

# Workload- and Process-Variation Aware Voltage/Frequency Tuning for Energy Efficient Performance Sustainability of NTC Manycores

Ioannis Stamelakos<sup>a</sup>, Sotirios Xydis<sup>b</sup>, Gianluca Palermo<sup>a</sup>, Cristina Silvano<sup>a</sup>

<sup>a</sup>*Politecnico di Milano - Dipartimento di Elettronica, Informazione e Bioingegneria*

<sup>b</sup>*Institute of Communication and Computer Systems - National Technical University of Athens*

---

## Abstract

The power-wall raised by the stagnation of supply voltage in deep-submicron technology nodes, is now the major scaling barrier for moving towards the many-core era. At the same time, the adoption of manycore architectures is considered to be crucial for satisfying the increasing computational power demands and throughput requirements imposed by the explosion in software complexity and volume. The rise of the so-called Dark Silicon, caused by the power budget violations that allow only a small portion of the available computational resources to be simultaneously exploited, points to the direction of energy efficient platforms. Near-Threshold voltage Computing (NTC) has emerged as a promising approach to overcome the manycore power-wall, at the expense of higher sensitivity to process variation and reduced performance which can be compensated with massive parallelization. Given that several application domains operate over specific performance constraints, the performance sustainability is considered a major issue for the wide adoption of NTC. In this work, assuming a feasible, low overhead Power Delivery Network (PDN) for NTC, we investigate how performance guarantees can be ensured when moving towards NTC manycores through a variability-aware voltage and frequency allocation methodology, showing that performance can be efficiently sustained at the NT region while reducing energy dramatically. Additionally, we propose an algorithm for balancing throughput un-

---

*Email addresses:* ioannis.stamelakos@polimi.it (Ioannis Stamelakos),  
sxydis@microlab.ntua.gr (Sotirios Xydis), gianluca.palermo@polimi.it  
(Gianluca Palermo), cristina.silvano@polimi.it (Cristina Silvano)

der process (and workload) variation that sustains performance while providing significant energy savings.

*Keywords:* Near-Threshold Computing, Manycore Architectures, Low Power, Energy Efficiency, Variability

---

## 1. Introduction

Technology scaling has enabled the chip industry to integrate more cores per processor leading to manycore designs. However, the stagnation of voltage scaling, known also as Dennard's law, and its thermal and power budget implications, do not let us exploit the full potential of those architectures. This problem, widely known as Dark Silicon, is expected to be exacerbated in the future. At the same time the demand for computational power is steadily increasing even though it is hidden from the average user through cloud computing and server applications. Manycore processors can provide high throughput for highly parallel workloads, making them ideal candidates for running multi-threaded parallel applications, however, they are constrained by their strict power/thermal budgets.

Near-Threshold Voltage Computing (NTC) [1] represents a promising technique to mitigate the effects of dark silicon in manycore architectures, allowing a large number of cores to operate simultaneously under a given power envelope. In comparison with the conventional Super-Threshold Voltage Computing (STC), computations at the NT region are performed in a very energy efficient manner, unfortunately at the expense of reduced performance and high susceptibility to parametric process variation. Thus, sustaining performance levels delivered in conventional STC is considered a key factor for enabling a wider adoption of NTC in manycores.

Several studies, e.g. [2], [3], have shown that NTC performance loss can be compensated at some level by squeezing all the available parallelism, however they adopt a best-effort strategy to mitigate process variability without providing any performance guarantees. Recently, voltage island management techniques for NTC performance sustainability have been presented [4], [5]. However, the analysis is performed assuming an ideal power delivery model, i.e. power conversion overheads are not taken into account, and only the effects of process variability are compensated. Moreover, parallel applications are exhibiting also workload variability, originated by unbalanced thread execution. In several studies [6], [7], the authors address workload variability by balancing power consumption while

31 maximizing throughput. However, the proposed techniques target the ST region,  
32 neglecting the increased effects of NTC process variability.

33 In this work, we address the problem of performance sustainability when an  
34 application moves from the ST to the NT region. Unlike previous approaches in  
35 NTC [2], [3], [4], [5], we adopt realistic system models accounting for irregulari-  
36 ties and overheads imposed by the Power Delivery Network (PDN). Specifically,  
37 we propose a voltage and frequency tuning strategy that enables energy efficient  
38 and performance sustainable NT manycores through the combined mitigation of  
39 workload- and process-dependent variability. The proposed tuning strategy first  
40 performs frequency upper bounding, guaranteeing performance sustainability in  
41 NTC as well as throughput balancing and then further fine-tunes voltage allo-  
42 cation, mitigating the variability effects. As shown, this technique is extremely  
43 lightweight, i.e. easily integrated in the runtime system and provides increased  
44 NTC energy efficiency.

45 Extensive experimental analysis showed that we can have energy gains rang-  
46 ing from 55% up to 76% while sustaining performance and eliminating variabil-  
47 ity in NTC. Additionally, we can mitigate significantly the throughput imbalance  
48 caused by the increased variability in NTC, obtaining in that way an extra 39%  
49 of energy savings. The rest of the paper is organized as follows: In Section 2  
50 we present the related research in NTC, in Section 3 we describe extensively the  
51 micro-architectural, process variation and power delivery model, in Section 4 we  
52 analyze our methodology for sustaining performance and balancing throughput  
53 in NTC, Section 5 presents the experimental results and finally 6 concludes this  
54 work.

## 55 **2. Relative Research**

### 56 *2.1. NTC Manycore Architectures*

57 The NTC paradigm has emerged as an extremely optimized solution for energy-  
58 efficient systems through low-voltage processing. Early studies [1],[2] have shown  
59 that near-threshold voltage operation provides an optimal design point in respect  
60 to energy and performance efficiency. Early NT voltage processor prototypes [8],  
61 [9] have been recently demonstrated validating the theoretical premises, while  
62 several studies show the high efficiency of NTC for cloud- [10] and server-based  
63 workloads [11].

64 Targeting mainly manycore architectures, NTC imposes several challenges re-  
65 garding the application mapping and the resource/power management due to its in-  
66 creased sensitivity to parametric variation [12]. In [3], Karpuzcu et al. proposed a

67 single-voltage multiple-frequency power management scheme for clustered Near-  
68 Threshold Voltage (NTV) many-cores, aiming mainly at the minimization of the  
69 on-chip power regulation overhead. In [4], the authors are addressing the prob-  
70 lem of sustaining performance when moving to NTC. They explore the efficiency  
71 of voltage island domains and propose a multiple-voltage single-frequency power  
72 management scheme to mitigate within-die variability at NTC voltages. In [5],  
73 they extend power management towards a multiple-voltage multiple-frequency, to  
74 further improve performance. Recently, the use of machine-learning techniques  
75 has been proposed for voltage/frequency allocation in wide-range processors [13],  
76 while in [6] the authors proposed a variability aware thread balancing scheme for  
77 throughput maximization under power and thermal constraints.

## 78 2.2. Power Delivery in NTC

79 Recently, many researchers have investigated the possibility of delivering near-  
80 threshold voltages on-chip: Hsieh et. al [14] designed a linear voltage regulator  
81 (VR) with a variable output voltage that ranges from 0.5 to 1 V in steps of 0.1  
82 V. Lee et al. [15] evaluated the technique known as “voltage stacking” for near-  
83 threshold voltages: multiple low-voltage blocks are powered from a single higher  
84 voltage by “stacking” the logic blocks and recycling charge between the layers.  
85 SuperRange [16] is a wide operational range (0.4 - 1.2 V) power delivery scheme  
86 that uses an off-chip VR to deliver the Super Threshold (ST) voltages and an on-  
87 chip VR for delivering the NT ones with an average of 70% power efficiency.  
88 Booster [17] uses two power supply rails at 400mV and 600mV respectively and  
89 uses hints provided by synchronization libraries to determine which cores should  
90 be “boosted” (run at higher frequency) maintaining an average per-core frequency  
91 and reducing the effect of process variation. In [18], the authors show how the  
92 workload characteristics of different applications can be analyzed at design-time  
93 and exploited at runtime in order to obtain optimal voltage allocation in NTC,  
94 while using standard, low-cost power delivery architectures. Interesting power  
95 delivery (PD) architectures with per-core DVFS support have been proposed as  
96 well even though they target multicore platforms that do not operate in the NT re-  
97 gion. VRCon [19] uses a switching network that allows the cores to share on-chip  
98 voltage regulators according to the DVFS intervals and load conditions, improv-  
99 ing the overall energy efficiency, but it is not highly scalable since the integration  
100 of hundreds of switching (SW) regulators and switches on-chip would impose  
101 a significant area overhead. Sinkar et al. [20], based on the observation that a  
102 Low-Dropout Regulator (LDO) can share its largest component with the Per-Core  
103 Power-Gating (PCPG) device and therefore can be integrated with a minimum

