μ Systems Research Group School of Electrical and Electronic Engineering



Extended Power and Energy Normalized Performance Models for Many-Core Systems

Mohammed A. Noaman Al-hayanni, Ashur Rafiev, Rishad Shafik, Fei Xia, Alex Yakovlev

> Technical Report Series NCL-EEE-MICRO-TR-2016-198

> > January 2016

Contact: m.a.n.al-hayanni@ncl.ac.uk, ashur.rafiev@ncl.ac.uk, rishad.shafik@ncl.ac.uk, fei.xia.@ncl.ac.uk, alex.yakovlev@ncl.ac.uk

Supported by EPSRC grant GR/2016

NCL-EEE-MICRO-TR-2016-198 Copyright © 2016 Newcastle University

μSystems Research Group School of Electrical and Electronic Engineering Merz Court Newcastle University Newcastle upon Tyne, NE1 7RU, UK

http://async.org.uk/

Extended Power and Energy Normalized Performance Models for Many-Core Systems

Mohammed A. Noaman Al-hayanni, Ashur Rafiev, Rishad Shafik, Fei Xia, Alex Yakovlev

January 2016

Abstract

Continued technology scaling in VLSI has enabled more and more computation cores to be integrated in the same chip. This has facilitated the parallelization of processing and the increase of performance whilst keeping energy consumption at reasonable levels. To study the potential improvement of performance in such many core systems, three existing models have been popular in both the research community and industry. Amdahl's law is the original speedup model that estimates the maximum performance improvement with fixed workloads. Gustafson's law is a popular model that introduces variable workloads and estimates fixed time speedup. Sun and Ni combined the above two models into one considering the memory-bounded situation. These models are further extended via Hill-Marty model through considerations of homogeneous and a limited assumption of heterogeneous core configurations to estimate performance computation via Pollack's rule. This report investigates into these models and extends them to cover a generalized assumption of core heterogeneity more relevant for contemporary many-core architectures. We also present power and energy models based on the extended heterogeneous models making them usable for power and energy normalized performance and similar system metrics. Our models, being entirely general, cover popular power and performance control methods such as Dynamic Voltage Frequency Scaling (DVFS), power gating, etc. A case study is performed with an ARM big.LITTLE architecture containing Cortex A7 and A15 cores, including a comprehensive analysis with different ratios of parallel and sequential workloads to identify the most energy-efficient system configuration based on these models.

1 Introduction

Technology scaling has facilitated significant performance improvement at reduced power consumption through increased operating frequency and smaller device geometries [1]. According to Dennard's CMOS scaling law [2] despite such smaller geometries the power density of these devices remains constant. This is because the number of transistors per unit of area is also increasing substantially, which also conforms to Moore's [3] and Koomey's laws [4]. Dennard's law further states that the performance per watt is growing exponentially, doubling every 1.5 years.

Over the years significant research has been carried out to understand the trend of performance growth with many interconnected cores. An examples of these models is Pollack's Rule, which suggests that performance is increasing approximately proportional to the square root of the complexity [5]. Following this rule, a twofold growth of the components in a double processor will provide double the performance, in contrast to a single

	Homogeneity	Heterogeneity	Power Consumption	Performance / Watt	Performance/Joule	Amdahl's Model	Gustafson's Model	Sun and Ni's Model
[6]	Yes	No	No	No	No	Yes	No	No
[7]	Yes	No	No	No	No	Yes	Yes	No
[8]	Yes	No	No	No	No	Yes	Yes	Yes
[14]	Yes	Simple	No	No	No	Yes	No	No
[15]	Yes	Simple	No	No	No	Yes	Yes	Yes
[16]	Yes	Simple	Yes	Yes	Yes	Yes	No	No
[17]	Yes	No	No	No	No	Yes	Yes	Yes
Extended	Yes	General	Yes	Yes	Yes	Yes	Yes	Yes
Model								

Table 1: Existing Speedup Models and the Proposed Model

processor [1]. Therefore, multicore systems will deliver further improvement in throughput and latency for the same die area.

The most appropriate metric to describe performance gain is speedup. The first scalable model in relation to the multicore processor model is explained by Amdahl's law [6]. It assumes that a fixed workload is executed in N processors of a multi- or many-core system and compares the throuput/perfromance with the same workload executed in a single processor. In 1988, Gustafson introduced the principle of scalable computing in multicore processors pertaining to the fixed time model. This model proposes a linear speedup model that increases the workload proportional to increasing machine scalability, while the execution time remains fixed [7].

In other words, more parallel processors complete larger workloads spending the same amount of time and the speedup is calculated according to how much larger the workload is in multiple cores compared with that in a single core. In 1990, Sun and Ni suggested a new model, which included extended workload calculations by considering the capability of the memory. It is important to note that the executed workload and time should change based on the capability of the system, while the performance calculations appeared linear within the increasing cores [8,9].

