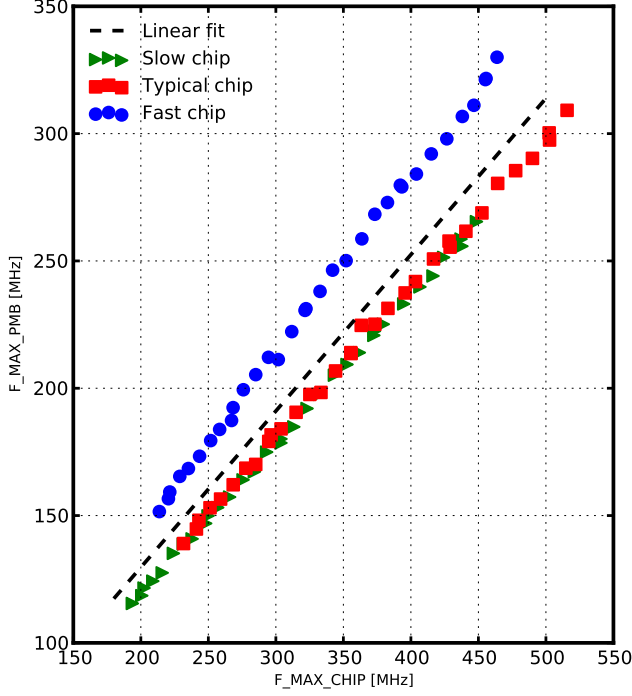


Fig. 9: In this plot are represented the data of 3 devices belonging to *Fast*, *Typical* and *Slow* process corners, the dashed line represents the model fitting the data-set

PULP 0.5 V	0.7 V	0.9 V
11%/100 mV	5%/100 mV	3%/100 mV

TABLE VI: Body bias induced performance gain

in the previous sections. Assuming to compensate the model uncertainty caused by temperature, which was the 4% at 0.7 V, we need to add a margin on V_{BB} of 100 mV.

The parameters of the linear model derived in this section change according to the level of knowledge of the operating condition. If the process and the temperature are known, the model has less uncertainty, and small margins are added to the VBB to compensate for the PMB measurement error (e.g. 3% at 0.7 V). If the temperature is known and the process is unknown, the uncertainty on the frequency reported by PMBs is higher (4% at 0.7 V), and possible slow process corners are compensated by wider FBB margins on the applied VBB (i.e. making the chip faster than what would be needed, if the process corner was not a slow one). If both temperature and process corners are unknown, the same principle applies, the uncertainty on the frequency measurements is further higher (9.7% at 0.7 V), and the controller assumes wider VBB margins to compensate possible slow process corners and the unknown temperature. Additionally, the different body bias margins used to compensate the model errors can be summed as independent contributions. The resulting 9.7% total error reported for VDD = 0.7 V at unknown process and temperature can be compensated by the sum of the body bias margins that would separately compensate each error contribution [29]. Note that, although the parameters of the model used in the

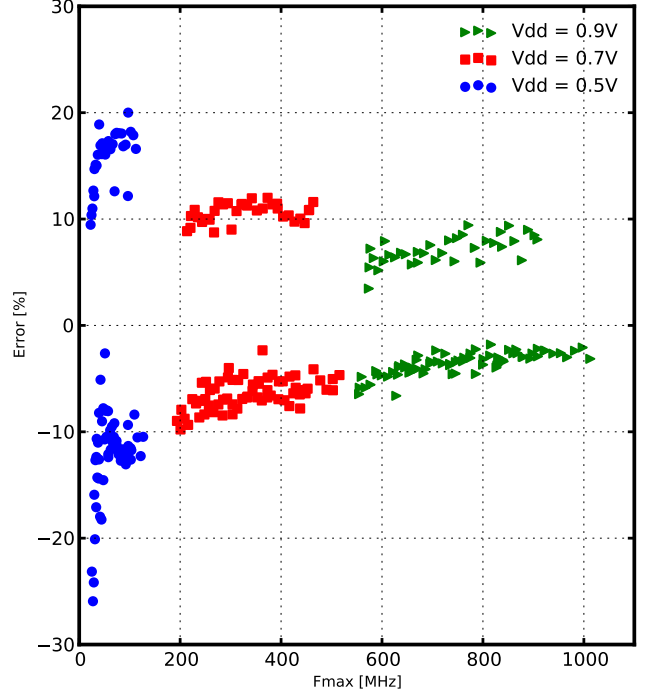


Fig. 10: This plot shows the distribution of the relative error versus the maximum frequency when the process-independent model is used to convert the PMB output value.

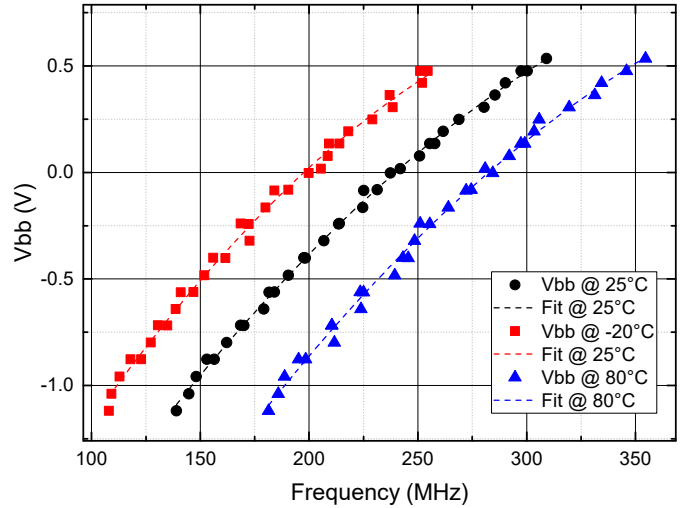


Fig. 11: This plot shows the relation between the maximum frequency of a typical device and the body-bias voltage at three operating points: Vdd = 0.7 V and $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$.

controller change according to the level of knowledge of the temperature and process, the PMBs keep always the same configuration.