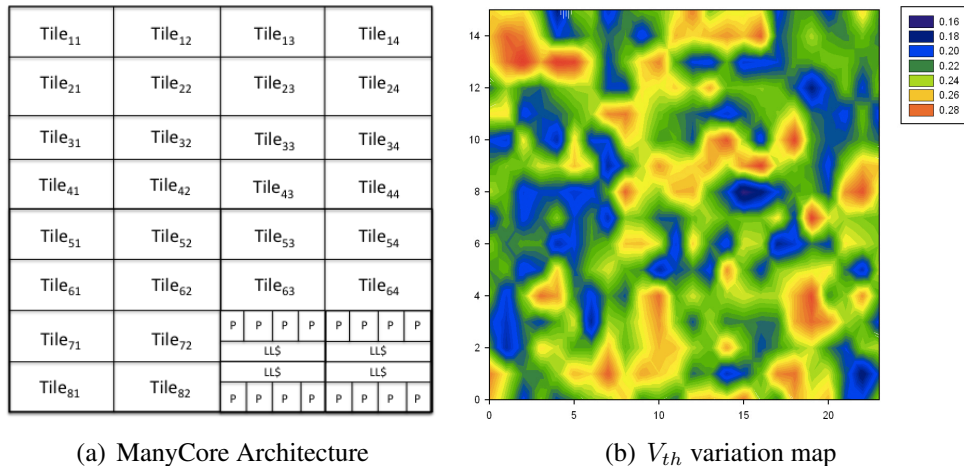


Figure 1: Tile-based architecture (a) and  $V_{th}$  variation map derived by [22] (b).

104 area overhead, demonstrated that a multicore platform with per-core voltage do-  
 105 mains is implementable, and under certain assumptions, can be as effective as  
 106 one having per core switching VRs (instead of LDOs). Godycki [14] et al., by  
 107 deploying on-chip switched-capacitor regulators, design a reconfigurable power  
 108 distribution network, offering a wide range of voltages, enabling in that way fast,  
 109 fine-grained, per-core voltage scaling. All of those techniques address part of  
 110 the problem but they do not offer a clear solution for a per-core, near-threshold  
 111 voltage power delivery architecture that could be applied to a manycore platform.

### 112 3. System Level Model of NTC Manycore

113 **Micro-architectural model:** We focus our study on tile-based architectures,  
 114 including the ones proposed in [1], [3] and [4]. Figure 1(a) shows an abstract view  
 115 of the tile-based manycore architecture, as well as the different intra-tile organi-  
 116 zations. The intra-tile architecture can be customized by varying the number of  
 117 cores per tile and the memory configuration of the last level cache (LLC). In this  
 118 paper, each core owns a private instruction and data cache (P\$). The LLC (LL\$)  
 119 is shared among 4 cores, composing a tile. The Intel Nehalem [21] processor  
 120 configuration for the core and the P\$ has been adopted.

#### 121 3.1. Process Variation Model

122 In order to capture the process variation at NT operation, we integrate the  
 123 Varius-NTV [22] microarchitectural model within the proposed framework. While

124 Varius-NTV reuses the spherical distance function in [23] for modeling the intra-  
 125 die spatial correlations, it heavily extends it by updating the STC micro-architectural  
 126 delay and SRAM cell models to reflect in a more accurate manner the higher sen-  
 127 sitivity of NTC to process variation. Specifically, it (i) calculates gate-delay ac-  
 128 cording to the EKV model [24], (ii) it incorporates a 8T SRAM cell model for  
 129 reliable read/write operations in NTC and (iii) it considers a larger set of memory  
 130 timing and stability failure modes. We used the ArchFP [25] tool to automatically  
 131 generate the floorplan of the target manycore architectures. Based on the provided  
 132 manycore floorplan, Varius-NTV generates the corresponding variation maps ac-  
 133 counting for the within-die (WID) and die-to-die (D2D) process variations. Figure  
 134 1(b) shows a sample instance of its  $V_{th}$  variation map.

135 Assuming  $B$  as the set of component blocks in the floorplan and  $D$  the set of  
 136 dies, we define  $V_{th}^{(i,j)}, L_{eff}^{(i,j)}, i \in B, j \in D$  that corresponds to the  $V_{th}$  and  $L_{eff}$   
 137 respectively of the architecture's component  $i$  in sample die  $j$ . Once extracted,  
 138 the two parameters are used for allocating to each component the lowest possible  
 139  $V_{dd}^{(i,j)}$ , for sustaining the  $f_{NTC}$  frequency constraint, as described in Section 4.1.

### 140 3.2. Power Delivery Model

141 In most of the research conducted around NTC, the power delivery architec-  
 142 ture is neglected and it is generally assumed that the Near-threshold voltages can  
 143 be delivered without any losses. However, the PDN can have a large overhead on  
 144 the power consumption if not designed and customized for the specific platform.

145 Generally, the power delivery network can be divided into two components:

- 146 1. Off-chip network: one or more power supply rails, powered by off-chip  
 147 voltage regulators, deliver the appropriate voltages to the chip.
- 148 2. On-chip network: a second layer, connected with the off-chip network, con-  
 149 sisting of voltage regulators that step down the voltage and deliver it to the  
 150 cores. The most common types of VRs are:
  - 151 • Switching Regulators: They have a very good efficiency (~90%) but  
 152 they consume a lot of area and it is hard to integrate them on chip.
  - 153 • Low Dropout Regulators: An LDO is kind of linear regulator and its  
 154 efficiency is calculated as follows:

$$\eta_{ldo} = \frac{P_{out}}{P_{in}} = \frac{I_{out}V_{out}}{I_{in}V_{in}} = \frac{I_{out}V_{out}}{(I_{out} + I_q)V_{in}} \quad (1)$$

155 where  $I_q$  is the quiescent current flowing to the ground:

$$I_{in} = I_{out} + I_q \quad (2)$$

156 In [20], the authors, based on the observation that Per-Core Power-Gating  
 157 (PCPG) devices augmented with feedback control circuitry can serve as low-cost  
 158 Low-Dropout Regulators (LDOs), show that power efficiency as high as that de-  
 159 livered by on-chip switching Voltage Regulators (VRs) can be achieved. The  
 160 per-core area overhead is only  $\sim 2\%$ , making the implementation of this design  
 161 affordable even for a manycore platform. An LDO, as mentioned above is a type  
 162 of linear regulator, and its efficiency is defined in Equation 1.

163 The LDOs exhibit current efficiencies up to 99%, making the difference be-  
 164 tween the input and the output voltage the dominant factor of Equation 1. The  
 165 lower the difference, the higher the efficiency, however the difference should not  
 166 drop below a certain threshold because the circuit ceases to regulate its output  
 167 voltage against any fluctuations in the input voltage. This threshold is defined as  
 168 the *dropout voltage* and is defined at design time. The main advantages of an LDO  
 169 are its fast transient response and its low cost, making it a potential candidate for  
 170 on-chip VR.

171 Although typical voltage regulation at low  $V_{dd}$  levels exhibit large conversion  
 172 inefficiencies [26], the authors in [27] have shown that designing and implement-  
 173 ing an LDO for near- threshold logic circuits can be performed with quite high  
 174 efficiency, i.e. around 99%, and fast transition times, i.e. hundreds of nanosec-  
 175 onds up to a few microseconds. In this work, we utilise the LDO configuration  
 176 adopted directly from [27] with the following design characteristics: 0.5V input  
 177 voltage, 0.34V - 0.45V output voltage range, 0.05V dropout voltage and 98.7%  
 178 current efficiency. The SRAM voltage cannot be scaled down after a certain point  
 179 because it leads to major defects and errors and we assume it has its separate  
 180 constant supply : 500mV in standby mode and 700mV when accessing.