On the other hand, power consumption management is a significant issue in scalable systems. For instance, DVFS, clock gating and power gating techniques are designed for this reason. The fine grain power management suggested by [1], [10] [11] [12] [13] are some of the scaling techniques used in order to decrease power consumption. Speedup models described in existing studies for the comprehensive understanding of core modeling are listed in Table 1. The Hill-Marty model extended Amdahl's law to cover not only homogeneous structures but also heterogeneous configurations fitting a limited simple assumption of core heterogeneity applicable to such practical systems as CPU-GPU structures. [14]. The study in [17] extended this analysis to all three major speedup models. The authors of [15] evaluated the homogeneous speedup models alone. The other important issue represented by energy efficiency is demonstrated by [16] for the homogeneous and simple heterogeneous Amdahl's model.

From Table 1, it can be seen that the existing models [6-13], however, have a general limitation of not

studying the energy-efficiency of computer system configurations, in addition to limiting any study of core heterogeneity to a simple assumption only applicable to CPU-GPU like configurations. To address these limitation, this report makes the following contributions:

- . extends the assumption of system core heterogeneity to a completely general case covering such modern configurations as FPGA-based acceleration schemes, complex structures with many types of cores, complex Systems on Chip (SoC) including mobile computing platforms, data centers with large numbers of heterogeneous processing units, etc.;
- . extends the three major speedup models (Table 1) to estimate power and energy normalized speedup metrics [14–17];
- . studies the comparative power/performance trade-offs of these models for energy-efficient computing based on homogeneous and heterogeneous configurations;
- . incorporates representations of the effects of such power and energy optimization techniques as DVFS and clock and power gating in the power models, i.e. heterogeneity in power control methods in addition to core structures;
- . uses a mobile computing platform centered around ARM big.LITTLE Cortex A7-A15 cores in the form of Odroid-XU3 as a case study covering all aspects of the new modeling.

To the best of our knowledge this is the first comprehensive power and energy normalized performance analysis of the major many-core speedup models. It also represents the first attempt to extend these models to cover a fully general assumption of core heterogeneity. The rest of the report is organized as follows. Section II gives the background on existing speedup models for homogeneous systems; Section III extends existing speedup models to cover the general assumption of core heterogeneity; Section IV derives the average power consumption models for all three extended models; Section V describes a method for power and energy normalized performance analysis of these extended models for homogeneous and heterogeneous configurations; Section VI describes the case study; Section VII gives the outcomes of power and energy normalized performance analysis of homogeneous and heterogeneous models; And Section VIII concludes the report.

2 Homogeneous Speedup Models

For a homogeneous system we consider a system consisting of N cores, each core having performance of IPS_1 instructions per second. This section describes various existing models for determining the system's speedup SP(N) in relation to a single core, which can be used to find the performance of the system:

$$IPS_N = SP(N) \cdot IPS_1. \tag{1}$$

The parallel part of a workload is P and the sequential part is (1 - P), the communication overhead is negligible.

2.1 Amdahl's Law (Fixed Workload)

The general idea of this model is to compare execution time for some fixed workload WL on a single core with the execution time for the same workload on the entire N-core system [6].

Time T(1) to execute workload WL on a single core is WL/IPS_1 , whereas T(N) adds up the sequential execution time on one core and the parallel execution time on all N cores:

$$T(N) = \frac{(1-P) \cdot WL}{IPS_1} + \frac{P \cdot WL}{N \cdot IPS_1},$$
(2)

thus the speed up can be found as follows:

$$SP(N) = \frac{T(1)}{T(N)} = \frac{1}{(1-P) + \frac{P}{N}}.$$
(3)

2.2 Gustafson's Model (Fixed Time)

Gustafson re-evaluated the fixed workload speedup model to derive a new fixed time model [7]. In this model, the workload increases with number of cores, while the execution time is fixed.

Let's denote the initial workload and extended workload as WL and WL' respectively. The time to execute initial workload and expanded workload are T(N) and T'(N) respectively. The workload scaling ratio can be found from:

$$T(1) = \frac{WL}{IPS_1},\tag{4}$$

$$T(N) = \frac{(1-P) \cdot WL}{IPS_1} + \frac{P \cdot WL'}{N \cdot IPS_1}.$$
(5)

and, since T(1) = T(N), the extended workload can be found as:

$$WL' = N \cdot WL. \tag{6}$$

$$T'(1) = \frac{(1-P) \cdot WL}{IPS_1} + \frac{P \cdot N \cdot WL}{IPS_1}.$$
(7)

From the relation of scaled and unscaled execution time the following equation for speedup can be calculated:

$$SP(N) = \frac{T'(1)}{T(1)} = (1 - P) + P \cdot N.$$
(8)

The sequential part of the workload uses one core to perform its calculation at the performance IPS_1 , and the parallel execution uses N cores to perform its calculation at the performance $N \cdot IPS_1$.