As a concluding remark for this section, it can be noted that to minimize the error of the model, a calibration procedure capable to discover the chip process corner is required. This calibration procedure reduces the error by a factor of 2.4 and is described in section V.

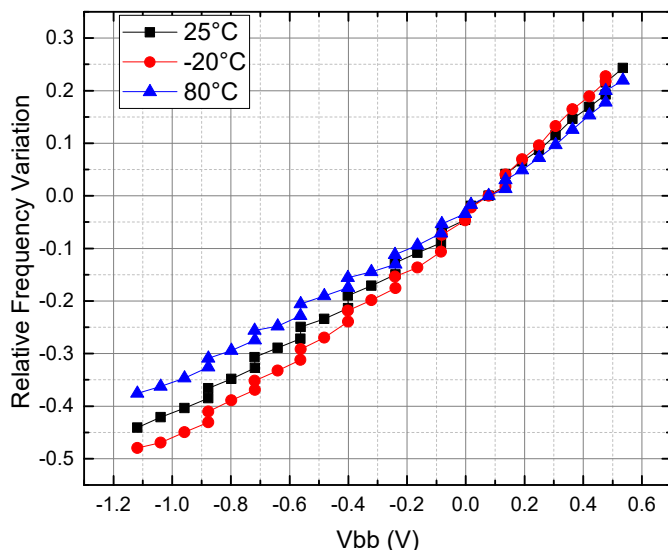


Fig. 12: In this plot is represented the relative frequency variation versus the body-bias voltage, with respect to the $V_{bb} = 0$ V condition. The operating points are: $V_{dd} = 0.7$ V and $T = \{-20^\circ\text{C}, 25^\circ\text{C}, 80^\circ\text{C}\}$.

F. Methodology Overhead

The compensation methodology we propose is based on the application of forward body bias to compensate process and temperature variations. As it is well known, FBB has a direct influence on the leakage current of a circuit, hence it is important to quantify the overhead associated with the additional body bias margins. Specifically, we defined the overhead as the ratio between the additional leakage caused by the body bias margin and the leakage current of the circuit when an ideal controller is regulating the body bias voltage (i.e. capable to apply the exact V_{BB} to achieve a given target frequency). How the additional leakage current affects the global power consumption strongly depends on the dynamic power consumption (i.e operating frequency) of the circuit, hence the overhead will be computed with respect to the leakage current under the *ideal* body bias regulation.

The overhead presented in this section refers to the three different models derived in the previous sections. More specifically, to the cases where the supply voltage is known and *i)* the process corner and the temperature are unknown, we define this model as *Process-Unaware Temperature-Unaware model ii)* the process corner is known and the temperature is unknown, we refer to it as *Process-Aware Temperature-Unaware model iii)* both temperature and process corner are known, which is the *Process-Aware Temperature-Aware model*.

When the *Process-Unaware Temperature-Unaware* model is used, we experience the highest overhead. As expected, the higher is the uncertainty of the model, the wider is the body bias safety margin required to use the model, and also the higher is the associated overhead because of the higher leakage current. However, having the possibility to determine the chip process corner and calibrate the model on the specific device (i.e. passing to the *Process-Aware Temperature-Unaware* model) the overhead can be significantly reduced. If we consider the leakage current associated with the body

bias margins needed to correct the errors of the *Process-Aware Temperature-Unaware* model, we can observe that the overhead is much lower (Table VII). Finally, probing also the temperature, it would be possible to use multiple, better correlated, *Process-Aware Temperature-Aware* PMB models at different temperatures. Table VII summarizes the results in terms of overhead.

V_{dd}	0.9 V	0.7 V	0.5 V
Proc-unaware/Temp-unaware			
F_{err}	6.6%	9.7%	25%
V_{BB} margin [mV]	150	150	200
Leakage overhead	33%	37%	66%
Proc-aware/Temp-unaware			
F_{err}	3%	4%	7%
V_{BB} margin [mV]	100	100	100
Leakage overhead	13%	14%	15%
Proc-aware/Temp-aware			
F_{err}	2%	3%	6%
V_{BB} margin [mV]	50	50	50
Leakage overhead	9%	10%	12%

TABLE VII: Frequency error of the three presented models and power consumption overheads with respect to the ideal compensation where no body bias margins are used.

V. MODEL UTILIZATION

A. Calibration Procedure

As we have demonstrated it in section IV, the model derived from the PMB frequency estimations of chips belonging to different process corners is affected by a significant error. However, this error can be reduced by adopting a calibration procedure. More specifically, the operation performed during the characterization phase on the testing equipment for a single temperature operating point can be replicated on the embedded board. The constraint to ensure proper calibration is to execute the procedure at a constant temperature. This goal can be achieved by exploiting an on-chip temperature sensor. For a correct calibration, the absolute value of the temperature at which the calibration is executed is not relevant, however the temperature must remain constant to not influence the characterization of the PMB sensor. Then, starting from a controller calibrated on a constant temperature value, two scenarios are possible: *i)* the temperature does not change with respect to the calibration point, and the controller has to compensate with VBB margins the 3% of model uncertainty at 0.7V (Fig. 5). *ii)* the temperature changes with respect to the calibration point, and the controller has to compensate a slightly higher model uncertainty, which takes into account the effects of temperature change, which is approximately 4% at 0.7V (Fig. 6).