#### 181 **4. Workload and Process Variability Management for Performance Sustain-** 182 **ability in NTC**

183 The effects of process variation are exacerbated in NTC, but except for that, in  
 184 order to exploit its energy efficiency potential, we should be able to provide per-  
 185 formance guarantees to the applications running in an NTC manycore platform,  
 186 sustaining their ST performance in the ideal case. This becomes more evident  
 187 if we consider the emerging paradigm of data centers and cloud computing. To  
 188 further motivate the aforementioned claim, Figure 2 shows the performance dis-  
 189 tribution for a 128-core NTC manycore that implements the best-effort EnergyS-

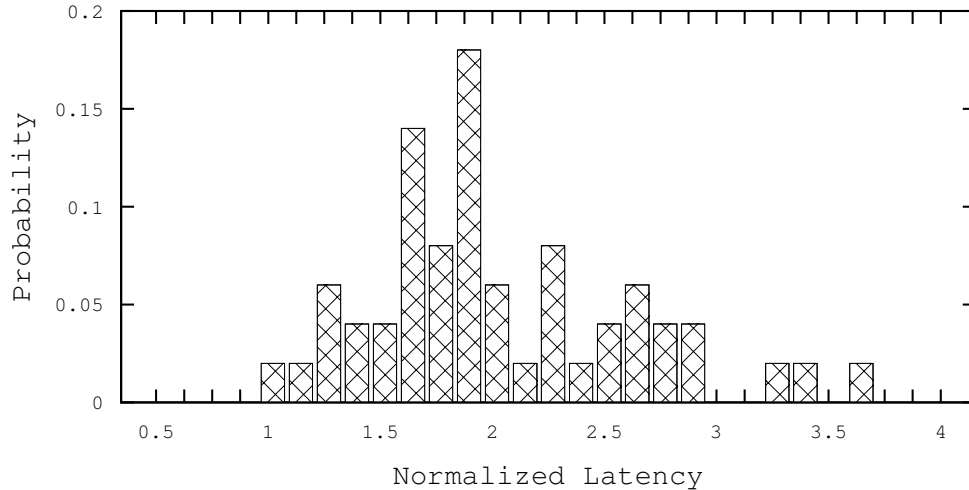


Figure 2: Performance distribution on a 128-core NTC manycore implementing the EnergySmart [3] approach.

190 mart power management Single-Voltage/Multiple-Frequencies (SVMF) approach  
 191 [3]. The results are obtained for the execution of the *Barnes* application over 100  
 192 different variation maps. The normalized performance value of 1 corresponds to  
 193 the *nominal* performance of the application. As shown, the performance of NTC  
 194 manycore platforms is not controllable and is spread out over a wide range of  
 195 normalized values (from 1 to 3.7) due to the underlying workload and process  
 196 variability. Thus, the adoption of NTC for applications, exhibiting specific per-  
 197 formance and/or throughput constraints, requires careful selection and tuning of  
 198 the power management scheme. In the following sections, we propose a fast run-  
 199 time voltage and frequency allocation strategy accounting for both workload- and  
 200 process- variability, enabling performance sustainability in NTC manycores. The  
 201 proposed strategy consists of two steps: the first step performs a performance sus-  
 202 tainability analysis for generating frequency and voltage upper bounds for NTC  
 203 operation, while the second step implements the voltage and frequency fine-tuning  
 204 that enables throughput balancing in NTC, thus mitigating performance variability  
 205 effects.

#### 206 4.1. Performance Sustainability in NTC

207 So far, application workloads have been originally developed and character-  
 208 ized for the ST region. In order to sustain ST performance figures (i.e. latency



209 or throughput) when moving to the NT region, the inherent parallelism of the ap-  
 210 plications should be exploited [28] to alleviate the impact of the reduced clock  
 211 frequencies in NTC. Assuming a minimum allowed latency  $L_{min}$  and maximum  
 212 core count constraint,  $C_{max}$  for the NTC manycore, we first calculate the clock  
 213 frequency of the platform at near-threshold,  $f_{NTC}$ , that satisfies the performance  
 214 constraint. Let  $L_{C_{max}}$  be the performance, in terms of latency, of an STC many-  
 215 core architecture with  $C_{max}$  number of cores, running at  $f_{STC}$ . When running in  
 216 STC,  $L_{min} - L_{C_{max}} > 0$  is the available latency slack due to the higher degree  
 217 of parallelism of the architecture, that can be exploited to run the application at  
 218 lower frequency. Utilizing this positive slack, the  $f_{NTC}$  is calculated as follows:

$$f_{NTC} = \frac{L_{C_{max}}}{L_{min}} \times f_{STC} \quad (3)$$

219 The calculated  $f_{NTC}$  refers to the target clock frequency of each core in NTC  
 220 for sustaining ST performance, without considering the spatial effects of process  
 221 variations. Next, we extract the values of  $V_{th}$  and  $L_{eff}$  according to the proce-  
 222 dure described in Section 3.1. Finally, the two parameters are used for allocating  
 223 to each component the lowest possible (provided by the power delivery system)  
 224  $V_{dd}^{(i,j)}$ , for sustaining the  $f_{NTC}$  frequency constraint given that:

$$f_{NTC} \propto \frac{(V_{dd}^{(i)} - V_{th}^{(i)})^\beta}{V_{dd}^{(i)} \times L_{eff}^{(i)}} \quad (4)$$

225 where  $\beta$  is a technology-dependent constant ( $\approx 1.5$ ). The extraction of the  $f_{NTC}$   
 226 and the per component  $V_{dd}^{(i,j)}$ , enables the adoption of a fine grained power man-  
 227 agement scheme for NT operation with guaranteed performance sustainability.

#### 228 4.2. Throughput Balancing in NTC

229 After presenting our methodology for sustaining performance in NTC, we  
 230 move on to dealing with a different but relevant subject: the throughput imbalance  
 231 observed in the threads of the same application when running in NTC. The per-  
 232 core/thread throughput imbalance that can be caused is usually neglected, how-  
 233 ever, this cannot be overlooked in NTC, because the underlying variability, which  
 234 is also exacerbated in each subsequent technology node, has a great impact on  
 235 both performance and power consumption. In [29], a 30% deviation in frequency  
 236 among the cores of an Intel 80-core manycore platform has been observed in real,  
 237 on-silicon measurements. The results are referring to a 65nm technology and a  
 238 much higher voltage than a near-threshold one, thus the variation is expected to

239 be much worse at the NT region. In multithreaded parallel applications, there  
 240 even exists an IPC imbalance caused by the access latency irregularities and the  
 241 per-thread data workload variation leading to different memory access patterns.  
 242 Given the experimental setup described in Section 5.1, the per-core, normalized  
 243 to the minimum value, throughput is shown in Figure 3. The IPCs are obtained  
 244 from simulations and no process variation is considered, i.e. all the cores run at the  
 245 same frequency. The standard deviation ( $\sigma$ ) is around 11%, but this can be wors-  
 246 ened by the actual maximum per-core frequencies caused by process variation.  
 247 The per-core throughput (normalized to the minimum value), when variability is  
 248 taken into account, is depicted in Figure 4 for the Cholesky benchmark run on 64  
 249 cores. Not only the core exhibiting the least throughput is now a different one but  
 250 also the variation in per-core throughput has increased dramatically, leading to a  
 standard deviation of 38%.

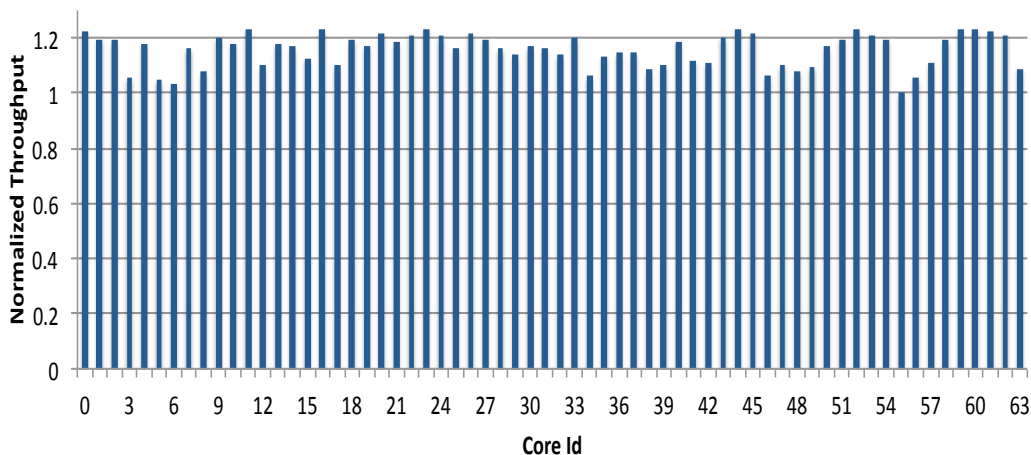


Figure 3: Per-core throughput: imbalance due to IPC variability.

251  
 252 The irregular behavior described above leads us to the conclusion that addi-  
 253 tional power is wasted from the faster cores even though not needed since they  
 254 will finish sooner and wait at the synchronization barrier. This observation, cre-  
 255 ates both a problem that has to be taken care of and an opportunity that should be  
 256 exploited. Our proposed algorithm manages to mitigate the throughput imbalance  
 257 and provides up to 43.5% of power savings, on average, for the reported results.  
 258 Additionally, we are able to identify for each benchmark which is the resource  
 259 allocation that maximizes energy efficiency, a decision that depends heavily on  
 260 their unique workload characteristics.

