

# AdaMD: Adaptive Mapping and DVFS for Energy-efficient Heterogeneous Multi-cores

Karunakar R. Basireddy, *Student Member, IEEE*, Amit Kumar Singh, *Member, IEEE*, Bashir M. Al-Hashimi, *Fellow, IEEE*, and Geoff V. Merrett, *Member, IEEE*

**Abstract**—Modern heterogeneous multi-core systems, containing various types of cores, are increasingly dealing with concurrent execution of dynamic application workloads. Moreover, the performance constraints of each application vary, and applications enter/exit the system at any time. Existing approaches are not efficient in such dynamic scenarios, especially if applications are unknown, as they require extensive offline application analysis and do not consider the runtime execution scenarios (application arrival/completion, and workload and performance variations) for runtime management. To address this, we present AdaMD, an adaptive mapping and dynamic voltage and frequency scaling (DVFS) approach for improving energy consumption and performance. The key feature of the proposed approach is the elimination of dependency on offline profiled results while making runtime decisions. This is achieved through a performance prediction model having a maximum error of 7.9% lower than the previously reported model and a mapping approach that allocates processing cores to applications while respecting performance constraints. Furthermore, AdaMD adapts to runtime execution scenarios efficiently by monitoring the application status, and performance/workload variations to adjust the previous DVFS settings and thread-to-core mappings. The proposed approach is experimentally validated on the Odroid-XU3, with various combinations of diverse multi-threaded applications from PARSEC and SPLASH benchmarks. Results show energy savings of up to 28% compared to the recently proposed approach while meeting performance constraints.

**Index Terms**—Heterogeneous multi-cores, Multi-threaded applications, Run-time management, Energy savings.

## I. INTRODUCTION

Modern mobile platforms are containing greater number of heterogeneous cores to support highly diverse and varying workloads (e.g., the Odroid-XU3 [1] and Mediatek X20 [2]). Such platforms often execute applications concurrently, which simultaneously contend for system resources and typically exhibit varying resource demands over time [3]. Each application may have different performance requirements and exhibit various workload phases during its execution [4]. To adapt to such dynamic scenarios, mobile platforms offer an increasing number of resource configurations, such as enabling and disabling cores of different types, defining the thread-to-core mapping for a multi-threaded application, and setting dynamic voltage and frequency (DVFS) operating points.

B. K. Reddy, B. M. Al-Hashimi, and G. V. Merrett are with the School of Electronics and Computer Science, University of Southampton, United Kingdom (e-mail: krb1g15@ecs.soton.ac.uk; bmah@ecs.soton.ac.uk; gvm@ecs.soton.ac.uk).

A. K. Singh is with the School of Computer Science and Electronic Engineering, University of Essex, Colchester CO43SQ, United Kingdom (email:a.k.singh@essex.ac.uk).

Manuscript received Jan xx, 2019; revised xxx xx, 2019.

The process of thread-to-core mapping and setting DVFS levels play a crucial role in exploiting the system properties such that applications can meet their, often diverse, demands on performance and energy consumption [3]. In general, for each application, the management process first finds a thread-to-core mapping, and then core DVFS level by inspecting the workload profile while satisfying the performance requirement. This problem becomes much more complex when dynamically mapping concurrently executing applications due to contention for resources, and when the mapping is coupled with DVFS, i.e., energy-efficient allocation of processing cores and selection of DVFS settings [5], [6].

The reported approaches for solving this problem fall into three categories: 1) offline, 2) online, and 3) hybrid approaches. Several offline approaches have been proposed targeting different application domains and hardware architectures [7], [8]. These typically use computationally intensive search methods to find the optimal or near-optimal mapping for the applications that may run on the system. Conversely, online approaches [4], [9]–[11] must not be computationally intensive, as they are required to make efficient application mapping/DVFS decisions at runtime. Therefore, these techniques generally use heuristics to find a suitable platform configuration. Design time approaches usually find solutions of higher quality compared to online techniques, due to extensive design space exploration of the underlying hardware and applications. To address the drawbacks of pure offline and online approaches, various hybrid approaches [8], [12]–[17] using offline analysis to make runtime decisions based on the current state of the system are proposed.