2.3 Sun and Ni's Model (Memory Bounded)

Sun and Ni mixed the previous two speedup models by consider the memory bounded constrains [8,9]. In this model the execution time and the workload change according to the memory capability. The parameter g(N) reflects the scaling of the workload in relation to scaling the memory with the number of cores:

$$WL' = g(N) \cdot WL. \tag{9}$$

A typical example g(N) is given for an $M \times M$ matrix multiplication, which has the memory requirement of M^2 and the computation cost (workload) of M^3 . In this case, $g(N) = N^{\frac{3}{2}}$. The time to execute the scaled



Figure 1: The proposed general structure of a heterogeneous system (c) compared to a homogeneous system (a) and the previous assumption [14] on heterogeneity (b).

workload can be found from (4) and (5).

$$T'(1) = \frac{(1-P) \cdot WL}{IPS_1} + \frac{P \cdot g(N) \cdot WL}{IPS_1},$$
(10)

$$T'(N) = \frac{(1-P) \cdot WL}{IPS_1} + \frac{P \cdot g(N) \cdot WL}{N \cdot IPS_1}.$$
(11)

The speedup is calculated as follows:

$$SP(N) = \frac{T'(1)}{T'(N)} = \frac{(1-P) + P \cdot g(N)}{(1-P) + \frac{P \cdot g(N)}{N}}.$$
(12)

Because the workload is scaled by g(N) according to (9), one of the important properties of this model is that for g(N) = 1 Sun and Ni's model (12) transforms into Amdahl's Law (3), and for g(N) = N it becomes Gustafson's Law (8).

3 Heterogeneous Speedup Models

Previous attempts to extend speedup laws to heterogeneous systems were mainly focused on a single highperformance core and many smaller cores [14], which is relevant for CPU+GPU systems. In this work we aim to fully generalize the models. Hence, for a heterogeneous system we consider a system consisting of X clusters (types) of homogeneous cores with number of cores defined as a vector $\overline{N} = (N_1, \dots, N_X)$. Vector $\overline{\alpha} = (\alpha_1, \dots, \alpha_X)$ defines the performance of each core by cluster (type) in relation to some *base core equivalent* (BCE), such that for some $1 \le i \le X$ we have $IPS_i = \alpha_i \cdot IPS_1$. The structure is shown in Figure 1. This section extends homogeneous speedup models for determining the heterogeneous system's speedup $SP(\overline{N})$ in relation to a single BCE, which can then be used to find the performance of the system using (1).

3.1 Heterogeneous Amdahl's Law (Fixed Workload)

For heterogeneous systems the combined performance of all cores while executing the parallel code P can be found as a weighted sum of all performances:

$$N_{\alpha} \cdot IPS_1 = \sum_{i=1}^{X} N_i \cdot \alpha_i \cdot IPS_1, \tag{13}$$

 N_{α} is called a *performance-equivalent number* of BCEs. In other words, this performance is equal to N_{α} BCE cores executing the same parallel code; N_{α} can be a fractional number. However, in the case of synchronized-parallel execution (i.e. if the parallel execution waits for the slowest core to finish), a different equation has to be used to find N_{α} :

$$N_{\alpha} = \min \overline{\alpha} \cdot \sum_{i=1}^{X} N_i.$$
(14)

For the use case calculations in this report we used (13).

We also assume that the sequential part is executed on a single core in the cluster X. Hence, the time to execute the fixed workload WL on the given heterogeneous system is:

$$T_X(\overline{N}) = \frac{(1-P) \cdot WL}{\alpha_X \cdot IPS_1} + \frac{P \cdot WL}{N_\alpha \cdot IPS_1}.$$
(15)

The speedup in relation to single BCE is:

$$SP\left(\overline{N}\right) = \frac{T\left(1\right)}{T_X\left(\overline{N},\overline{\alpha}\right)} = \frac{1}{\frac{(1-P)}{\alpha_X} + \frac{P}{N_\alpha}}.$$
(16)

3.2 Heterogeneous Gustafson's Model (Fixed Time)

Because of the workload scaling, we cannot directly compare speedup while executing the sequential code in the core *X* to single BCE execution. Let's first find the speedup $SP_X(\overline{N})$ relative to a single core *X*. This is done similarly to Gustafson's derivation (Section 2.2).