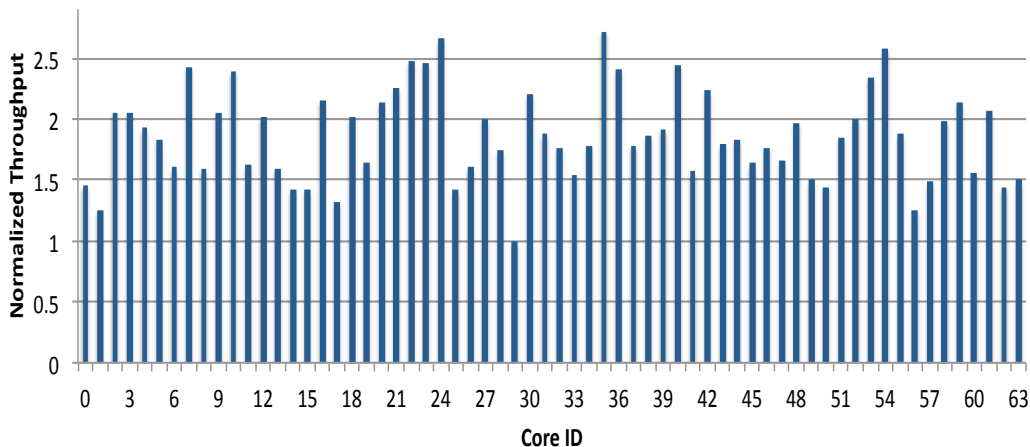


Figure 4: Per-core throughput: imbalance due to process variability.

261 For this reason we designed an algorithm for dealing with this discrepancy  
 262 (Algorithm 1). The IPC of each core is provided from our simulations in Sniper  
 263 and can be obtained at runtime by the performance counters provided by the pro-  
 264 cessor. The relationship between frequency and the transistor characteristics is  
 265 depicted in Equation 4, where  $\beta$  is a technology-dependent constant ( $\approx 1.5$ ),  $V_{th}$   
 266 and  $L_{eff}$  is the transistors' threshold voltage and the effective length respectively.

267 As shown in Figure 4, the throughput achieved among the different cores dif-  
 268 fers significantly.

269 Initially, we calculate the maximum per-core throughput that can be achieved  
 270 assuming that the  $V_{dd}^i$  and  $f_{NTC}$ , extracted from the previous step, is allocated to  
 271 them. The throughput of the  $i_{th}$  core is calculated as in Equation 5.

$$Throughput_i = IPC_i \cdot F_i \quad (5)$$

272 where  $IPC_i$  and  $F_i$  are the Instructions Per Cycle (IPC) and Frequency of the  $i_{th}$   
 273 core respectively.

274 In the pseudocode, this first phase is shown in lines 1-4: given the LDO volt-  
 275 age level that satisfies the  $V_{dd}^i$  and the  $f_{NTC}$  are assigned to the core (line 2-3) and  
 276 the equivalent throughput is calculated (line 4) for each core. The core exhibiting  
 277 the minimum throughput is now found and becomes the reference value and bot-  
 278 tleneck (line 6). For the rest of the cores, the minimum frequency (and hence volt-  
 279 age) for decreasing their throughput as close as possible to the minimum (found in  
 280 the previous step), is calculated using Equations 5 and 4 respectively (lines 8-9).

---

**Algorithm 1:** Algorithm for Balancing Throughput

---

**Input** : Per-core *IPC* ( $IPC_i$ ), *LDO* voltage levels ( $LDO_{out}^j$ ), Variability modeling parameters ( $V_{th}^i, L_{eff}^i$ )

**Output:** Per-core *Voltage* and *Frequency* allocation for balancing throughput

```
1 foreach  $i \in Cores$  do
2    $Vdd_i \leftarrow LDO_{out}(LDO, V_{dd}^i)$ 
3    $F_i \leftarrow f_{NTC}$ 
4    $T_i \leftarrow F_i \cdot IPC_i$ 
5 end
6  $T_{min} \leftarrow \min(T_i)$ 
7 foreach  $i \in Cores$  do
8    $F_i \leftarrow \frac{T_{min}}{IPC_i}$ 
9    $Vdd_i \leftarrow f^{-1}(F_i, V_{th}^i, L_{eff}^i)$ 
10 end
```

---

281 Assuming that the LDOs have discrete output voltage levels, we select the lowest  
282 level possible that is providing a higher value than the ideal  $V_{dd}$  calculated, ensur-  
283 ing in that way that the desired frequency can be reached. For our experiments  
284 we assumed that the LDOs have a 12.5 mV output voltage resolution, meaning  
285 that they can provide the desired voltage in steps of the aforementioned voltage  
286 value. As mentioned in the previous subsection (3.2), the LDOs provide an output  
287 voltage range of 100mV (0.35V - 0.45V), offering, as a consequence, 8 discrete  
288 voltage levels.

289 The algorithm described above sustains and guarantees the throughput exhib-  
290 ited from each application (in the sense that its goal is not to maximize it like  
291 in other research attempts), while minimizing the overall power consumption. It  
292 introduces a low overhead and can be run at runtime while it scales linearly with  
293 the number of cores, i.e.  $\mathcal{O}(n)$  complexity. Additionally, each per-core compu-  
294 tation (throughput calculation and voltage/frequency allocation) is orthogonal to  
295 each other and can be computed fully in parallel. In this work, it is applied to ho-  
296 mogeneous and symmetrical workloads but it can be applied to unbalanced ones  
297 considering each time the target platform (e.g. cloud, server etc.), the application  
298 scenario (e..g single/multiple applications/instances, number of threads etc.) and  
299 the designer’s goals, but this is left as future work.

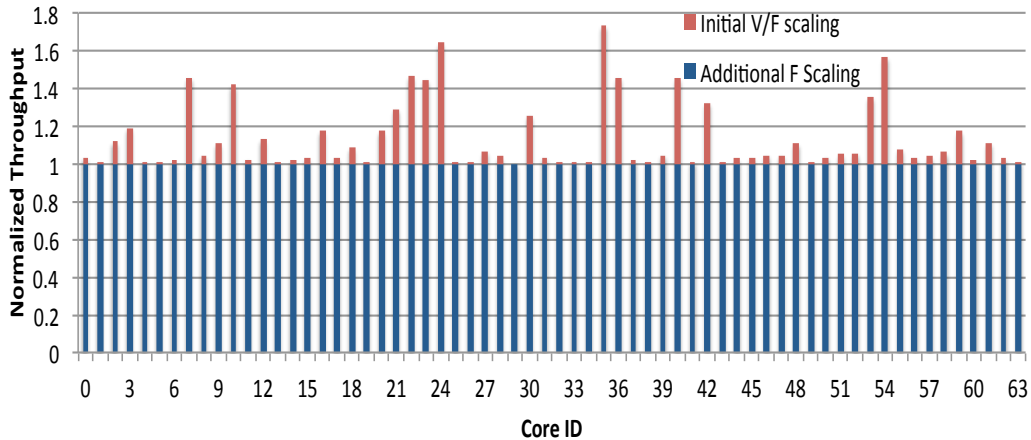


Figure 5: Per-core throughput: balancing obtained by the proposed algorithm.

## 300 5. Experimental Results

### 301 5.1. Experimental Setup

302 We use the Sniper manycore simulator [30] for obtaining the performance  
 303 counters and McPAT [31] for the initial power estimation. Since McPat is not  
 304 validated against near-threshold voltages, we use it for obtaining the ST values  
 305 and then we scale accordingly in order to calculate the near threshold ones, us-  
 306 ing the models for NTC provided in [24]. In order to characterize the process  
 307 variation at the NT region we deployed the Varius-NTV microarchitectural model  
 308 [22]. The applications simulated are taken from the Splash-2 [32] and Parsec  
 309 [33] benchmark suites, which provide high performance parallel workloads. Each  
 310 application has different characteristics and exhibits a different behavior, i.e. per-  
 311 formance, depending on the number of allocated resources, i.e. cores. The adopted  
 312 on-chip voltage regulation scheme [27] enables voltage scaling response times on  
 313 the order of hundreds of nanoseconds up to a few microseconds. The proposed  
 314 voltage/frequency tuning framework is applied over epochs of 1ms, thus even in  
 315 the case that the throughput balancing algorithm is used and triggered constantly  
 316 due to regular workload imbalances, the overhead of the transition is negligible,  
 317 i.e. several orders of magnitude lower than the aforementioned epoch.

318 Since we are simulating homogeneous multi-threaded applications, we con-  
 319 sider a discrete number of configurations with the following number of cores:  
 320 4, 8, 16, 24, 32, 64. A summary of the experimental setup used to evaluate the

Table 1: Experimental Setup: Platform Parameters

Parameters	Value
Process Technology	22nm
STC Supply Voltage	1.05V
NTC Supply Voltage	0.45V
Nominal $V_{th}/\sigma_{V_{th}}$	0.23V/0.025
Number of Cores/Core Area	64/6mm <sup>2</sup>
Tile/	4cores
Private Cache Size/Area	320KB/4.14mm <sup>2</sup>
Last Level Cache Size – Area	8 MB / 15.52mm <sup>2</sup>

321 methodology is presented in Table 1. Intel’s Nehalem architecture (22nm technol-  
 322 ogy node) has been adopted for determining the core and cache types and sizes.  
 323 The maximum number of cores is 64, divided in tiles consisting of 4 cores and a  
 324 shared last level cache (LL\$) of 8MB and each core owns a private instruction and  
 325 data cache (P\$).