However, a review of the prior arts (see section VI) shows that the existing approaches, targeting heterogeneous multi-cores, have the following shortcomings. They use heavy application-dependent profile data and thus are not efficient in managing dynamic workloads when unknown applications with different performance constraints are executing concurrently. For example, the number of different frequency and core configurations for the Odroid-XU3 platform [1] (four big and four LITTLE cores that can operate at 13 and 19 different frequencies, respectively) is 4080 ( $(4 \times 13 \times 4 \times 19) + (4 \times 13) + (4 \times 19)$ ). Most importantly, all these approaches do not perform adaptations (changing the mappings and/or DVFS settings) at an application arrival/completion, and performance variations. To this end, this paper presents AdaMD, an adaptive mapping approach coupled with DVFS for performance-constrained multi-threaded applications, executing on heterogeneous multi-cores. AdaMD selects an resource combination (number of cores and their type) that meets the application's performance requirement while minimising energy consump-

tion. This is achieved by employing performance prediction models for resource combination enumeration and selection. Furthermore, the application workload, performance and its status (finished or newly arrived) are monitored for adaptive resource allocation and DVFS. The key contributions of this paper are:

- 1) A performance prediction model that has a maximum percentage error of 8.1%, which is 7.9% lower than the previously reported model [17].
- 2) An online mapping approach that allocates processing cores to application(s) based on performance constraints without using any application-dependent offline results.
- 3) To adapt to application arrival or completion times, and workload/performance variations, an adaptive approach that adjusts the existing thread-to-core mappings and DVFS settings during application execution is presented.
- 4) Experimental validation of the proposed approach on the Odroid-XU3 [1], using several multi-threaded applications from PARSEC [18] and SPLASH [19] benchmarks.

The remainder of this article is organised as follows. Section II presents a motivational example for our work, while section III presents the problem definition for this work. A detailed description of the proposed AdaMD approach is given in Section IV. The experimental setup and validation of our approach are explained in Section V. Section VI discusses the related work and highlights the difference between the proposed approach and exiting works. Finally, Section VII concludes the paper.

## II. MOTIVATION

A heterogeneous computing system with two types of cores, executing multiple performance-constrained applications concurrently, is illustrated in Fig. 1. Dotted squares colored in white/black represent processing cores. For example, such scenarios could be observed when a smartphone user simultaneously runs a music player, Facebook, background email service, downloading a file, etc. As shown in Fig. 1(a), the initial mapping for each application (App1, App2, and App3) is decided based on its performance constraints while considering the energy as an optimization goal. This requires finding an energy-efficient resource combination (number of cores and their type). While these applications are executing, there are primarily three runtime execution scenarios possible: i) any application(s) may finish executing, ii) an application(s) may experience performance degradation due to contention for shared resources, and iii) a new application(s) may arrive into the system. In the first case, if application App1 finishes execution, its resources can be allocated to App2 and App3, which may help them execute faster (and hence put them into a low-power mode sooner), as shown in Fig. 1 (b). This may result in increased performance and lower energy consumption, because power is dissipated for a shorter duration.

For case ii), as reported by previous work [5], [20], applications go through different workload phases during their execution. For example, some workload phases could be more compute-intensive than others or vice versa. Furthermore, in case of concurrent execution, an application may experience

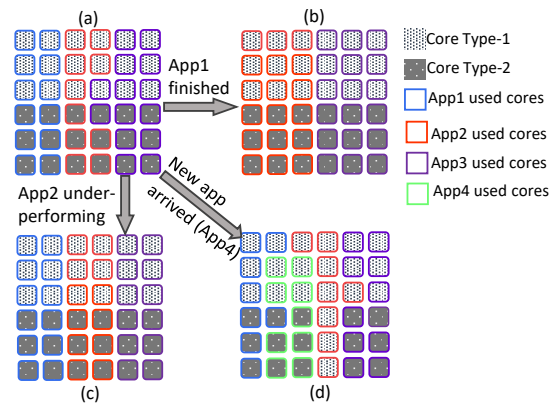


Fig. 1. A motivational example showing three possible runtime execution scenarios (b, c & d) when a system, having two types of cores - Type-1 and Type-2, starts with executing three performance-constrained applications (a). Cores running the same application are encircled with a line of the same color. App1, App2, App3, and App4 represent user applications.

interference from other applications due to shared resources such as Last Level Cache, Memory, etc. All the factors above culminate into variation in an application's workload, subsequently leading to variation in application performance. Therefore, the application's performance has to be monitored periodically, and appropriate action (changing the DVFS setting or remapping) taken to avoid/minimize performance violations. Fig. 1 (c) demonstrates such a case, where more resources are allocated to App2 to mitigate the performance degradation experienced during runtime. If there are no free cores available, as in our case, the cores are taken from the over-performing App3.