Here we describe the calibration procedure. At the boot, the benchmark application mentioned in section IV is loaded in the device memory by an external microcontroller connected through an SPI interface. Then the application is executed in a loop which increases the operating frequency at every iteration by 1MHz. The frequency is increased until the chip starts to

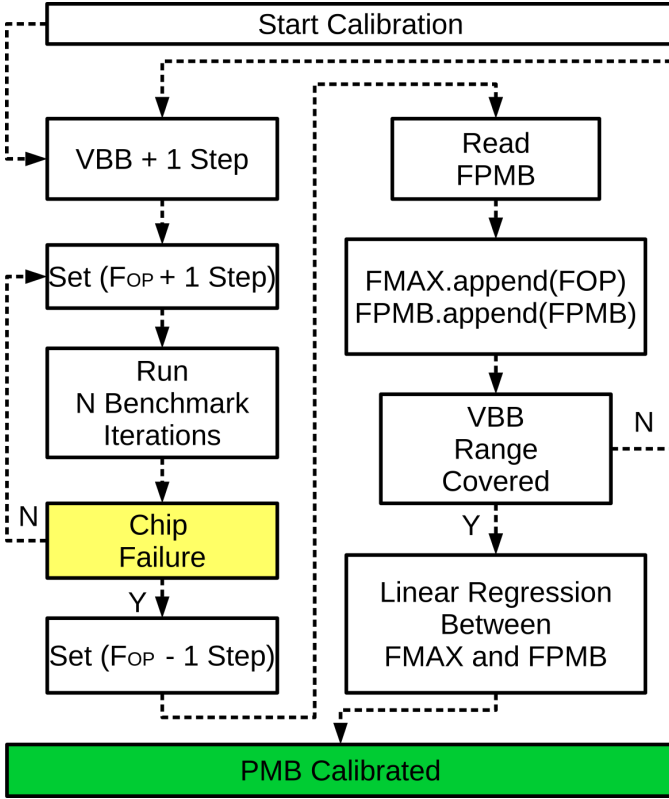


Fig. 13: Calibration procedure executed on the embedded board to determine the process corner and tune the PMB model.

fail, either returning wrong arithmetic results or completely failing. This test is executed on multiple points covering the entire body bias range to correlate the maximum frequency of the chip with the PMB frequency estimation.

Using a model calibrated on the specific process corner, when $V_{DD} = 0.7V$, it is possible to pass from the *Process-unaware Temperature-unaware* to the *Process-aware Temperature-unaware* (Table VII), reducing the error by 2.4X. The duration of a complete calibration at a single VDD operating point depends on the minimum step of the sweep performed on VBB and the maximum frequency. Additionally, the maximum frequency search method may increase the duration of the calibration procedure (i.e. linear sweep vs binary search). Finally the duration of the benchmark executed by the processor also changes the calibration time. In our example, we used a 50 mV body bias voltage step to span the 1.5 V range, resulting in 30 VBB points, and a linear sweep with a 1 MHz step for the maximum frequency search. At every operating condition we executed 10000 benchmark iterations. The overall duration, starting from a 100 MHz initial frequency for the maximum frequency search, lasted approximately 6 seconds. The procedure to perform this operation is illustrated as block diagram representation in Figure 13.

B. Body-Bias Controller

In this section we will show how the models derived so far and selected with the procedure described in the previous section can be used in practical applications. More specifically, an accurate on-chip performance feedback, as well as

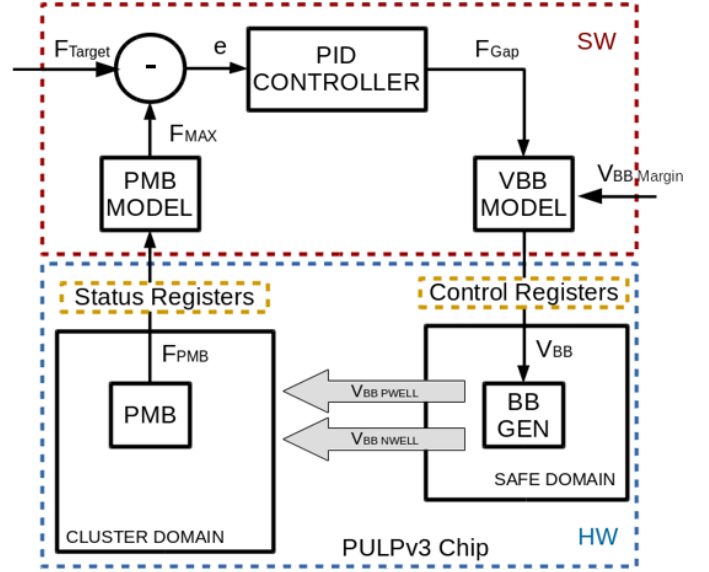


Fig. 14: Block diagram representation of the body-bias controller.

the modeled behavior of the body-bias generator enable the building of software control systems; the structure we propose is a mixed Hardware-Software solution which properly set the cores cluster Body-Bias voltage depending on the clock frequency set-point provided as input.