### 326 5.2. Energy Savings: NTC vs STC

327 Figure 6 shows the energy savings obtained when going from 4 cores at ST  
 328 voltage to 64 cores at NT voltage for each benchmark, assuming the power de-  
 329 livery system described in 3.2. All the applications are benefited significantly by  
 330 running in NTC mode, i.e the minimum energy reduction is 55%. However, there  
 331 is a large variation in the saving between different applications. This is caused  
 332 from the distinct workload characteristics and the inherent parallelism of each  
 333 benchmark, i.e. the performance scalability w.r.t the increase of the number of  
 334 cores (4 to 64), since it impacts the minimum frequency to be sustained and thus  
 335 the minimum  $V_{dd}$  to be deployed. As depicted in Figure 6, a maximum of 76%  
 336 decrease in energy is observed for the Raytrace application, whereas the minimum  
 337 gain is around 55% for the Cholesky application. The average savings are 66%.

### 338 5.3. Throughput Balancing

339 Figure 7, depicts the per-core throughput, normalized to the minimum value,  
 340 after running the algorithm proposed in the previous section. In order to evaluate  
 341 its efficiency we define a measure similar to standard deviation ( $\sigma$ ) by substitut-  
 342 ing the mean with the minimum. Its purpose is to measure the distance of the  
 343 throughput values from the minimum for each core. This is because our goal is

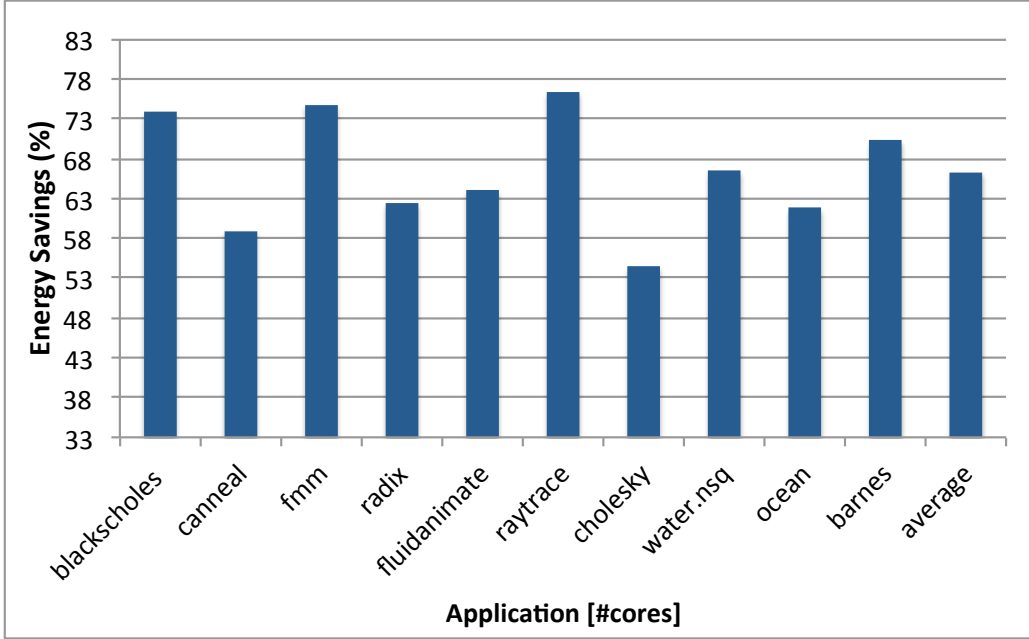


Figure 6: Energy savings: 4-core STC chip versus 64-core NTC

344 to balance per-core throughput as much as possible w.r.t. the core with the lowest  
 345 throughput. We define the following measure:

$$\zeta = \sqrt{\sum_i \frac{(T_i - T_{min})^2}{\#cores}} \quad (6)$$

346 Whereas the throughput  $\zeta$ -deviation, as defined in Equation 6, for the 16-  
 347 threaded Cholesky application is 77.8% before balancing, it drops to 10.5% after  
 348 applying our algorithm. For 64 cores, the same measure is increased to 95.8%  
 349 before balancing but it is reduced down to 30% after applying our technique (69%  
 350 reduction). As we can see in Figure 7, most of the cores exhibit a similar through-  
 351 put, except for some of them, represented by the peaks in red color in the graph,  
 352 e.g. 24, 35, 54, that are much faster and with the specific voltage range provided  
 353 by the regulators it is impossible to reduce their throughput any further. For each  
 354 application the  $\zeta$ -deviation is reduced by 70% on average. The previous results  
 355 are obtained assuming that, given the allocated voltage, the maximum achievable  
 356 frequency is assigned to each core. As mentioned before, throughput in Figure 7  
 357 is normalized w.r.t. minimum value, therefore, the faster cores can exhibit a quite  
 358 higher throughput (represented as red colored peaks in the graph) due to the higher

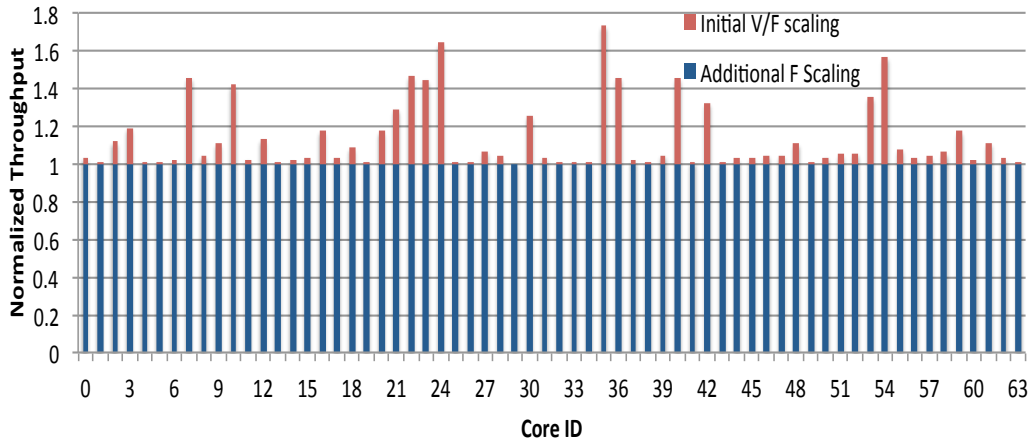


Figure 7: Per-core throughput: balancing obtained by the proposed algorithm.

359 frequency that can be assigned to them. As an extra, fine tuning optimization, we  
 360 can slow down those cores, by reducing the frequency in order to perfectly bal-  
 361 ance the throughput and eliminate the variation completely, obtaining some extra  
 362 power savings (reported in section 5.5).

363 The technique proposed, may look general and not specifically customized for  
 364 NTC. However, the important observation here is that 1) the effect of variability  
 365 is greatly increased in NTC causing a much greater imbalance than in STC and 2)  
 366 the number of resources (cores) is much larger in NTC, creating the potential of  
 367 greater power/energy savings. Specifically, when assuming a 4-core platform in  
 368 STC the average savings for all applications is just just 16%, whereas even for a  
 369 16-core platform the average savings don't exceed 19%.

#### 370 5.4. Quantitative comparison with state-of-art

371 Having described our methodologies for sustaining performance and balanc-  
 372 ing throughput in NTC, as well as the results obtained in each case, we take a step  
 373 further to see how they are connected and how they compare to each other.

374 If we consider a single voltage scenario (similar to the approach presented in  
 375 [3]), then the only degree of freedom we have is frequency scaling. This is also  
 376 mentioned as Single Voltage/Multiple Frequencies (SVMF) approach. Assum-  
 377 ing a performance sustainability scenario, we compare this approach to the one  
 378 presented in Section 4. We also consider that the voltage provided to the SVMF  
 379 platform is the maximum one delivered by the power delivery architecture used in



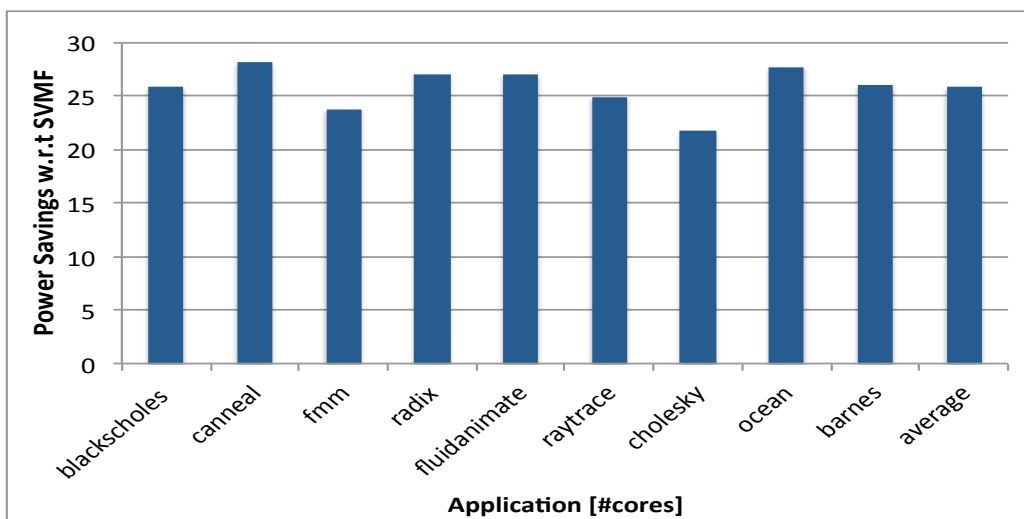


Figure 8: Power savings obtained when compared with the SVMF approach.