$$T_X(1) = \frac{WL}{\alpha_X \cdot IPS_1},\tag{17}$$

$$T_X(N) = \frac{(1-P)WL}{\alpha_X \cdot IPS_1} + \frac{P \cdot WL'}{N_\alpha \cdot IPS_1},$$
(18)

 $T_X(1) = T_X(N)$, hence the extended workload can be found as:

$$WL' = \frac{N_{\alpha}}{\alpha_X} \cdot WL. \tag{19}$$

$$T_X'(1) = \frac{(1-P) \cdot WL}{\alpha_X \cdot IPS_1} + \frac{P \cdot N_\alpha \cdot WL}{\alpha_X \cdot IPS_1},$$
(20)

$$SP_X\left(\overline{N}\right) = \frac{T'_X(1)}{T_X(1)} = (1-P) + \frac{P \cdot N_\alpha}{\alpha_X}.$$
(21)

The speedup of a single core X in relation to BCE is α_X , thus the total speedup relative to BCE is:

$$SP(\overline{N}) = \alpha_X \cdot SP_X(\overline{N}) = (1 - P) \cdot \alpha_X + P \cdot N_\alpha.$$
⁽²²⁾

The sequential part of the workload uses one core in the cluster X to perform its calculation at the performance $\alpha_X \cdot IPS_1$, and the parallel execution uses all cores to perform its calculation at the performance $N_{\alpha} \cdot IPS_1$.

3.3 Heterogeneous Sun and Ni's Model (Memory Bounded)

Similarly to Amdahl's and Gustafson's cases, we can extend Sun and Ni's model to the general heterogeneous case as follows:

$$T_X'(1) = \frac{(1-P) \cdot WL}{\alpha_X \cdot IPS_1} + \frac{P \cdot g(N) \cdot WL}{\alpha_X \cdot IPS_1},$$
(23)

$$T'_{X}(N) = \frac{(1-P) \cdot WL}{\alpha_{X} \cdot IPS_{1}} + \frac{P \cdot g\left(\overline{N}\right) \cdot WL}{N_{\alpha} \cdot IPS_{1}},$$
(24)

$$SP_X(\overline{N}) = \frac{T'_X(1)}{T'_X(N)} = \frac{1}{\alpha_X} \cdot \frac{(1-P) + P \cdot g(\overline{N})}{\frac{(1-P)}{\alpha_X} + \frac{P \cdot g(\overline{N})}{N_\alpha}}.$$
(25)

The speedup in a heterogeneous system relative to BCE is calculated as follows:

$$SP(\overline{N}) = \alpha_X \cdot SP_X(\overline{N}) = \frac{(1-P) + P \cdot g(\overline{N})}{\frac{(1-P)}{\alpha_X} + \frac{P \cdot g(\overline{N})}{N_\alpha}}.$$
(26)

When $g(\overline{N}) = 1$, this model transforms into heterogeneous Amdahl's Law (16), and for $g(\overline{N}) = \frac{N_{\alpha}}{\alpha_{\chi}}$ it becomes heterogeneous Gustafson's Law (22), as expected from (19).

For all heterogeneous models, substitution $\alpha_X = 1, N_\alpha = N$ will give the homogeneous versions of respective models. In other words, homogeneity is a special case of heterogeneity.

4 Average Power Consumption Models

The power consumption models are built under the assumption that the cores consume power when idle. When idle power is zero, this assumption covers the special case of power gating.

Let's the active power of a core in the homogeneous system (Section 2) be W_A and the idle power of a core be W_i respectively. Active power can also be expressed as a sum of idle power and effective power W_1 (used for computation), $W_A = W_1 + W_i$. In the total power consumption of the system, the constant term of total idle power W_{idle} does not benefit to the model and can be added later. The power models W(N) are focused on the effective power, and the total power of the system can be calculated as follows:

$$W_{total} = W(N) + W_{idle}, \tag{27}$$

In the heterogeneous system (Section 3), the difference between power consumptions of the cores is ex-

pressed by the vector $\overline{\beta} = (\beta_1, ..., \beta_X)$, which defines the effective power in relation to a BCE's effective power, such that for some $1 \le j \le X$ we have effective power $W_j = \beta_j \cdot W_1$. All idle powers of heterogeneous cores are combined into W_{idle} . In general case, we say that:

$$W_{idle} = N_i \cdot W_i, \tag{28}$$

where N_i is idle power equivalent number of BCEs and W_i is the idle power of a single BCE.