C. Controller Modules

The proposed body-bias controller architecture is composed of: *i*) a feedback, *ii*) a subtractor which computes the mismatch between feedback and set-point, *iii*) a Proportional-Integrative-Derivative (PID) controller, *iv*) an actuator fed with the PID controller output (VBB MODEL + BB GEN). Fig. 14 shows a block diagram of the body bias control system.

The control system *feedback* module is divided in two main building blocks: *i*) a hardware component, the previously mentioned PMB, which returns a raw maximum frequency estimation; *ii*) the software model derived for the PMB, which allows to properly convert the raw output of the sensor in a clock frequency value.

The *actuator*, as in the case of feedback module, is composed of two building blocks: the on-chip Body-Bias generator, paired with a software model which allows to determine the right amount of Body-Bias needed to fill the frequency gap between feedback and set-point.

The *PID controller* is a standard well-known control mechanism used in industrial control systems [39], the proposed body bias control system is entirely implemented in software and it is fed with the mismatch between feedback and set-point frequency. The PID controller has to be tuned accordingly to the feedback and the actuator, this process can be done empirically. In the context of this control task we decided to adopt a parameter set which minimizes the settling time while keeping under control undershoots. Note that, for this kind of system, undershoots are more critical than overshoots. During an overshoot, the chip is biased with more than necessary FBB. Therefore, For a short amount of time, the maximum

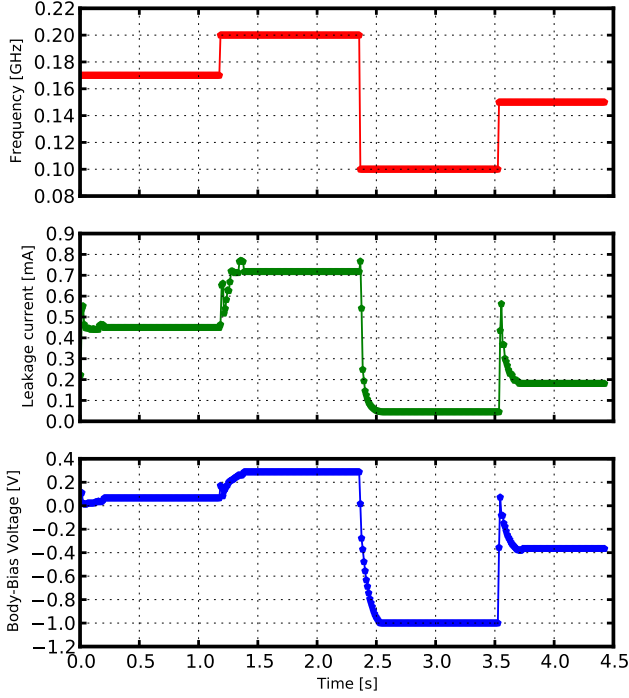


Fig. 15: Body bias voltage regulation operated by the controller during boosting and leakage reduction phase, respectively. Supply voltage is 0.7V

frequency of the chip is faster than than expected . On the contrary, during undershoots, the controller applies less FBB than required (or even RBB). Therefore, the actual maximum frequency of the chip can be lower than the requested frequency. This condition is critical and very likely determines a chip failure

As shown in Fig. 14, the entry point of the control system is the frequency set-point F_{Target} . The first operation performed by the controller is to measure the maximum frequency of the system. Once the output of the PMB module is ready, it is converted in a frequency value, and it is compared with the input set-point obtaining the frequency mismatch. Then, the PID module is fed with the frequency mismatch and its output is sent to the body-bias software model. Finally, once the controller computes the body-bias voltage to apply, it is used to set the body-bias generator to the new V_{BB} value.

D. Controller operation

1) *Frequency Tracking*: In the following we present a controller working example. Fig. 15 shows the V_{BB} regulation operated by the controller when the frequency set-point is changed to different values ⁴.

The first set-point, 175MHz, is very close to the maximum frequency of the device. In this condition the controller applies a very small amount of body-bias, since the system is already capable to run at the requested frequency, without

⁴Every time a new frequency is requested, the controller reset the body bias to $V_{BB} = 0V$ in a single step, and uses this voltage as a starting point. This choice has been taken to easy the implementation and reduce on the average the regulation time when changing from one frequency set-point to the next

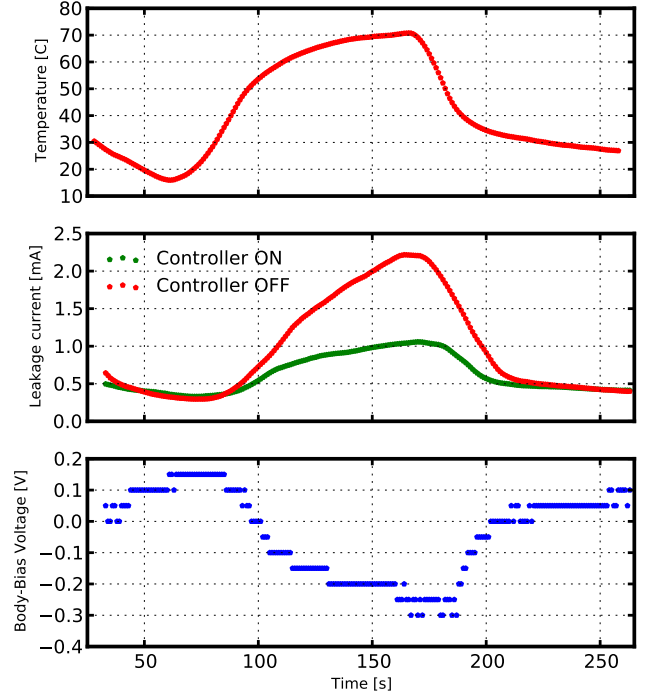


Fig. 16: This plot shows in parallel the body-bias voltage modulation, the environmental temperature variation and the leakage current when the body bias controller is active.

additional forward body-bias. The only body bias voltage applied ($V_{BB} \neq 0$) is the body bias margin. When 200MHz are requested as set-point, the controller applies a forward body-bias voltage of approximately 300mV; applying a strong forward body bias, the system is now capable to run at 200MHz. On the contrary, when a slow frequency is set as set-point, 100MHz, the controller applies an aggressive reverse body-bias to reduce leakage power as much as possible. Finally, when the set-point is set to 150MHz, the system reduce the reverse body bias, applying to -400mV.