For case iii), considering the processing capabilities of the underlying hardware, the user may launch a new application while other applications are running. If all the processing cores have been allocated to the already running applications, the runtime management software should check if there are possibilities to re-adjust the current mapping and allocate resources to the newly arrived application without violating performance constraints. This is shown in Fig. 1 (d), where App4 is added to the system while App1, App2, and App3 are executing. The resources of over-performing applications App1 and App3 are allocated to App4 while keeping the same number of cores for App2.

As discussed before, existing approaches do not consider the above execution scenarios (case i, ii and iii) for adaptation and moreover, they also depend on extensive offline characterisation and/or instrumentation of the chosen applications. As experimentally demonstrated in Section V, adaptation at application arrival and completion, and workload/performance variations would lead to better utilisation of the system resources, and higher energy savings and performance.

## III. PROBLEM FORMULATION

Earlier studies have shown that the thread-to-core mapping problem alone is NP-complete [3]. Therefore, combining it with DVFS would increase the complexity of mapping problem due to the huge design space, thereby making the runtime management significantly inefficient. Similarly, if the number of cores or heterogeneity or frequency levels increases, the



TABLE I  
PARAMETERS USED IN THE PROPOSED APPROACH

Number of Active Cores
Frequency of the Cores
L1 I-Cache Misses
L1 D-Cache Misses
L2 Cache Misses
Instructions Retired
Branch Misses
CPU Cycles
Per Core CPU Utilisation
Memory Reads Per Instruction

types of cores in a platform or by estimating the performance of application for different types of cores by running only on one core type. The former approach requires the migration of the application across various core types. As observed experimentally in [21], migration cost across clusters on a big.LITTLE architecture is relatively high: 2.10 ms to move a thread from a LITTLE cluster to a big cluster, and 3.75 ms to move from a big cluster to a LITTLE cluster. This overhead grows with the number of cores and types. Considering the runtime overheads and scalability, this is not an efficient approach. However, this approach would not need offline analysis as everything is measured at runtime. On the other hand, a performance prediction-based approach avoids thread migration by using the performance models built offline or online. Previous approaches have shown that learning performance models at runtime would make the approach non-scalable and has its overheads in terms of execution time and power [5], [15]. Therefore, AdaMD builds the performance models at design time through a generalized methodology, which can easily be adopted to a new platform/architecture.

**Performance models:** Application performance is usually measured in terms of IPS or IPC, and the relative improvement in the performance is referred to as *speedup*. We define speedup  $\eta$  as

$$\eta = \frac{IPC_{CoreType1}}{IPC_{CoreType2}} \quad (1)$$

where,  $IPC_{CoreType1}$ ,  $IPC_{CoreType2}$  are the IPC of the application achieved on core type-1 and core type-2, respectively. The performance model estimates the speedup, which is used for computing the application performance on a second core type ( $IPC_{CoreType2}$ ), by running the application on one core type and collecting the runtime parameters, and measuring its performance ( $IPC_{CoreType1}$ ) (line 16, Algorithm 1).

To build the performance models, three steps are followed. The first step is identifying the parameters/metrics that capture the most performance-limiting factors by analysing the correlation between various metrics and speedup. Modern processors support monitoring of various architectural events which can be used for analysing the performance, power, etc. However, not all metrics that contribute to performance can be monitored simultaneously due to the limited number of hardware PMCs provided by the platform. For example, on an Odroid-XU3/XU4, the Cortex-A15 processor allows monitoring of seven events, including the cycle counter, at a time. Therefore, metrics that contribute more to the speedup have to be identified. Based on our analysis and the infor-

