

# Exploration of Fine-Grain Body Bias Control in Many-Core Processor Arrays

Bevan Baas, Brent Bohnenstiehl, and Jin Cui  
 Department of Electrical and Computer Engineering  
 University of California, Davis  
 Davis, CA USA  
 {bbaas,bvbohnen,tkxcui}@ucdavis.edu

**Abstract**—Although fine-grain many-core processor arrays have demonstrated great increases in performance, energy efficiency, and area efficiency across many workload domains, methods to integrate the control and optimization of body bias into these arrays has not been well explored. We investigate circuits to implement per-core body bias, and simulate complex workloads to estimate the net benefits in a many-core processor array when body bias is jointly optimized along with per-core supply voltage and clock frequency scaling. Compared to a system utilizing per-core clock frequency and supply voltage tuning, adding per-core body-bias voltage tuning decreases energy dissipation by a mean of 22%–30% at the same throughput for four 900+ core applications.

**Keywords**—DVFS; FD-SOI; GALS; Globally asynchronous locally synchronous; many core; parallel processor; SOI.

## I. INTRODUCTION AND BACKGROUND

Increasing levels of integration and the performance gains possible through parallel processing have motivated the design of single die containing large numbers of processors (many-core processor arrays). To maintain reasonable levels of power dissipation in these chips, new methods are needed to increase efficiency. One fruitful direction in this search is to adapt the clock frequency and power supply voltage ( $V_{dd}$ ) over small domains such as per core or per a portion of each core [1].

In addition, fully-depleted silicon on insulator (FD-SOI) technologies in combination with the control of the body-bias voltage ( $V_{bb}$ ) has shown great promise in enabling circuits able to adapt their efficiency versus performance [2].

To explore the results of using a many-core processor array with adaptable per-core clock frequency, power supply voltage, and body bias, we model the 1000-processor KiloCore chip [3] which is a massively parallel computing platform. The many-core chip is energy efficient for a wide variety of workloads, capable of high performance, easily scalable to higher processor counts, and suitable for a range of applications and critical kernels, acting either alone or as a coprocessor in a heterogeneous system [4].

## II. PROPOSED BODY BIAS CIRCUITS

Because the overhead of implementing an independent on-die  $V_{bb}$  generator per core (or portion of a core) would be prohibitive, we propose using parallel power grids to provide pseudo-independent  $V_{dd}$  and  $V_{bb}$  voltages throughout 100s or 1000s of domains per chip. The main advantage of this method is that circuits can obtain many of the benefits of arbitrarily-tunable voltages while the voltage converters can be built more efficiently at the chip level or off chip.

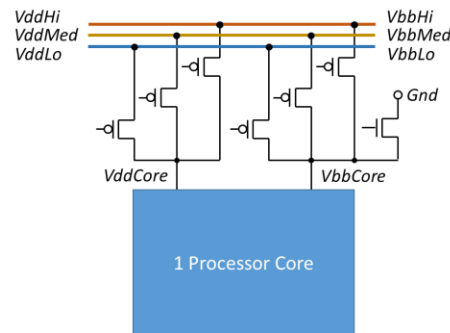


Fig. 1. Schematic illustrating the first proposed architecture with parallel power grids shared between  $V_{dd}$  and  $V_{bb}$  (three grids, one domain per core, and only one  $V_{bb}$  polarity shown in this example).

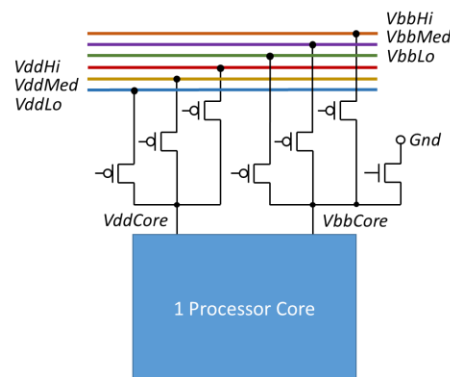


Fig. 2. Schematic illustrating the second proposed architecture with independent parallel power grids for  $V_{dd}$  and  $V_{bb}$  (three grids per core voltage and only one  $V_{bb}$  polarity shown in this example).