2) *Temperature tracking*: Here we show another working example. In this case the controller is used to guarantee a 170MHz frequency target against environmental temperature changes. Fig. 16 shows the behaviour of temperature, leakage and body-bias voltage when the controller is active.

From the leakage current plot it is clear that the controller is compensating the environmental temperature increase applying a reverse body-bias voltage. On the contrary, when the temperature decreases, to guarantee the requested target frequency the controller applies a forward body-bias voltage to boost the maximum performance of the chip.

VI. RESULTS

In this section we present results in terms of energy efficiency gain and leakage reduction when the body-bias controller is turned on.

Fig. 17 shows a comparison between the leakage current when *i)* the controller is off and the leakage is not compensated against temperature *ii)* the controller is on and it is regulating the V_{BB} using the *Process-aware Temperature-unaware* model

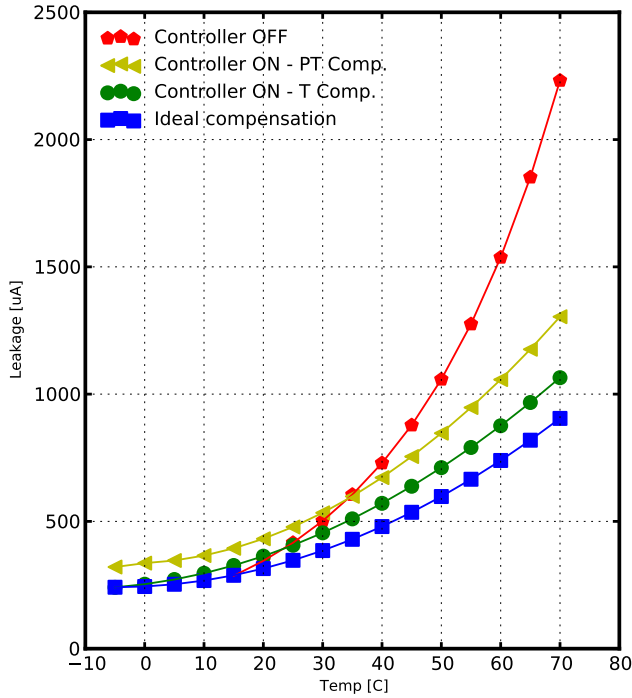


Fig. 17: This plot shows the leakage current versus temperature, the operating frequency is 170MHz and the supply voltage is 0.7V. The red curve reports the Leakage current when the controller is off. The green curve shows the Leakage current when the controller is on; the process corner is known. The yellow curve refers to the case where the controller is on but the process corner is not known. To guarantee the operating frequency, additional FBB margins are used, resulting in a leakage current increase. Blue curve shows an ideally compensated leakage current, i.e. when no FBB margins are used by the controller. Note that below 17 °C it is not possible to sustain the 170MHz target frequency without applying forward body biasing, hence the circuit stops to work.

iii) the controller is on and it is regulating the V_{BB} using the *Process-unaware Temperature-unaware* model iv) the ideal case where the leakage is compensated without additional body bias margins, hence assuming no error in the model converting the frequency estimation of the PMB. Note that the influence of the leakage power on the global power consumption depends on the operating frequency of the device, in our measurements the controller is regulating the body bias to achieve a target frequency of 170MHz, which represents an optimal operating point, since the working frequency is close to the maximum one.

Fig. 17 also shows that the margin required to compensate process variations causes a big increase in the leakage current. In this context the effect on the global power consumption is not negligible. However this problem can be solved by running a simple calibration procedure which determines if the chip is a *Slow*, *Typical* or a *Fast* one.

The benefit of the calibration procedure is to reduce the error of the model, allowing to use the *Process-aware Temperature-unaware* model instead of the *Process-unaware Temperature-unaware*. As a consequence, the body bias margin to apply is reduced from 150mV to 100mV.

Fig. 18 compares the power consumption of the chip, at different temperatures, both in the case when the body bias