380 this work. Figure 8 depicts the power savings we get when we compare the 2 ap-  
 381 proaches. The savings we get are more or less similar for all applications with an  
 382 average value of 25%. This is because, even though the frequency can be scaled  
 383 down to the one needed for sustaining the performance requirement, the voltage  
 384 cannot be scaled accordingly, wasting extra power when not needed.

385 Next, we move on to comparing 3 different approaches as well as evaluat-  
 386 ing the impact of throughput balancing discussed in the previous section. First,  
 387 we evaluate a power management technique similar to the one described above  
 388 (SVMF) [3]: This time we allow the voltage to be adjustable, but after calculat-  
 389 ing (according to Section 4) the frequency that satisfies the required con-  
 390 straint, the minimum voltage provided by the power delivery system that guar-  
 391 antees the aforementioned frequency, while taking into account process varia-  
 392 tion, is allocated to the whole platform, i.e. the slowest core is determining the  
 393 voltage. Then a Multiple-Voltages/Multiple-Frequencies (MVMF) power man-  
 394 agement technique is deployed, taking advantage of the multiple voltage levels  
 395 and the fine-grained, per-core voltage scaling to reduce the per-core variability  
 396 and maximum per-core frequency. Last, the methodology presented in this pa-  
 397 per, Multiple-Voltages/Single-Frequency (MVSF) combined with the throughput  
 398 balancing heuristic is used. In that way we can sustain performance, mitigate  
 399 variability and reduce energy consumption, showing that this technique is very  
 400 efficient when dealing with this kind of applications.

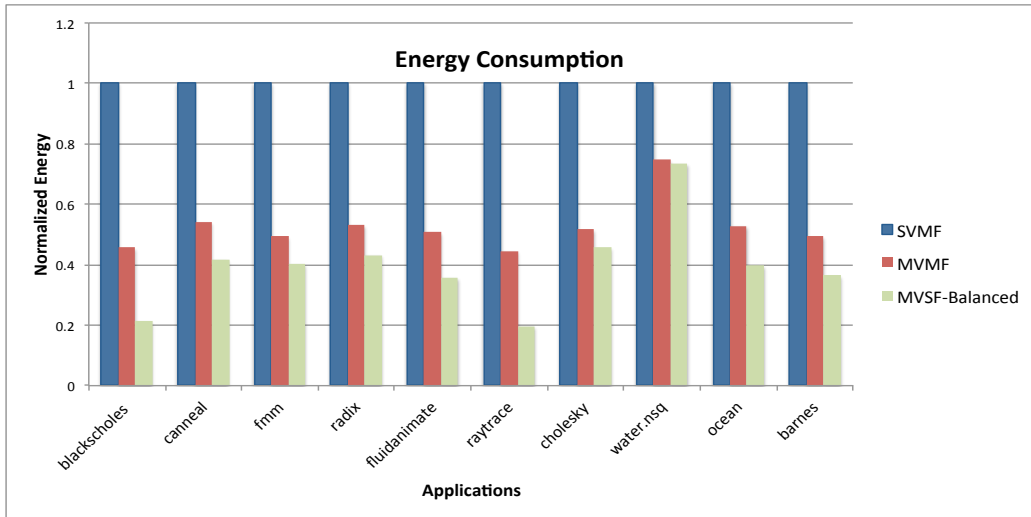


Figure 9: Energy consumption comparison of the SVMF, MVMF and MVSF (w.throughput-balancing) techniques.

401 The results are depicted in Figure 9. As we can see, providing an adjustable  
 402 per-core voltage decreases dramatically the energy consumption, 50% on average  
 403 (ranging from 35% up to 55% depending on the scalability of each application.  
 404 However, when throughput balancing is taken into account, we can reduce en-  
 405 ergy even further by just using one per-chip frequency, i.e. the one that meets the  
 406 performance constraint. As we can see, the results vary depending again on the  
 407 workload characteristics. For example, *Blackscholes*, which is an embarrassingly  
 408 parallel application, can reduce the energy consumption by 50%, whereas *Wa-*  
 409 *ter.nsq* does not have essentially any benefit. On average, the technique delivers  
 410 an extra energy reduction of 26%.

411 Finally, we take the previous analysis one step further and we compare the  
 412 SVMF with the MVSF approach after balancing throughput while running at the  
 413 minimum frequency for sustaining the ST performance. Additionally, we compare  
 414 these two techniques with a best effort throughput balancing technique: we try  
 415 to run our applications as fast as the resources provided can allow us, relaxing  
 416 the performance constraint, but still applying the balancing throughput algorithm  
 417 described in Section 4.2

418 Figure 10 shows that the MVSF technique provides better results compared  
 419 to the SVMF one. This can be easily explained by the fact that while in the first  
 420 case after calculating the platform frequency, we can tune in a fine-grained, per-

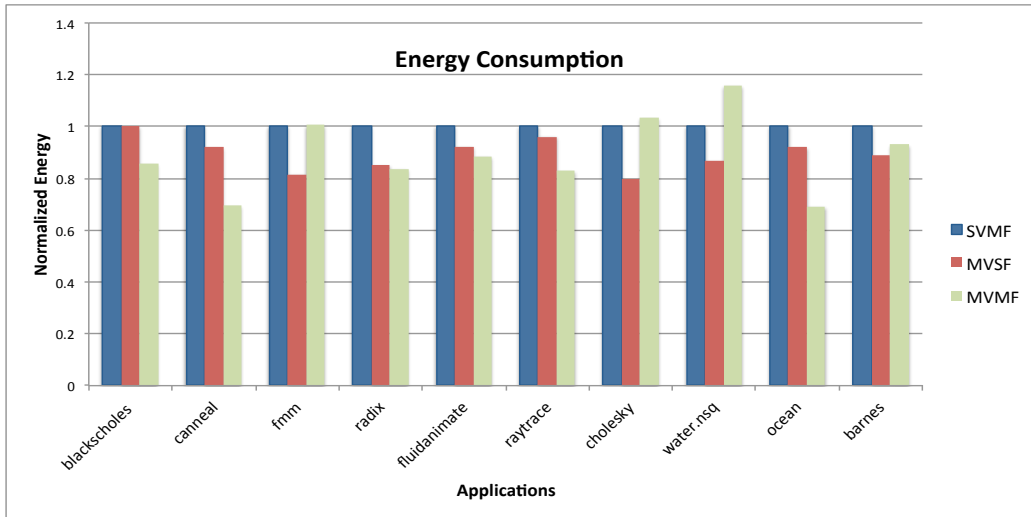


Figure 10: Energy consumption comparison of the different techniques.

421 core manner the voltage, allowing the faster cores to operate under a lower  $V_{dd}$ , in  
 422 the SVMF approach the common voltage (which is the one enabling the slowest  
 423 core to run in the required frequency for sustaining performance), increases the  
 424 power consumption and eliminates the potential voltage scaling-induced power  
 425 savings delivered from the faster cores. As we can observe, it provides always  
 426 better results in terms of energy consumption, 35% on average. The least benefit  
 427 is observed for the *Water.nsq* application, which scales poorly with the number of  
 428 cores, not allowing in that way great savings because of the increased frequency  
 429 that is needed in order to sustain its ST performance.