The schematic in Fig. 1 shows the first architecture where both  $V_{dd}$  and  $V_{bb}$  share the same power grids requiring a joint optimization for voltage selection, but simplifying the power grid requirements.

Similarly, the schematic in Fig. 2 shows the second architecture where  $V_{dd}$  and  $V_{bb}$  each utilize independent power grids enabling greater flexibility and range of values.

### III. EXPERIMENTAL METHODS AND RESULTS

To evaluate the effectiveness of the proposed  $V_{dd}$  and  $V_{bb}$  selection methods, four complex workloads executing on the KiloCore processor array were simulated on a cycle-accurate simulator with detailed sub-instruction energy calculations using data measured from silicon. The simulator accurately measures active and idle periods while taking into account all architectural features within each core including artifacts of differing clock frequencies, as well as inter-core transfers [5].

Simulated data are scaled to STMicroelectronics' 28 nm FD-SOI [6] using data from spice simulations of simple circuits for a wide variety of supply voltages and body-bias voltages, and LVT and RVT transistor types.

The goal of the optimization is to find the optimum chip-level Hi/Med/Lo voltages for  $V_{dd}$  (and  $V_{bb}$  for the second architecture). As part of this process, the optimal  $V_{dd}$  and  $V_{bb}$  selection of each core in the array must be determined.

Optimization begins with application profiling. A series of tests determine the lowest frequency at which each core may run without reducing application throughput, along with the core's leakage and switching energies. Possible voltage grid selections are then explored, where each core is assigned a combination of  $V_{dd}$  and  $V_{bb}$  which will give the greatest energy reduction while meeting the core's frequency requirement. Running at optimum voltages, energy for the entire array is minimized while maintaining full throughput.

#### A. Simulated workloads

Four complete software applications which execute on the KiloCore array [7] were analyzed. They are:

- Advanced Encryption Standard (AES) with 128-bit keys, using 968 cores;
- 4096-point complex Fast Fourier Transform (FFT), using 975 cores;
- Low Density Parity Check (LDPC) decoder with a 4095-bit code length, using 963 cores; and
- Sort of 100-byte data records with 10-byte keys, using 1000 cores.

An important parameter to optimize the efficiency of workloads on a many-core array is the range of activity levels across cores. As an example, Fig. 3 shows that the LDPC application allows many cores to run as low as 20% of the maximum frequency.

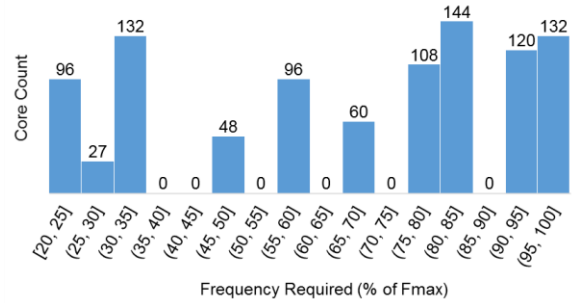


Fig. 3. Histogram of core frequency requirements when operating at the maximum attainable throughput for the 963-processor LDPC application.

#### B. Results

Voltage optimization was performed for five supply circuit options: (1) single  $V_{dd} = 1.1$  V grid with  $V_{bb} = 0$ , (2) single variable  $V_{dd}$  grid with  $V_{bb} = 0$ , (3) three  $V_{dd}$  grids with  $V_{bb} = 0$ , (4) three grids shared between  $V_{dd}$  and  $V_{bb}$  (architecture 1), and (5) three grids each for  $V_{dd}$  and  $V_{bb}$  (architecture 2). To explore optimization effectiveness when the core array is operating below maximum throughput, overall application throughput is swept from 20% to 100% of maximum.