The effective power model can be found as a time-weighted average of the sequential power W_S the and parallel power W_P :

$$W(N) = \frac{W_{S} \cdot T_{S}(N) + W_{P} \cdot T_{P}(N)}{T_{S}(N) + T_{P}(N)},$$
(29)

where $T_S(N)$ and $T_P(N)$ are speedup-dependent times to execute sequential and parallel parts respectively. In the homogeneous system:

$$W_S = W_1, W_P = N \cdot W_1.$$
 (30)

In the heterogeneous system, if we execute the sequential code on a single core X:

$$W_S = \beta_X \cdot W_1,$$

$$W_P = W_1 \cdot \sum_{j=1}^X \beta_j \cdot N_j = N_\beta \cdot W_1.$$
(31)

 N_{β} is called a *power-equivalent number* of BCEs. Heterogeneous power models will transform into homogeneous if $\alpha_X = \beta_X = 1$ and $N_{\alpha} = N_{\beta} = N$.

4.1 Power Model for Amdahl's Law (Fixed Workload)

From (15) we know that:

$$T_{S}\left(\overline{N}\right) = \frac{(1-P) \cdot WL}{\alpha_{X} \cdot IPS_{1}}, \ T_{P}\left(\overline{N}\right) = \frac{P \cdot WL}{N_{\alpha} \cdot IPS_{1}}.$$
(32)

By substituting (32) and (31) into (29) we have a power model for the heterogeneous system:

$$W\left(\overline{N}\right) = \left(\frac{\beta_X}{\alpha_X} \cdot (1-P) + \frac{N_\beta}{N_\alpha} \cdot P\right) \cdot SP\left(\overline{N}\right) \cdot W_1,\tag{33}$$

where the speedup $SP(\overline{N})$ is calculated using (16). For homogeneous system, this will transform into:

$$W(N) = SP(N) \cdot W_1, \tag{34}$$

thus for Amdahl's Law the power scales with the speedup.

4.2 Power Model for Gustafson's Model (Fixed Time)

In this model we have a fixed time T, so the workload splits execution into:

$$T_{S}\left(\overline{N}\right) = (1-P) \cdot T, \ T_{P}\left(\overline{N}\right) = P \cdot T.$$
(35)

Thus, we can find a power model for the heterogeneous system:

$$W\left(\overline{N}\right) = \left(\frac{\beta_X \cdot (1-P) + N_\beta \cdot P}{\alpha_X \cdot (1-P) + N_\alpha \cdot P}\right) \cdot SP\left(\overline{N}\right) \cdot W_1,\tag{36}$$

For homogeneous system, this will transform into:

$$W(N) = SP(N) \cdot W_1, \tag{37}$$

where the speedup $SP(\overline{N})$ is calculated using (22).

4.3 Power Model for Sun and Ni's Model (Memory Bounded)

From (24) we can find:

$$T_{S}\left(\overline{N}\right) = \frac{(1-P) \cdot WL}{\alpha_{X} \cdot IPS_{1}}, \ T_{P}\left(\overline{N}\right) = \frac{P \cdot g\left(\overline{N}\right) \cdot WL}{N_{\alpha} \cdot IPS_{1}}.$$
(38)

Thus, we can find a power model for the heterogeneous system:

$$W\left(\overline{N}\right) = \left(\frac{\frac{\beta_X}{\alpha_X} \cdot (1-P) + \frac{N_\beta}{N_\alpha} \cdot P \cdot g\left(\overline{N}\right)}{(1-P) + P \cdot g\left(\overline{N}\right)}\right) \cdot SP\left(\overline{N}\right) \cdot W_1,\tag{39}$$

where the speedup $SP(\overline{N})$ is calculated using (26). This model will transform into (33) if $g(\overline{N}) = 1$, or (36) for $g(\overline{N}) = \frac{N_{\alpha}}{\alpha_{\nu}}$. For homogeneous system, (39) will also transform into:

$$W(N) = SP(N) \cdot W_1. \tag{40}$$

All power models – (33), (36), and (39) – can be represented using *power scaling* $PS(\overline{N})$, which can be derived from the respective model equations:

$$W\left(\overline{N}\right) = PS\left(\overline{N}\right) \cdot SP\left(\overline{N}\right) \cdot W_1.$$
(41)

5 Power-normalized and Energy-normalized Performance

The power model explains the total power consumption in this model during workload execution. It is likewise represents the cooling capacity. Furthermore, it is simple to model the performance achievable at the same cooling capacity from calculating performance per watt (Perf/Watt). This model is reciprocal of energy per instruction (EPI_N) because performance is the reciprocal of execution time.

 EPI_N can be found from dividing the total power (27) by the system's performance (1):

$$EPI_N = \frac{W_{total}}{IPS_N} = \frac{W(N) + W_{idle}}{IPS_1 \cdot SP(\overline{N})},$$
(42)

which is true for all cases of $W(\overline{N})$: Amdahl's, Gustafson's, or Sun and Ni's.