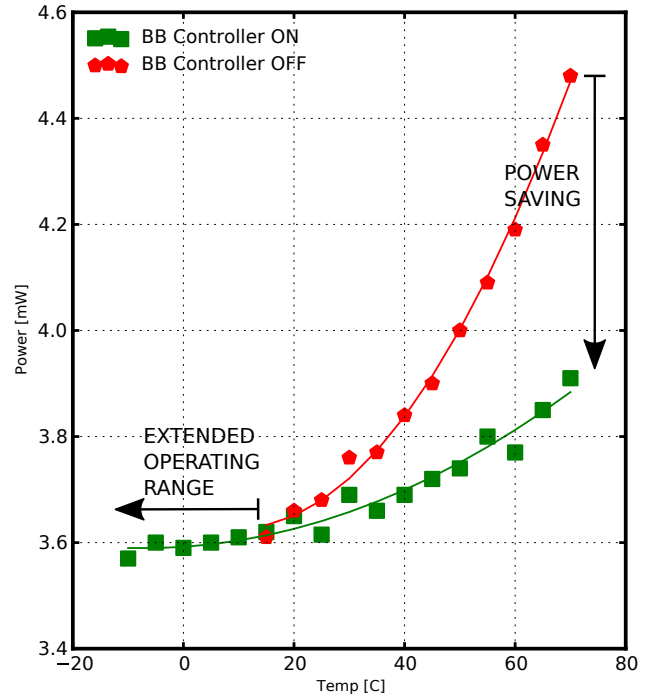


Fig. 18: Comparison between power consumption with and without the performance aware body bias controller performed at 170MHz, $V_{DD} = 0.7V$. Note that below 17 °C is not possible to run at 170MHz without compensating the performance degradation with additional body biasing (i.e. with the controller on).

controller is turned on and off. As it is evident from the plot, the use of the body bias controller allows to improve the energy efficiency by 15% (at high temperatures), and extend the operating range to very low temperatures necessary to deal with the wide range of environmental conditions of IoT devices, leveraging a low-margins design methodology for energy efficient implementation.

VII. CONCLUSION

In this work we demonstrated the body-biasing capabilities of the UTBB FD-SOI technology in compensating temperature and process variations.

We performed an analysis on Process Monitoring Blocks to correlate the on-chip estimated maximum frequency to the actual performance of the device, hence we derived a calibration model for the PMB sensors. We extended the model obtained at a single temperature to multiple temperatures, to compensate the effect that environmental temperature changes have on the performance of the device. Then, we further generalized the model to the situation where the process corner is unknown, considering it as an uncertainty of the model. We found that at 0.7V the frequency estimation error is in the order of 3% when both process and temperature are determined, 4% if the temperature is unknown and 9.7% when also the process corner is unknown. We developed a strategy to eliminate such frequency estimation error by over-compensating with selective body bias margins. We also derived an on-board calibration procedure capable

to determine the process corner of the chip and to select the proper *Process-aware Temperature-unaware* model. Once the models have been obtained, we developed a closed loop control strategy for the body bias voltage based on a simple PID controller. We tested the controller on a real embedded platform, demonstrating that with a minimum overhead, an automatic body bias regulation can reduce by a factor of 2 the leakage power consumption caused by temperature changes. Finally we evaluated the effect of the leakage reduction in a common operating context, demonstrating that our approach can introduce a 15% energy efficiency improvement on the global power consumption.

ACKNOWLEDGMENT

We thank STMicroelectronics for chip fabrication. This work was supported by European project ExaNoDe (H2020-671578) and European project OPRECOMP (732631).