### Algorithm 1 AdaMD Mapping and Adaptation

**Input:** Applications and performance constraints (Apps)

**Output:**  $\forall$ Apps, mappings and DVFS settings

```

1: PMU_initialize() // initialises PMCs
2: while (1) do
3:   if (NewApp) then
4:     Update the Application Queue 'Apps';
5:     NewApp = 0;
6:   end if
7:   for  $\forall i \in$  Apps do
8:     if (unmapped) then
9:       Allocate a free core 'l' to 'i' and execute;
10:      Measure MRPI and move onto an appropriate core (j);
11:      /*Data collection for performance model*/;
12:      Wait until ROI begins;
13:      pmcs = pmcs_data_collect(j);
14:      f = cpufreq_get_freq_hardware(j);
15:      pmcs.push_back(f);
16:       $\eta$  = speedup_estimate(pmcs, j);
17:      Compute possible resource combinations and resource
        combination with minimum energy  $t_h$  (Eq. (4), (5) &
        (6));
18:      Allocate resources as per  $t_h$ ;
19:    end if
20:  end for
21:  /*Distribute the free cores to active applications*/
22:  Sort the applications by  $\eta$  (list);
23:  while (freecores > 0) do
24:    Increase the resources of app  $i \in$  list by  $y$ ;
25:    freecores = freecores -  $y$ ;
26:     $i++$ ;
27:  end while
28:  /*Application performance and workload adaptation*/
29:  If application workload changes call DVFS(); //Algorithm 2
30:  for  $i \in$  Apps do
31:    if App 'i' under-performs then
32:      Increase frequency or allocate more cores;
33:    end if
34:  end for
35:  /*Application completion detection and adaptation*/
36:  if  $p \in$  Apps finishes then
37:    Distribute freed resources of 'p' to under-performing apps;
38:    Allocate remaining resources to apps equally by sorting
      them based on  $\eta$ ;
39:  end if
40:  /*if stop_governor is set, process exits*/
41:  if (stop_governor) then
42:    PMU_terminate(); //Terminates PMC collection
43:    exit(0);
44:  end if
45: end while

```

mation given in [21], [22], we have identified that cache misses (L1 I/D-Cache & L2 Cache), branch misses, CPU cycles and instructions retired are the appropriate PMCs for estimating the speedup on our chosen platform (listed in Table I). The second step is the collection of characterisation data for a diverse set of applications. As part of this, we have created a diverse set of workloads, containing single and multi-threaded applications from SPEC CPU2006 [23], LMBench [24], RoyLongbottom [25], PARSEC 3.0 [18], SPLASH [19], and MiBench [26]. The Odroid-XU3 platform has four Cortex-

328  
329  
330  
331  
332  
333  
334  
335  
336  
337

338 A7 and four Cortex-A15 cores that can operate at 19 and  
 339 13 different DVFS levels, respectively. For each application,  
 340 data has been collected for every 50 ms at all available  
 341 frequencies on the platform. Furthermore, in the case of  
 342 multi-threaded applications, the number of threads/cores are  
 343 varied from one to four (number of available cores for each  
 344 type). In each case, six PMCs, frequency of the big and  
 345 LITTLE CPUs, execution time of the application on the big  
 346 cluster and LITTLE cluster, and the number of active cores,  
 347 are all used in the modelling. For consistent results, each  
 348 experiment is repeated ten times, and corresponding average  
 349 values are considered while create the model. To create a  
 350 more general approach for deriving performance models, we  
 351 explored several statistical and machine learning techniques.  
 352 Using the open source WEKA workbench [27] to verify the  
 353 relationship between input features/attributes and output/target  
 354 variables. Of all the explored methods, we found that additive  
 355 regression of decision stumps, using boosting for a regression  
 356 problem, resulted in good accuracy as shown in Section V-B.

357 The problem of function estimation usually consists of a  
 358 random *output* variable  $y$  and a set of random *input* features  
 359  $X = \{x_1, x_2, \dots, x_n\}$ . Given a training sample  $\{y_i, X_i\}_1^N$   
 360 of known  $(y, X)$  values, the objective is to identify a function  
 361  $\hat{f}(X)$  that relates  $X$  to  $y$ , such that the expected value  $(E_{y,X})$   
 362 of some specified error function  $\psi(y, f(X))$  is minimized.