For a single BCE we can denote energy per instruction as a sum of effective energy EPS_1 and idle energy EPS_i :

$$EPI_{BCE} = \frac{W_1}{IPS_1} + \frac{W_i}{IPS_1}.$$
(43)

Applying the power model (41) to (42) and also considering (28), we find:

$$EPI_{N} = EPI_{1} \cdot PS\left(\overline{N}\right) + \frac{N_{i} \cdot EPI_{i}}{SP\left(\overline{N}\right)}.$$
(44)

This equation shows that the effective component of the energy increases with the power scaling $PS(\overline{N})$, and the idle energy decreases with the speedup $SP(\overline{N})$.

Energy-normalized performance represents how much performance one can gain if willing to increase energy per operation. This gain in relation to BCE can be found as:

$$\left(\frac{IPS_N}{EPI_N}\right) \cdot \left(\frac{IPS_1}{EPI_{BCE}}\right)^{-1} = SP\left(\overline{N}\right) \cdot \frac{EPI_N}{EPI_{BCE}}.$$
(45)

The equation shows that the increment factor scales with the speedup, and this is true for all three models.

6 Case Study

We carried out an extensive case study demonstrating the use of these models. This study is based on a multi-core mobile platform, the Odroid-XU3 board [18]. The main part of it is the 28nm Application Processor Exynos 5422. It is an SoC hosting an ARM big.LITTILE heterogeneous octa-core processor consisting of four Cortex A7 cores and four Cortex A15 cores. The big Cortex-A15 is a high performance 32-bit core having 32 KB instruction and 32KB data L1 caches and 2 MB L2 cache and the maximum frequency of 2.0 GHz. The LITTLE Cortex-A7 is a low power 32-bit core including the same L1 cache size and 512 KB L2 cache, and the maximum frequency of 1.4 GHz.

This SoC also has four power domains: A7 power domain, A15 power domain, GPU and memory power domains. The Odroid-XU3 board allows per-domain DVFS using voltage-frequency pairs, however for frequencies within the range of 200MHz to 800MHz, the voltage stays constant (DFS-only).

The traditional simple assumption for heterogeneous architectures, shown in ure 1(b), cannot describe systems such as big.LITTLE. Hence, the presented use case is a perfect application for our generalized heterogeneous models.

A set of characterization experiments was carried out to determine power consumptions and performances for each core type. The main parameters for this study can be arranged into the following points.

System's heterogeneity

Following the general heterogeneous system structure assumption proposed in Section III, we can set the constants for Odroid-XU3. Two types of cores (A7 and A15) give us X = 2. In our experiments, in order to improve measurement accuracy, one of the A7 cores was reserved for exclusive use by the operating system. Therefore the numbers of cores by type are $N_{A7} = 3$ and $N_{A15} = 4$.

Core active powers

Theoretical active power calculations are derived from the experiments. The power is measured while executing a full workload on the processor's cores and sweeping through all DVFS points [19]. In general, the theoretical dynamic power estimation can be calculated by the power equation [20]:

$$Power_{Dynamic} = C \cdot V^2 \cdot F, \tag{46}$$

where V is the voltage, F is the frequency, and C is the constant, which relates to the combined capacitance of the switching logic. We used the experimental data to curve-fit by MATLAB and derive the Cortex A7 and A15 power equations. The result is the following values for C: in A7 it is equal to 1.0nF and in A15 it is 6.0nF. Considering A7 as BCE, we have $\beta_{A7} = 1$, $\beta_{A15} = 6.0$ to supply our power models from Section IV.

Core idle powers

Theoretical idle power calculations supported by the experiments [19] give the value of A15 idle power as $W_{iA15} = 0.021W$. Having one of the A7 cores occupied by the operating system prevents measuring idle power for that domain. In this study we accept the minimum measured power of 0.008W as the domain's idle power W_{iA7} , which is also our BCE's idle power W_i . We also do not switch off the cores, so the total idle power of the system remains constant: $W_{idle} = N_i \cdot W_i$. From the above measurements, we calculated $N_i = 14.5$.

Relative performance of cores

We did not use performance counters to find the actual number of clocks per instruction (*CPI*) for different types of instructions. Given these are RISC processors, we assume the general value of *CPI* = 1 without losing generality. From the characterization experiments, we found that on the average an A15 core has 1.5 times the throughput of an A7 core when both are running at the same frequency. We also want to calculate the models for different DVFS points. However, for frequency values when A7 cannot be run and A15 can, or for the DFS-only region, the performance specifying vector $\overline{\alpha}$ changes. Therefore, for each DVFS point the value for α_{A15} is computed as follows:

$$\alpha_{A15} = 1.5 \cdot \frac{freq_{A15}}{freq_{A7}}.$$
(47)

We assume $\alpha_{A7} = 1$ considering A7 is our BCE.