Figures 4 and 5 show the benefits of  $V_{bb}$  adjustment. Metrics are normalized against the case of three  $V_{dd}$  voltage grids with a fixed  $V_{bb} = 0$ .

Figure 4 shows energy reductions by the applications and a linear fit of their geometric means. Mean energy dissipation reductions range from 16% at 100% activity to 29% at 30% activity.

Figure 5 shows the same data for circuit architecture 2. Mean energy dissipation reductions range from 22% at 100% activity to 30% at 30% activity. This 30% activity point is notable in that it is when applications begin to reach the minimum allowed voltage in our tests, 0.4 V. Lower activity conditions are penalized with increased energy consumption due to this constraint.

Figure 6 explores the benefit in moving from architecture 1 to 2. Mean energy savings vary from 2% to 11% across the activities. Variance across sample points is due to their sensitivity to the voltage selections by the optimizer, with the largest gains found when architecture 2 selects  $V_{bb}$  voltages significantly different from  $V_{dd}$  voltages—a situation architecture 1 cannot match.

Figure 7 shows the expected improvements in energy efficiency for each architectural variant, using a geometric mean of the four applications. Improvements are calculated relative to energy spent with  $V_{dd}$  at maximum voltage and  $V_{bb} = 0$ .

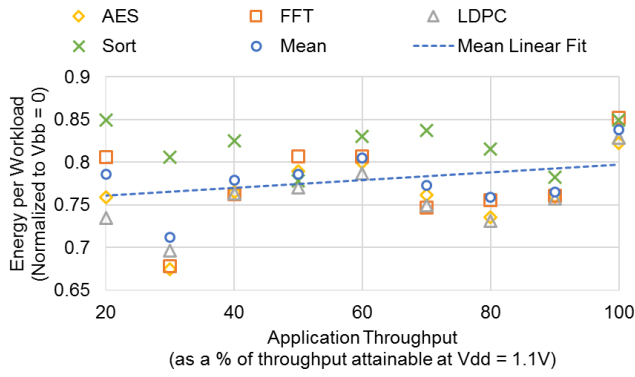


Fig. 4. Reduction in energy usage for the four applications described in Section III.A operating with per-core  $V_{bb}$  selected using Architecture 1 compared to the case with  $V_{bb} = 0$  V. In both cases,  $V_{dd}$  is selected from one of three grids.

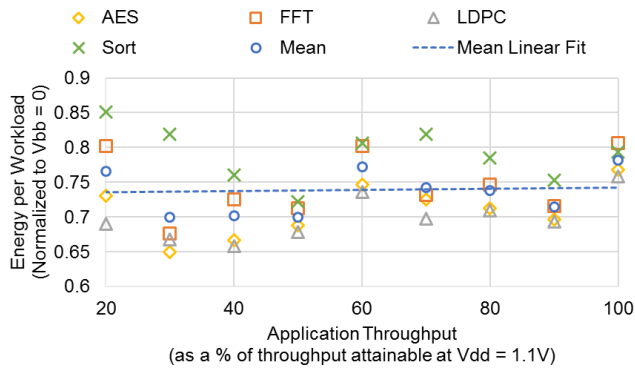


Fig. 5. Reduction in energy usage for the four applications described in Section III.A operating with per-core  $V_{bb}$  selected using Architecture 2 compared to the case with  $V_{bb} = 0$  V. In both cases,  $V_{dd}$  is selected from one of three grids.