$$\hat{f}(X) = \arg \min_{f(X)} E_{y,X} \psi(y, f(X)) \quad (2)$$

363 In general, boosting approximates  $\hat{f}(X)$  by an additive expansion  
 364 of the form, i.e., adding a set of base learners [28], as  
 365 shown below:

$$f(X) = \sum_{k=0}^M \alpha_k h(X; \beta_k) \quad (3)$$

366 Here, the base learner functions  $h(X; \beta)$  are simple functions  
 367 of  $X$  with parameters  $\beta = \{\beta_1, \beta_2, \dots, \beta_M\}$  and  $\{\alpha_k\}_0^M$   
 368 are expansion coefficients. Owing to simplicity, decision stump  
 369 (one-level decision tree) is used as a base learner in our  
 370 work. In brief, additive regression takes an initial guess for  
 371 the speedup (the average speedup observed by all applications  
 372 in the training set) and estimates the speedup by summing positive  
 373 and negative additive-regression factors to  $f_0(X)$ . Each  
 374 additive-regression factor is associated with an input feature  
 375 the factor depends on. As the base learner is a decision stump,  
 376 the input feature is associated with two regression factors, i.e.,  
 377 each of  $\{h(X; \beta_k)\}_0^M$  produces one positive/negative additive-  
 378 regression factor depending on the value of the input feature.  
 379 Additive-regression factors are computed in a forward stage-  
 380 wise manner to minimize the squared error of the predictions  
 381 after  $M$  iterations, which decides the number of base learners.  
 382 Readers can refer to [28] for more details on additive regression.  
 383

384 **3) Resource Combination Enumerator:** For each application,  
 385 a set of all possible resource combinations (number  
 386 of cores and their type) meeting performance constraints has  
 387 to be computed to choose the one that minimizes the overall

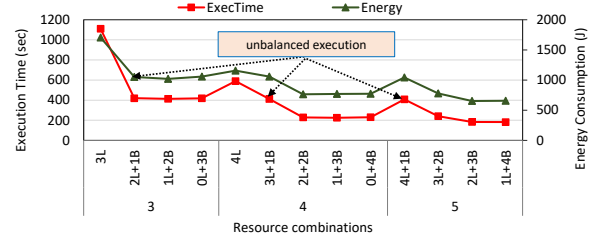


Fig. 3. Energy and execution time at different resource combinations of big (B) and LITTLE (L) for the application *BodyTrack* from PARSEC [18], executing on the Odroid-XU3.

energy consumption (line 17, Algorithm 1). Let  $R$  be the set of  
 possible resource combinations on a platform, and  $PerfApp_i$   
 is the performance constraint for an application  $App_i$ , then the  
 performance meeting thread-to-core mappings ( $T_{map_i}$ ) can be  
 defined as follows:

$$T_{map_i} = \{r \in R \mid perf(r) \leq PerfApp_i\} \quad (4)$$

Here,  $perf(r)$  defines the performance of an application when  
 executed on the resource combination  $r$ . For simplicity, let us  
 take our chosen platform, the Odroid-XU3, with two types of  
 cores: big (B) and LITTLE (L);  $N_b$  and  $N_l$  are set of big and  
 LITTLE cores, respectively. Then,  $perf(r)$  is computed as:

$$perf(r) = n_b \times \eta \times IPC_b + n_l \times IPC_l + IPC_o \quad (5)$$

where,  $\eta = IPC_b/IPC_l$ , performance on the big and  
 LITTLE core is denoted by  $IPC_b$  and  $IPC_l$ , respectively.  
 Furthermore,  $n_b \in N_b$ ,  $n_l \in N_l$  and  $r = n_l \cup n_b$ .  $IPC_o$   
 is the performance overhead incurred when an application is  
 mapped onto cores that do not share a cache. For instance,  
 the big and LITTLE clusters in the Odroid-XU3 do not  
 share caches, which results in an inter-cluster communication  
 overhead when the threads of an application run on both  
 the big and LITTLE clusters. As shown in Equation 5, for  
 our chosen platform with eight cores, near linear speedup is  
 expected with increase in number of cores [29]. Even if there  
 is an error in estimation, this would anyway be compensated  
 by performance monitor (Section IV-B2).

**4) Resource Selector:** The job of *resource selector* is to  
 minimize the energy consumption by selecting a resource combination  
 with minimum energy from the performance meeting  
 thread-to-core mappings  $T_{map_i} = \{3L, 4L, 1L + 1B, \dots\}$ ,  
 where  $L$  and  $B$  refers to big and LITTLE cores, respectively.  
 This can be achieved by selecting a thread-to-core mapping  
 $t_h \in T_{map_i}$  that has the highest performance per watt (PPW)  
 (line 17, Algorithm 1).