Parallelization parameter

The speedup models take parameter *P* from the nature of the executed workload. In our theoretical calculations we investigated a number of values for *P*. In the next section we provide results for two example values (P = 0.9 and P = 0.1) covering highly parallelizable and not parallelizable cases.

7 Outcomes

In this section we present selected calculation results organized in three groups.



Figure 2: Metric explorations on a fixed DVFS point (1400MHz)



Figure 3: Energy per Instruction - P = 0.9.



Figure 4: Energy per Instruction - P = 0.1.

7.1 Metric explorations on a fixed DVFS point

In the first group of results we present the following metrics of interest calculated for a various combination of active and idle cores on a fixed DVFS point ($freq_{A7} = freq_{A15} = 1400$ MHz):

- Performance IPS_N according to (1),
- Average power consumption W_{total} according to (27),
- Energy per instruction EPI_N according to (44),
- Energy-normalized performance according to (45).

These parameters have been estimated for all presented heterogeneous speedup and power models. The graphs for different models display similar trends, hence we only present Sun and Ni's model for its generality;



Figure 5: Energy Normalized Performance - P = 0.9.



Figure 6: Energy Normalized Performance - P = 0.1.



Figure 7: Energy per instruction for a homogeneous system.

g(N) was set to the matrix multiplication example presented in Section II:

$$g\left(\overline{N}\right) = \left(\frac{N_{\alpha}}{\alpha_{A15}}\right)^{\frac{3}{2}}$$

NCL-EEE-MICRO-TR-2016-198, Newcastle University





Figure 8: Performance of homogeneous A7 cores



Figure 9: Average power consumption of homogeneous A7 cores



Figure 10: Energy per instruction of homogeneous A7 cores

Figure 8 shows the graphs for the listed parameters for P = 0.9 and P = 0.1.

It can be seen from the data that although the power consumption increases with the number of cores participating in the computation, the performance also increases with the cumulative effect being that the performance per unit of power spent still improving with more cores used. This is mainly because of the influence of idle power. If you don't use a core, the idle power is wasted.

The higher the parallelization factor, the better the performance and energy-normalized performance, as expected.

More interestingly, from the energy per instruction metric one can see it increasing when the number of A15 cores increase but decrease when the number of A7 cores increase. This is on account of the much higher efficiency of A7 cores in terms of energy per instruction.



Figure 11: Power normalized performance of homogeneous A7 cores



Figure 12: Energy normalized performance of homogeneous A7 cores



Figure 13: Performance of heterogeneous cores - P = 0.9

7.2 Frequency scaling

This group of results illustrate the scaling of the energy per instruction EPI_N and the energy-normalized performance with the system's frequency. The values for frequencies have been selected within the allowed range of 200MHz to 2000MHz and the same frequency have been set for A7 and A15 cores if possible (for values above 1400MHz, the frequency for A7 is set to the allowed maximum of 1400MHz). This point causes a



Figure 14: Performance of heterogeneous cores - P = 0.1



Figure 15: Average power consumption of heterogeneous cores - P = 0.9

non-smooth change in α_{A15} leading to a peculiar non-smooth behaviour of the metrics. There are two other less obvious behaviour boundaries: 800MHz, where DVFS switches to DFS, and 1900MHz, above which the A15 cores experience throttling because of thermal issues. All these points are reflected by our models. Figure 3 shows the graphs for the energy per instruction EPI_N in different heterogeneous core combinations for P = 0.9. Figure 5 shows the graphs for the energy-normalized performance, also for P = 0.9.

7.3 Homogeneous example

Figure 7 presents an example of applying the presented models to a homogeneous system for completeness, demonstrating that X = 1 also works. From this figure, we can make an interesting observation: if you put more cores to solving a problem with a low parallelization capability (P = 0.1), energy per instruction suffers, especially at the lower frequencies.



Figure 16: Average power consumption of heterogeneous cores - P = 0.1



Figure 17: Power normalized performance of heterogeneous cores - P = 0.9

8 Conclusions

This report addresses the emerging issue of the system heterogeneity becoming more common and diverse in its structure well beyond the traditional CPU+GPU assumption. This is done by introducing the general model for system heterogeneity. The three known speedup models (Amdahl's Law, Gustafson's model, Sun and Ni's model) are extended to cover this general heterogeneity. In addition to performance speedup, the report presents the models for power and energy related system metrics.

The derived theoretical models have been applied to a real-life heterogeneous system, whose structure does not fit into the traditional heterogeneity assumption. The model parameters have been characterized from a set of experiments, and the metrics of interest have been calculated to demonstrate the model capabilities. These metric include performance scaling, average power scaling, energy per instruction, and energy-normalized performance. However, so far we have not done cross-validation of the results against any other set of experimental results. This cross-validation is a candidate for immediate future work.