REFERENCES

- [1] R. G. Dreslinski, M. Wiecekowsky, D. Blaauw, D. Sylvester, and T. N. Mudge, "Near-threshold computing: Reclaiming moore's law through energy efficient integrated circuits," *Proceedings of the IEEE*, vol. 98, no. 2, pp. 253–266, 2010.
- [2] D. Markovic, C. C. Wang, L. P. Alarcón, T. Liu, and J. M. Rabaey, "Ultralow-power design in near-threshold region," *Proceedings of the IEEE*, vol. 98, no. 2, pp. 237–252, 2010.
- [3] D. Rossi, A. Pullini, I. Loi, M. Gautschi, F. K. Gürkaynak, A. Teman, J. Constantin, A. Burg, I. Miro-panades, E. Beigné, P. Flatresse, and L. Benini, "Near-Threshold Parallel Computing : The PULPv2 Cluster," 2017.
- [4] A. Pahlevan, J. Picorel, A. P. Zarandi, D. Rossi, M. Zapater, A. Bartolini, P. G. D. Valle, D. Atienza, L. Benini, and B. Falsafi, "Towards near-threshold server processors," in *2016 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2016, pp. 7–12.
- [5] M. Alioto, Ed., *Enabling the Internet of Things*. Springer International Publishing, 2017.
- [6] W. Lee, Y. Wang, T. Cui, S. Nazarian, and M. Pedram, "Dynamic thermal management for FinFET-based circuits exploiting the temperature effect inversion phenomenon," in *Proceedings of the 2014 International Symposium on Low Power Electronics and Design*, ser. ISLPED '14. New York, NY, USA: ACM, 2014, pp. 105–110. [Online]. Available: <http://doi.acm.org/10.1145/2627369.2627608>
- [7] Y. Pu, X. Zhang, J. Huang, A. Muramatsu, M. Nomura, and K. Hirairi, "Misleading Energy and Performance Claims in Sub / Near Threshold Digital Systems," *Technology*, pp. 4–6, 2010.
- [8] K. Han, J.-J. Lee, J. Lee, W. Lee, and M. Pedram, "TEI-NoC: Optimizing Ultra-Low Power NoCs Exploiting the Temperature Effect Inversion," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 0070, no. c, pp. 1–1, 2017. [Online]. Available: <http://ieeexplore.ieee.org/document/7896585/>
- [9] M. Alioto, "Ultra-Low Power {VLSI} Circuit Design Demystified and Explained: A Tutorial," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 59, no. 1, pp. 3–29, jan 2012.
- [10] K. Sundaresan, K. Brouse, K. U-Yen, F. Ayazi, and P. Allen, "A 7-MHz process, temperature and supply compensated clock oscillator in 0.25 μ m CMOS," in *Proceedings of the 2003 International Symposium on Circuits and Systems*, 2003. *ISCAS '03*. IEEE, 2003.
- [11] G. Gammie, A. Wang, M. Chau, S. Gururajao, R. Pitts, F. Jumel, S. Engel, P. Royannez, R. Lagerquist, H. Mair, J. Vaccani, G. Baldwin, K. Heragu, R. Mandal, M. Clinton, D. Arden, and U. Ko, "A 45nm 3.5g baseband-and-multimedia application processor using adaptive body-bias and ultra-low-power techniques," in *2008 IEEE International Solid-State Circuits Conference - Digest of Technical Papers*. IEEE, feb 2008.
- [12] S. Clerc, M. Saligane, F. Abouzeid, M. Cochet, J.-M. Daveau, C. Bottoni, D. Bol, J. De-Vos, D. Zamora, B. Coeffic *et al.*, "8.4 a 0.33 v/-40 c process/temperature closed-loop compensation soc embedding all-digital clock multiplier and dc-dc converter exploiting fdsoi 28nm back-gate biasing," in *Solid-State Circuits Conference-(ISSCC), 2015 IEEE International*. IEEE, 2015, pp. 1–3.
- [13] D. Rossi, A. Pullini, M. Gautschi, I. Loi, F. K. Gürkaynak, P. Flatresse, and L. Benini, "A 60 gops/w,-1.8 v to 0.9 v body bias ulp cluster in 28nm utbb fd-soi technology," in *2015 IEEE SOI-3D-Subthreshold Microelectronics Technology Unified Conference (S3S)*. IEEE, oct 2015.
- [14] D. Rossi, I. Loi, F. Conti, L. Benini, C. Müller, and A. Burg, "A Self-Aware Architecture for PVT Compensation and Power Nap in Near-Threshold Processors," pp. 46–53, 2017.
- [15] M. Seok, G. K. Chen, S. Hanson, M. Wiecekowsky, D. T. Blaauw, and D. Sylvester, "{CAS-FEST} 2010: Mitigating Variability in Near-Threshold Computing," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 1, no. 1, pp. 42–49, 2011.
- [16] D. Ernst, N. S. Kim, S. Das, S. Pant, R. Rao, T. Pham, C. Ziesler, D. Blaauw, T. Austin, K. Flautner, and T. Mudge, "Razor: A low-power pipeline based on circuit-level timing speculation," *Proceedings of the Annual International Symposium on Microarchitecture, MICRO*, vol. 2003-January, pp. 7–18, 2003.
- [17] S. Das, S. Pant, D. Roberts, S. Lee, D. Blaauw, T. Austin, T. Mudge, and K. Flautner, "A self-tuning DVS processor using delay-error detection and correction," *IEEE Symposium on VLSI Circuits, Digest of Technical Papers*, vol. 2005, no. 4, pp. 258–261, 2005.
- [18] D. Blaauw, S. Kalaiselvan, K. Lai, W.-H. Ma, S. Pant, C. Tokunaga, S. Das, and D. Bull, "Razor II: In situ error detection and correction for PVT and SER tolerance," in *2008 IEEE International Solid-State Circuits Conference - Digest of Technical Papers*. IEEE, feb 2008.
- [19] D. Bull, S. Das, K. Shivashankar, G. S. Dasika, K. Flautner, and D. Blaauw, "A Power-Efficient 32 bit {ARM} Processor Using Timing-Error Detection and Correction for Transient-Error Tolerance and Adaptation to {PVT} Variation," *IEEE Journal of Solid-State Circuits*, vol. 46, no. 1, pp. 18–31, jan 2011.
- [20] M. Fojtik, D. Fick, Y. Kim, N. Pinckney, D. M. Harris, D. Blaauw, and D. Sylvester, "Bubble razor: Eliminating timing margins in an ARM cortex-M3 Processor in 45 nm CMOS using architecturally independent error detection and correction," *IEEE Journal of Solid-State Circuits*, vol. 48, no. 1, pp. 66–81, 2013.
- [21] A. Drake, R. Senger, H. Deogun, G. Carpenter, S. Ghiasi, T. Nguyen, N. James, M. Floyd, and V. Pokala, "A distributed critical-path timing monitor for a 65nm high-performance microprocessor," *Digest of Technical Papers - IEEE International Solid-State Circuits Conference*, 2007.
- [22] J. Tschanz, K. Bowman, S. Walstra, M. Agostinelli, T. Karnik, and V. De, "Tunable replica circuits and adaptive voltage-frequency techniques for dynamic voltage, temperature, and aging variation tolerance," *2009 Symposium on VLSI Circuits*, pp. 112–113, 2009.
- [23] A. Mhz, V. Ghz, E. Beigné, A. Valentian, I. Miro-panades, R. Wilson, P. Flatresse, F. Abouzeid, T. Benoist, C. Bernard, S. Bernard, O. Billoint, S. Clerc, B. Giraud, A. Grover, J. L. Coz, J.-p. Noel, O. Thomas, and Y. Thonnart, "A 460 mhz at 397 mv, 2.6 ghz at 1.3 v, 32 bits vliw dsp embedding fmax tracking," vol. 50, no. 1, pp. 125–136, 2015.
- [24] M. Zandrahimi, Z. Al-Ars, P. Debaudand, and A. Castillejo, "Challenges of using on-chip performance monitors for process and environmental variation compensation," in *Proceedings of the 2016 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. Research Publishing Services, 2016.
- [25] J. Constantin, A. Bonetti, A. Teman, C. Müller, L. Schmid, and A. Burg, "DynOR: A 32-bit microprocessor in 28 nm FD-SOI with cycle-by-cycle dynamic clock adjustment," *European Solid-State Circuits Conference*, vol. 2016-October, pp. 261–264, 2016.
- [26] K. A. Bowman, J. W. Tschanz, S. L. L. Lu, P. A. Aseron, M. M. Khellah, A. Raychowdhury, B. M. Geuskens, C. Tokunaga, C. B. Wilkerson, T. Karnik, and V. K. De, "A 45 nm resilient microprocessor core for dynamic variation tolerance," *IEEE Journal of Solid-State Circuits*, vol. 46, no. 1, pp. 194–208, 2011.
- [27] M. Miyazaki, G. Ono, T. Hattori, K. Shiozawa, K. Uchiyama, and K. Ishibashi, "A 1000-MIPS/w microprocessor using speed adaptive threshold-voltage CMOS with forward bias," in *2000 IEEE International Solid-State Circuits Conference. Digest of Technical Papers (Cat. No.00CH37056)*. IEEE, 2000.
- [28] J. W. Tschanz, J. T. Kao, S. G. Narendran, R. Nair, D. A. Antoniadis, A. P. Chandrakasan, and V. De, "Adaptive body bias for reducing impacts of die-to-die and within-die parameter variations on microprocessor frequency and leakage," *IEEE Journal of Solid-State Circuits*, vol. 37, no. 11, pp. 1396–1402, 2002.
- [29] S. Kumar, C. Kim, and S. Sapatnekar, "Body bias voltage computations for process and temperature compensation," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 16, no. 3, pp. 249–262, mar 2008.