$$t_h = \arg \max_{t \in T_{map_i}} PPW(t) \quad (6)$$

where,  $PPW(t)$  is computed as the ratio between IPC  
 achieved for the resource combination ' $t \in T_{map_i}$ ' and its  
 power consumption. This requires measuring the power consumption  
 using on-chip power sensors or employing a power model when a  
 platform does not have power sensors [30]. However, the power  
 model would also require the collection

of various PMCs data at regular intervals of time, and its PMCs may be different than the ones used by performance models [21]. This would need multiplexing PMCs, leading to runtime overheads. To address this, the estimated speedup  $\eta$  can be used as a proxy for identifying the energy-efficient resource combination when power sensors are not available. This is achieved by choosing a resource combination with the ratio between the minimum number of big cores to the minimum number of LITTLE cores ( $C_r$ ) is higher/close to the speedup. As big core can execute  $\eta$  times faster than LITTLE core, above resource selection strategy leads to balanced workload sharing between big and LITTLE cores by executing  $\eta$  times more threads on big than LITTLE. This would lead to efficient utilisation of big cores and supports the balanced execution of an application. For example, if the speedup of an application is  $2\times$ , then the algorithm initially tends to allocate 2-big cores and 1-LITTLE core. This is also demonstrated in Fig. 3, where unbalanced execution resulted in increased execution time and energy consumption. This figure also shows that applications with a speedup greater than one will benefit in terms of energy and performance from allocating more number of cores, as  $C_r$  reaches one or higher.

Furthermore, if  $\eta$  is less than 1, all LITTLE cores are allocated as the application does not benefit from executing on big cores in terms of performance/power. This makes the proposed algorithm effective for single-threaded applications as well, where it maps memory-intensive applications ( $\eta \leq 1$ ) onto LITTLE cores, and compute-intensive ( $\eta > 1$ ) onto big cores. Finally, the output of the resource selector is a resource combination with lower energy consumption and minimum resources that are required for meeting the performance constraints. The information about minimum resources is used by the resource manager.

### B. Resource Manager/Runtime Adaptation

The *Resource Manager*, shown in Fig. 2, is responsible for adapting to application arrival/completion, performance/workload variation, and managing resources at runtime. It consists of the Resource Allocator/Reallocator, Performance Monitor and DVFS governor. These are discussed in detail in the following sections.

1) **Resource Allocator/Reallocator:** The Resource Allocator manages finding free cores and allocating them to the application based on its selected resource combination (line 18, Algorithm 1). This is done by keeping track of allocated cores and free cores available in the platform. The allocated cores are maintained per application, which are used by the performance monitor for measuring application performance and for releasing the resources when the application finishes. While allocating the resources to an application, the resource allocator keeps the knowledge of cores that are leading to over-performance of an application, called extra cores. After finishing the allocation of resources to the applications in application queue (Apps), if there are still free resources available, these are allocated to the running applications if the energy consumption can be minimized by reducing the application execution time. The allocation of extra resources

is done by first creating a sorted list of active applications in descending order of their speedup. Then, application  $i$  at the top of the list is selected, and its allocated cores are increased by one. This process is repeated for remaining applications in the list until no free cores are left (lines 22-27, Algorithm 1). Note that applications with  $\eta < 1$  in the list are given only LITTLE cores as they do not benefit from big cores in terms of energy efficiency.

The Resource Reallocator keeps track of application completion and arrival of new applications into the system. When an application completes execution, it invokes the reallocation routine after releasing the allocated resources (lines 36-39, Algorithm 1). The reallocation routine then distributes the freed resources to the active applications. First, it measures the performance of each application (IPC or IPS) to check if any application is under-performing, i.e., measured performance is lower than the given performance constraint. If an application is under-performing, it then computes the amount of performance loss (the difference between achieved performance and given performance constraint), and then estimates the required resources using Eq. 5 to compensate it. If any resources are remaining after allocating the freed resources to under-performing applications, these resources are distributed among the applications as described in the previous paragraph. As discussed in Section IV-B2 and IV-B3, application performance/workload adaptation is also performed to avoid performance violations as application may experience contention from other applications or workload may change over the time. This may occur at any time during application execution. Therefore, to increase the resource utilisation, free cores are distributed to active applications first. Furthermore, when a new application arrives into the system, the resource reallocator tries to identify and allocate the resources as per  $t_h$  (Eq. 6). This is done by checking if there are enough free resources available in the platform to satisfy the application requirements. In case free resources are not available for meeting performance constraints, the extra cores of over-performing applications are used. After doing this, if the application requirements are still not met, application execution is continued using the available resources until any running application completes and releases allocated resources.