Figure 18: Power normalized performance of heterogeneous cores - P = 0.1



Figure 19: Energy normalized performance of homogeneous A7 cores

Acknowledgment The authors wish to thank Rem Gensh and Ali M. M. Aalsaud for providing the experimental data for Odroid-XU3.

References

- [1] S. Borkar, "Thousand core chips: a technology perspective," in *Proceedings of the 44th annual Design Automation Conference*. ACM, 2007, pp. 746–749.
- [2] R. H. Dennard, V. Rideout, E. Bassous, and A. Leblanc, "Design of ion-implanted mosfet's with very small physical dimensions," *Solid-State Circuits, IEEE Journal of*, vol. 9, no. 5, pp. 256–268, 1974.
- [3] G. E. Moore *et al.*, "Cramming more components onto integrated circuits," *Proceedings of the IEEE*, vol. 86, no. 1, pp. 82–85, 1998.

- [4] J. G. Koomey, S. Berard, M. Sanchez, and H. Wong, "Implications of historical trends in the electrical efficiency of computing," *Annals of the History of Computing, IEEE*, vol. 33, no. 3, pp. 46–54, 2011.
- [5] F. J. Pollack, "New microarchitecture challenges in the coming generations of cmos process technologies (keynote address)," in *Proceedings of the 32nd annual ACM/IEEE international symposium on Microarchitecture*. IEEE Computer Society, 1999, p. 2.
- [6] G. M. Amdahl, "Validity of the single processor approach to achieving large scale computing capabilities," in *Proceedings of the April 18-20, 1967, spring joint computer conference.* ACM, 1967, pp. 483–485.
- [7] J. L. Gustafson, "Reevaluating amdahl's law," *Communications of the ACM*, vol. 31, no. 5, pp. 532–533, 1988.
- [8] X.-H. Sun and L. M. Ni, "Another view on parallel speedup," in *Supercomputing* '90., *Proceedings of*. IEEE, 1990, pp. 324–333.
- [9] —, "Scalable problems and memory-bounded speedup," *Journal of Parallel and Distributed Computing*, vol. 19, no. 1, pp. 27–37, 1993.
- [10] J. W. Tschanz, S. G. Narendra, Y. Ye, B. Bloechel, S. Borkar, V. De *et al.*, "Dynamic sleep transistor and body bias for active leakage power control of microprocessors," *Solid-State Circuits, IEEE Journal of*, vol. 38, no. 11, pp. 1838–1845, 2003.
- [11] S. Eyerman and L. Eeckhout, "Fine-grained dvfs using on-chip regulators," ACM Transactions on Architecture and Code Optimization (TACO), vol. 8, no. 1, p. 1, 2011.
- [12] A. Das, M. Schuchhardt, N. Hardavellas, G. Memik, and A. Choudhary, "Dynamic directories: A mechanism for reducing on-chip interconnect power in multicores," in *Proceedings of the Conference on Design*, *Automation and Test in Europe*. EDA Consortium, 2012, pp. 479–484.
- [13] T. S. Muthukaruppan, A. Pathania, and T. Mitra, "Price theory based power management for heterogeneous multi-cores." in ASPLOS, 2014, pp. 161–176.
- [14] M. D. Hill and M. R. Marty, "Amdahl's law in the multicore era," Computer, no. 7, pp. 33–38, 2008.
- [15] N. Ye, Z. Hao, and X. Xie, "The speedup model for manycore processor," in *Information Science and Cloud Computing Companion (ISCC-C)*, 2013 International Conference on. IEEE, 2013, pp. 469–474.
- [16] D. H. Woo and H.-H. S. Lee, "Extending amdahl's law for energy-efficient computing in the many-core era," *Computer*, no. 12, pp. 24–31, 2008.
- [17] X.-H. Sun and Y. Chen, "Reevaluating amdahlâs law in the multicore era," *Journal of Parallel and Dis*tributed Computing, vol. 70, no. 2, pp. 183–188, 2010.
- [18] (2015) Odroid platform. [Online]. Available: http://www.hardkernel.com/main/products/
- [19] R. Gensh, A. Aalsaud, A. Rafiev, F. Xia, A. Iliasov, A. Romanovsk, and A. Yakovlev, "Experiments with odroid-xu3 board," Newcastle University, Computing Science, Claremont Tower, Claremont Road,Newcastle Upon Tyne, NE1 7RU, England., Tech. Rep., 2015.

[20] J. Rabaey and M. Pedram, Low Power Design Methodologies, ser. The Springer International Series in Engineering and Computer Science. Springer US, 2012. [Online]. Available: https: //books.google.co.uk/books?id=9IzuBwAAQBAJ