430 The results regarding the third technique (MVMF) are also interesting. The  
 431 throughput balancing technique, can provide significant energy reduction, up to  
 432 45% and 17% compared to the SVMF and MVSF techniques (*Canneal*, *Ocean*),  
 433 but not in all cases. For example, the energy consumed in the *Water.nsq* bench-  
 434 mark is 15% higher than in the SVMF approach and in some other benchmarks  
 435 like *Radix* and *Barnes* we have similar energy consumption with the MVSF ap-  
 436 proach. This initially strange behavior can be explained by the fact that in this  
 437 scenario we deploy a best effort approach, not considering performance sustain-  
 438 ability, but going even faster when possible. This can be an interesting trade-off,  
 439 especially when we have an optimization problem, like maximizing performance  
 440 given an energy constraint. In the case of *Water.nsq*, the speedup is not big enough,  
 441 so we end up consuming more energy. The speedup of the various applications is

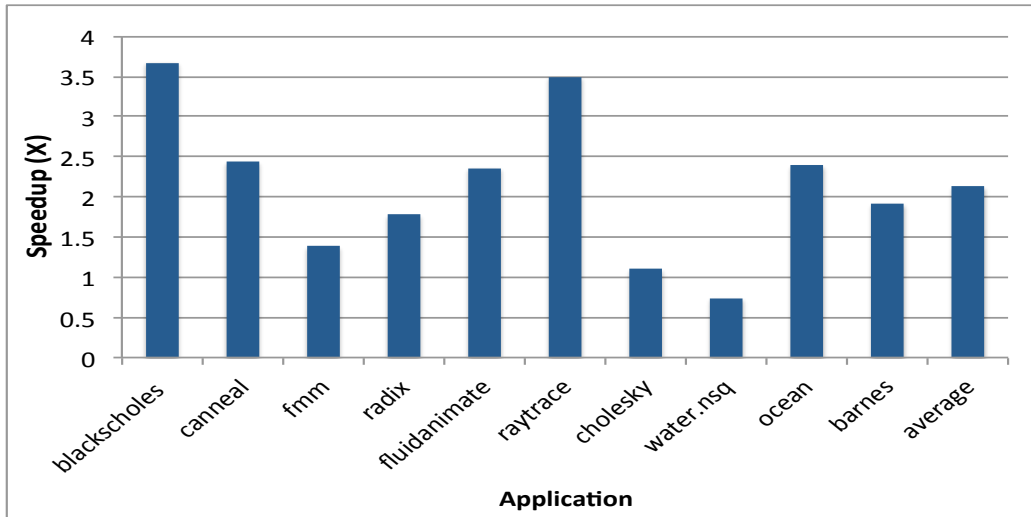


Figure 11: Speedup of the MVMF technique

442 shown in Figure 11. Some benchmarks (*Blackscholes*, *raytrace*) get a significant  
 443 speedup (3x), while the average value is 2x. The most important point, that we  
 444 need to stress out is, that by using the throughput balancing algorithm proposed in  
 445 this work, we can have energy savings even while improving performance com-  
 446 pared to the first two techniques (SVMF, MVSF).

#### 447 5.5. Impact of Resource Allocation on NTC Energy Efficiency

448 Till now, we considered that granted a platform that consists of a specific num-  
 449 ber of cores (in our case 64), we allocate all of the resources each time when we  
 450 want to run in NTC. However, this may not be the case if we just want to run in  
 451 the most energy efficient configuration and we do not care about meeting a spe-  
 452 cific performance constraint, e.g. STC performance sustainability. In this section,  
 453 for each application, we characterize each configuration in terms of energy effi-  
 454 ciency and we find the most efficient one. Energy efficiency, as shown in Equation  
 455 7, is defined as Throughput over Power and is measured in Instructions/Joule or  
 456 MIPS/Watt equivalently. After running our algorithm, we can determine the most  
 457 energy efficient configuration for each application.

$$Energy\ Efficiency = \frac{Throughput}{Power} \quad (7)$$

458 Table 2 summarizes the results. Applications have different optimal configu-  
 459 ration points (#cores) depending on their workload characteristics, i.e. after a cer-

tain point there is no sense in allocating more resources because the performance-throughput is not increased significantly in order to compensate for the increased power consumed by the extra cores. For each optimal configuration, we calculate the power savings achieved by applying our optimization algorithm. The power savings are quite high, 39% on average, (varying from 34% up to 44%) and they become even more significant when considering that the overall performance was not impacted. As mentioned in the previous subsection, by performing an extra frequency scaling for the faster cores, we can eliminate the throughput variation completely, increasing ideally the average power savings from 39% to 43.5% (varying from 37% up to 51%).

Table 2: Energy Efficient Optimal Configuration

App	#cores	App	#cores
blackscholes	48	canneal	48
fmm	32	ocean	32
fluidanimate	24	raytrace	24
streamcluster	48	cholesky	64
swaptions	8	water.nsq	8
radix	8	barnes	8

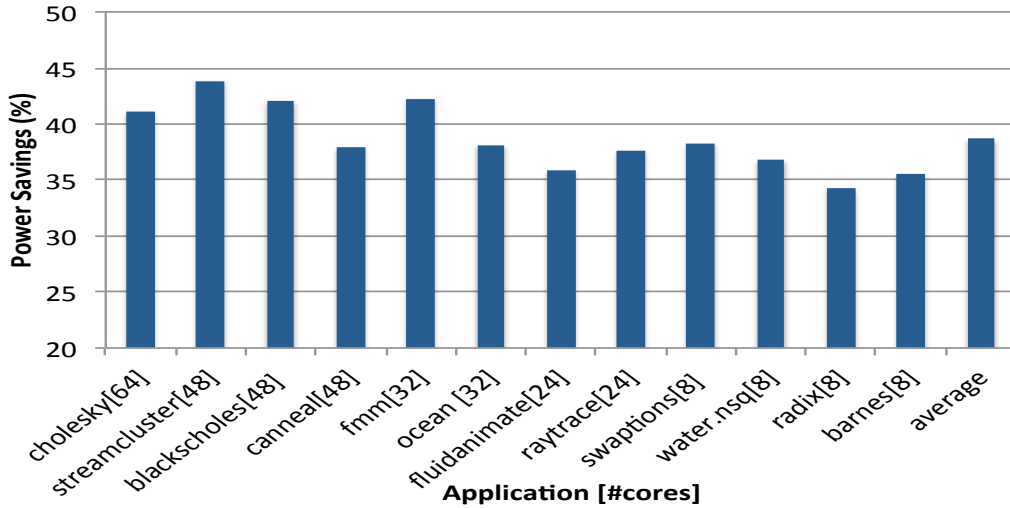


Figure 12: Power savings obtained by the proposed technique.

## 470 6. Conclusion

471 In this work, we addressed the problem of energy-efficient performance sus-  
472 tainability in NTC manycores under workload- and process-variability. We in-  
473 vestigated the benefits of NTC for highly parallel applications, while considering  
474 the power inefficiencies of a scalable and low overhead NT power delivery net-  
475 work. We showed that NTC can provide performance guarantees while delivering  
476 significant energy savings (66% on average). Finally, we addressed throughput  
477 imbalance caused by the excessive impact of process variability at near-threshold  
478 (NT) voltages, by proposing a runtime algorithm that manages to reduce both the  
479 variation of the per-core throughput w.r.t. the minimum and the power consump-  
480 tion, on average, by 70% and 43.5% respectively for the reported results, while  
481 not impacting the overall performance.

## 482 References

- 483 [1] R. G. Dreslinski, M. Wieckowski, D. Blaauw, D. Sylvester, T. N. Mudge,  
484 Near-threshold computing: Reclaiming moore’s law through energy efficient  
485 integrated circuits, *Proceedings of the IEEE 98* (2010) 253–266.
- 486 [2] N. Pinckney, K. Sewell, R. G. Dreslinski, D. Fick, T. Mudge, D. Sylvester,  
487 D. Blaauw, Assessing the performance limits of parallelized near-threshold  
488 computing, in: *Proceedings of the 49th Annual Design Automation Confer-*  
489 *ence, ACM*, pp. 1147–1152.
- 490 [3] U. R. Karpuzcu, A. Sinkar, N. S. Kim, J. Torrellas, Energysmart: Toward  
491 energy-efficient manycores for near-threshold computing, in: *High Perform-*  
492 *ance Computer Architecture (HPCA2013)*, 2013 IEEE 19th International  
493 Symposium on, IEEE, pp. 542–553.
- 494 [4] I. Stamelakos, S. Xydis, G. Palermo, C. Silvano, Variation-aware voltage  
495 island formation for power efficient near-threshold manycore architectures,  
496 in: *Design Automation Conference (ASP-DAC)*, 2014 19th Asia and South  
497 Pacific, IEEE, pp. 304–310.
- 498 [5] C. Silvano, G. Palermo, S. Xydis, I. Stamelakos, Voltage island manage-  
499 ment in near threshold manycore architectures to mitigate dark silicon, in:  
500 *Proceedings of the conference on Design, Automation & Test in Europe,*  
501 *European Design and Automation Association*, p. 201.