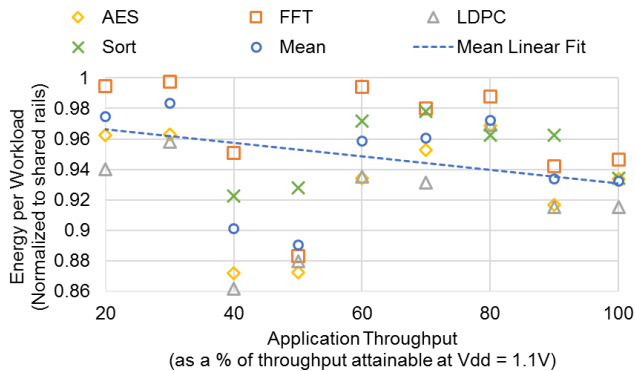


Fig. 6. Reduction in energy usage using Architecture 2 for  $V_{bb}$  voltage selection compared to using Architecture 1.

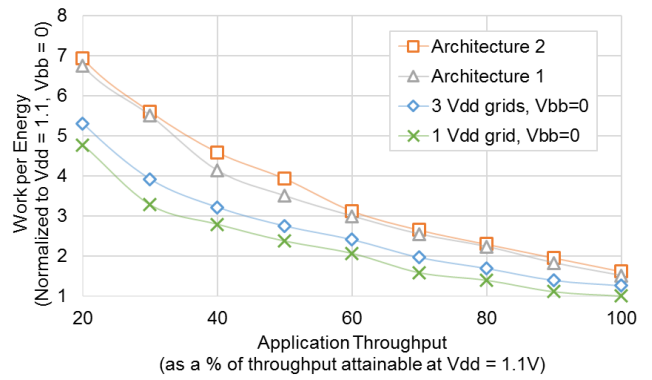


Fig. 7. The geometric mean of the factors of increase in energy efficiency for the four applications described in Section III.A operating with per-core  $V_{dd}$  selection from one of three power grids. Data are compared to the case with fixed  $V_{dd} = 1.1$  V, and zero body bias.

#### IV. FUTURE WORK

While results presented in this work provide valuable insight, further gains in efficiency are expected with the modeling of more complex applications—in particular those with widely-varying activity levels across cores, e.g., workloads containing multiple interconnected applications. In addition, algorithms and circuits that enable the dynamic adjustment of  $V_{bb}$  will provide real-time adaptation to changes in workload throughputs.

#### ACKNOWLEDGMENTS

The authors gratefully acknowledge the support of STMicroelectronics, P. Flatresse, and P. Cogez.

#### REFERENCES

- [1] D. N. Truong et al., "A 167-Processor Computational Platform in 65 nm CMOS" IEEE Journal of Solid-State Circuits, vol. 44, no. 4, pp. 1130–1144, April 2009.
- [2] D. Jacquet et al., "A 3 GHz dual core processor ARM Cortex-A9 in 28 nm UTBB FD-SOI CMOS with ultra-wide voltage range and energy efficiency optimization," IEEE Journal of Solid-State Circuits, vol. 49, no. 4, pp. 812–826, Apr. 2014.
- [3] B. Bohnenstiehl et al., "A 5.8 pJ/Op 115 Billion Ops/sec, to 1.78 Trillion Ops/sec 32nm 1000-Processor Array," Symposium on VLSI Circuits, June 2016.
- [4] B. Bohnenstiehl et al., "KiloCore: A 32-nm 1000-Processor Computational Array," IEEE Journal of Solid-State Circuits, vol. 52, no. 4, pp. 891–902, April 2017.
- [5] Z. Yu and B. Baas, "A Low-Area Multi-Link Interconnect Architecture for GALS Chip Multiprocessors," IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol.18, no.5, pp.750–762, May 2010.
- [6] F. Arnaud et al., "Switching energy efficiency optimization for advanced CPU thanks to UTBB technology," in IEEE Int. Electron Devices Meeting (IEDM) Digest, 2012, pp. 3.2.1–3.2.4.
- [7] B. Bohnenstiehl et al., "KiloCore: A Fine-Grained 1,000-Processor Array for Task-Parallel Applications," IEEE Micro, vol. 37, no. 2, pp. 63–69, March-April 2017.