- [30] J. Tschanz, N. Kim, S. Dighe, J. Howard, G. Ruhl, S. R. Vangal, S. Narendra, Y. Hoskote, H. Wilson, C. Lam, M. Shuman, C. Tokunaga, D. Somasekhar, S. Tang, D. Finan, T. Karnik, N. Borkar, N. A. Kurd, and V. De, "Adaptive frequency and biasing techniques for tolerance to dynamic temperature-voltage variations and aging," in *2007 IEEE International Solid-State Circuits Conference, ISSCC 2007, Digest of Technical Papers, San Francisco, CA, USA, February 11-15, 2007*. IEEE, 2007, pp. 292–604.
- [31] K. Kang, S. P. Park, K. Kim, and K. Roy, "On-chip variability sensor using phase-locked loop for detecting and correcting parametric timing failures," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 18, no. 2, pp. 270–280, 2010.
- [32] S. V. Kumar, C. H. Kim, and S. S. Sapatnekar, "Adaptive techniques for overcoming performance degradation due to aging in CMOS circuits," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 19, no. 4, pp. 603–614, 2011.
- [33] G. Ono, M. Miyazaki, H. Tanaka, N. Ohkubo, and T. Kawahara, "Temperature referenced supply voltage and forward-body-bias control (TSFC) architecture for minimum power consumption [ubiquitous computing processors]," in *Proceedings of the 30th European Solid-State Circuits Conference*. IEEE, 2004.
- [34] P. Flatresse, B. Giraud, J. P. Noel, B. Pelloux-Prayer, F. Giner, D. K. Arora, F. Arnaud, N. Planes, J. L. Coz, O. Thomas, S. Engels, G. Cesana, R. Wilson, and P. Urard, "Ultra-wide body-bias range LDPC decoder in 28nm UTBB FDSOI technology," *Digest of Technical Papers - IEEE International Solid-State Circuits Conference*, vol. 56, no. 424, pp. 424–425, 2013.
- [35] F. Conti, D. Rossi, A. Pullini, I. Loi, and L. Benini, "PULP: A ultra-low power parallel accelerator for energy-efficient and flexible embedded vision," *Journal of Signal Processing Systems*, vol. 84, no. 3, pp. 339–354, nov 2015.
- [36] M. Blagojevic, M. Cochet, B. Keller, P. Flatresse, A. Vladimirescu, and B. Nikolic, "A fast, flexible, positive and negative adaptive body-bias generator in 28nm FDSOI," *IEEE Symposium on VLSI Circuits, Digest of Technical Papers*, vol. 2016-Sept, pp. 9–10, 2016.
- [37] E. Beigne, A. Valentian, I. Miro-Panades, R. Wilson, P. Flatresse, F. Abouzeid, T. Benoist, C. Bernard, S. Bernard, O. Billoint, S. Clerc, B. Giraud, A. Grover, J. L. Coz, J.-P. Noel, O. Thomas, and Y. Thonnart, "A 460 MHz at 397 mV, 2.6 GHz at 1.3 v, 32 bits VLIW DSP embedding f MAX tracking," *IEEE Journal of Solid-State Circuits*, vol. 50, no. 1, pp. 125–136, jan 2015.
- [38] R. S. Mackay, H. Kamberian, and Y. Zhang, "Methods to reduce lithography costs with reticle engineering," *Microelectronic Engineering*, vol. 83, no. 4-9, pp. 914–918, 2006. [Online]. Available: <http://linkinghub.elsevier.com/retrieve/pii/S0167931706001742>
- [39] G. Y. Wei and M. Horowitz, "Fully digital, energy-efficient, adaptive power-supply regulator," *IEEE Journal of Solid-State Circuits*, vol. 34, no. 4, pp. 520–528, 1999.