- 502 [6] A. A. Sinkar, H. Wang, N. S. Kim, Maximizing throughput of  
503 power/thermal-constrained processors by balancing power consumption of  
504 cores, in: Quality Electronic Design (ISQED), 2014 15th International Sym-  
505 posium on, IEEE, pp. 633–638.
- 506 [7] Y. Turakhia, G. Liu, S. Garg, D. Marculescu, Thread progress equaliza-  
507 tion: Dynamically adaptive power and performance optimization of multi-  
508 threaded applications, CoRR abs/1603.06346 (2016).
- 509 [8] D. Fick, R. G. Dreslinski, B. Giridhar, G. Kim, S. Seo, M. Fojtik, S. Satpathy,  
510 Y. Lee, D. Kim, N. Liu, et al., Centip3de: A 3930dmips/w configurable  
511 near-threshold 3d stacked system with 64 arm cortex-m3 cores, in: Solid-  
512 State Circuits Conference Digest of Technical Papers (ISSCC), 2012 IEEE  
513 International, IEEE, pp. 190–192.
- 514 [9] S. Khare, S. Jain, Prospects of near-threshold voltage design for green com-  
515 puting, in: 2013 26th International Conference on VLSI Design and 2013  
516 12th International Conference on Embedded Systems, IEEE, pp. 120–124.
- 517 [10] J. Wang, J. Zhang, W. Zhang, K. Qiu, T. Li, M. Wu, Near threshold cloud  
518 processors for dark silicon mitigation: the impact on emerging scale-out  
519 workloads, in: Proceedings of the 12th ACM International Conference on  
520 Computing Frontiers, ACM, p. 4.
- 521 [11] A. Pahlevan, J. Picorel, A. P. Zarandi, D. Rossi, M. Zapater, A. Bartolini,  
522 P. G. Del Valle, D. Atienza, L. Benini, B. Falsafi, Towards near-threshold  
523 server processors, in: 2016 Design, Automation & Test in Europe Confer-  
524 ence & Exhibition (DATE), IEEE, pp. 7–12.
- 525 [12] H. Kaul, M. Anders, S. Hsu, A. Agarwal, R. Krishnamurthy, S. Borkar,  
526 Near-threshold voltage (ntv) design: opportunities and challenges, in: Pro-  
527 ceedings of the 49th Annual Design Automation Conference, ACM, pp.  
528 1153–1158.
- 529 [13] D.-C. Juan, S. Garg, J. Park, D. Marculescu, Learning the optimal operating  
530 point for many-core systems with extended range voltage/frequency scaling,  
531 in: Hardware/Software Codesign and System Synthesis (CODES+ ISSS),  
532 2013 International Conference on, IEEE, pp. 1–10.

- 533 [14] W.-C. Hsieh, W. Hwang, All digital linear voltage regulator for super-to  
534 near-threshold operation, *Very Large Scale Integration (VLSI) Systems*,  
535 *IEEE Transactions on* 20 (2012) 989–1001.
- 536 [15] S. K. Lee, D. Brooks, G.-Y. Wei, Evaluation of voltage stacking for near-  
537 threshold multicore computing, in: *Proceedings of the 2012 ACM/IEEE*  
538 *international symposium on Low power electronics and design*, ACM, pp.  
539 373–378.
- 540 [16] X. He, G. Yan, Y. Han, X. Li, Superrange: wide operational range power  
541 delivery design for both stv and ntv computing, in: *Design, Automation and*  
542 *Test in Europe Conference and Exhibition (DATE)*, 2014, IEEE, pp. 1–6.
- 543 [17] T. N. Miller, X. Pan, R. Thomas, N. Sedaghati, R. Teodorescu, Booster:  
544 Reactive core acceleration for mitigating the effects of process variation and  
545 application imbalance in low-voltage chips, in: *High Performance Computer*  
546 *Architecture (HPCA)*, 2012 IEEE 18th International Symposium on, IEEE,  
547 pp. 1–12.
- 548 [18] I. Stamelakos, A. Khajeh, A. Eltawil, G. Palermo, C. Silvano, F. Kurdahi, A  
549 system-level exploration of power delivery architectures for near-threshold  
550 manycores considering performance constraints, in: *VLSI (ISVLSI)*, 2016  
551 *IEEE Computer Society Annual Symposium on*, IEEE, pp. 484–489.
- 552 [19] W. Lee, Y. Wang, M. Pedram, Vrcon: dynamic reconfiguration of voltage  
553 regulators in a multicore platform, in: *Proceedings of the conference on*  
554 *Design, Automation & Test in Europe*, European Design and Automation  
555 Association, p. 365.
- 556 [20] A. A. Sinkar, H. R. Ghasemi, M. J. Schulte, U. R. Karpuzcu, N. S.  
557 Kim, Low-cost per-core voltage domain support for power-constrained high-  
558 performance processors, *Very Large Scale Integration (VLSI) Systems*,  
559 *IEEE Transactions on* 22 (2014) 747–758.
- 560 [21] D. Kanter, Inside Nehalem: Intel’s future processor and system.  
561 <http://www.realworldtech.com>, 2008.
- 562 [22] U. R. Karpuzcu, K. B. Kolluru, N. S. Kim, J. Torrellas, VARIUS-NTV: A  
563 microarchitectural model to capture the increased sensitivity of manycores  
564 to process variations at near-threshold voltages, in: *IEEE/IFIP International*  
565 *Conference on Dependable Systems and Networks*, DSN, pp. 1–11.



- 566 [23] S. Sarangi, B. Greskamp, R. Teodorescu, J. Nakano, A. Tiwari, J. Torrel-  
567 las, Varius: A model of process variation and resulting timing errors for  
568 microarchitects, *Semiconductor Manufacturing*, IEEE Transactions on 21  
569 (2008) 3–13.
- 570 [24] D. Marković, C. C. Wang, L. P. Alarcon, T.-T. Liu, J. M. Rabaey, Ultralow-  
571 power design in near-threshold region, *Proceedings of the IEEE 98* (2010)  
572 237–252.
- 573 [25] G. G. Faust, R. Zhang, K. Skadron, M. R. Stan, B. H. Meyer, Archfp: Rapid  
574 prototyping of pre-rtl floorplans, in: *VLSI-SoC*, pp. 183–188.
- 575 [26] B. H. Calhoun, A. Wang, N. Verma, A. Chandrakasan, Sub-threshold design:  
576 The challenges of minimizing circuit energy, in: *Proceedings of the 2006*  
577 *International Symposium on Low Power Electronics and Design, ISLPED*  
578 '06, ACM, New York, NY, USA, 2006, pp. 366–368.
- 579 [27] Y. Okuma, K. Ishida, Y. Ryu, X. Zhang, P.-H. Chen, K. Watanabe,  
580 M. Takamiya, T. Sakurai, 0.5-v input digital ldo with 98.7% current effi-  
581 ciency and 2.7- $\mu$ a quiescent current in 65nm cmos, in: *Custom Integrated*  
582 *Circuits Conference (CICC), 2010 IEEE*, IEEE, pp. 1–4.
- 583 [28] N. Pinckney, K. Sewell, R. G. Dreslinski, D. Fick, T. Mudge, D. Sylvester,  
584 D. Blaauw, Assessing the performance limits of parallelized near-threshold  
585 computing, in: *Proceedings of the 49th Design Automation Conference*, pp.  
586 1147–1152.
- 587 [29] S. Dighe, S. Vangal, P. Aseron, S. Kumar, T. Jacob, K. Bowman, J. Howard,  
588 J. Tschanz, V. Erraguntla, N. Borkar, et al., Within-die variation-aware  
589 dynamic-voltage-frequency scaling core mapping and thread hopping for an  
590 80-core processor, in: *Solid-State Circuits Conference Digest of Technical*  
591 *Papers (ISSCC), 2010 IEEE International*, IEEE, pp. 174–175.
- 592 [30] T. E. Carlson, W. Heirman, L. Eeckhout, Sniper: Exploring the level of ab-  
593 straction for scalable and accurate parallel multi-core simulations, in: *Inter-*  
594 *national Conference for High Performance Computing, Networking, Storage*  
595 *and Analysis (SC)*.
- 596 [31] S. Li, J. H. Ahn, R. D. Strong, J. B. Brockman, D. M. Tullsen, N. P. Jouppi,  
597 McPAT: an integrated power, area, and timing modeling framework for mul-  
598 ticore and manycore architectures, in: *Proceedings of the 42nd Annual*

- 599 IEEE/ACM International Symposium on Microarchitecture, MICRO 42, pp.  
600 469–480.
- 601 [32] J. M. Arnold, D. A. Buell, E. G. Davis, Splash 2, in: Proceedings of  
602 the fourth annual ACM symposium on Parallel algorithms and architectures,  
603 ACM, pp. 316–322.
- 604 [33] C. Bienia, S. Kumar, J. P. Singh, K. Li, The parsec benchmark suite: Char-  
605 acterization and architectural implications, in: Proceedings of the 17th in-  
606 ternational conference on Parallel architectures and compilation techniques,  
607 ACM, pp. 72–81.