2) **Performance Monitor:** Applications usually exhibit varying workload profiles (e.g., compute-intensive to memory-intensive and vice versa) during execution. When multiple applications are executing simultaneously, the workload profile of each application gets affected due to contention on shared resources [20]. As a result of this, application performance will vary over time, and may lead to the violation of performance constraints. To address this, each application's performance is periodically monitored to detect and compensate when performance constraint is violated (line 30-34, Algorithm 1). An application performance is measured by collecting PMCs corresponding to instructions retired and CPU cycles on all the cores that the application is currently running on. When an application's performance constraint is violated, either the operating frequency is increased, or more cores are allocated. Raising the operating frequency is given priority over assigning more cores as the latter incurs a migration

**Algorithm 2** DVFS governor (DVFS ())

---

```

1:  $MRPI_p = 0, util_p = 0, e_m = 0, e_u = 0;$ 
2: /*Per-core DVFS supporting platforms
Input: for each core 'i',  $MRPI[i]$  and  $f_{req}[i]$ 
Output: voltage-frequency ( $V-f[i]$ ) for next epoch
3:  $pmcs = get\_pmc\_data(i);$ 
4: compute actual MRPI ( $MRPI_a = \frac{instructions\ retired}{L2\ cache\ misses}$ 
5: compute actual utilisation ( $util_a = \frac{active\ CPU\ cycles}{Total\ CPU\ cycles}$ 
6:  $MRPI_p = predict\_mrpi(mrpi_p, mrpi_a, e_m);$ 
7: MRPI prediction error ( $e_m = mrpi_a - mrpi_p;$ 
8:  $util_p = predict\_utilisation(util_p, util_a, e_u);$ 
9: utilisation prediction error ( $e_u = util_a - util_p;$ 
10:  $V-f[i] = bin\_classify(util_p, mrpi_p);$ 
11: if ( $V-f[i] < f_{req}[i]$ ) then
12:    $V-f[i] = f_{req}[i];$ 
13:    $cpufreq\_set\_frequency(i, V-f[i]);$ 
14: end if
15: /*cluster-wide DVFS supporting platforms*/
16: for each cluster 'j' do
17:   Measure MRPI and utilisation of each core  $i \in j;$ 
18:   Compute the minimum MRPI ( $mrpi_a$ ) and utilisation ( $util_a$ );
19:   Repeat steps 6 to 13.
20: end for

```

---

539 overhead which is relatively large compared to the DVFS  
540 transition latency [21]. The operating frequency is increased in  
541 steps of 200 MHz until the performance constraint is satisfied  
542 and this frequency ( $f_{req}$ ) is communicated to DVFS governor  
543 (discussed in the next section) to make sure it does not scale  
544 down the frequency below this value. After the above step,  
545 if any of the applications are still under-performing, as the  
546 last solution, more cores are allocated from the available  
547 free cores or extra cores of over-performing applications.  
548 This allocation is done by computing the performance loss  
549 and corresponding required cores using Eq. 5. As already  
550 explained in Section IV-B1, for applications with  $\eta < 1$ ,  
551 LITTLE cores are preferred over big cores.

552 3) **DVFS governor:** Applications go through different  
553 workload phases (e.g., compute-intensive, memory-intensive,  
554 etc.) and this necessitates choosing a different frequency for  
555 each workload phase to reduce the power consumption while  
556 maintaining application performance within the bounds. For  
557 example, a memory-intensive workload can be executed at  
558 a lower frequency than a compute-intensive workload with  
559 no/negligible performance loss [20]. To this end, AdaMD  
560 adopts the technique proposed in [31], modified to take  $f_{req}$   
561 into account. Algorithm 2 presents the pseudocode of the  
562 DVFS governor.

563 This approach employs a binning-based approach with two  
564 classification layers (line 10). The first layer, consisting of util-  
565 isation bins, classifies the compute-intensity, and the second  
566 layer classifies the memory-intensity using MRPI bins. The  
567 classification bins are computed through an offline analysis  
568 of 81 diverse workloads, including: 25 from SPEC CPU2006  
569 [23], 20 from LMBench [24], 11 from RoyLongbottom [25],  
570 11 from PARSEC 3.0 [18] and 14 from MiBench [26]. For  
571 each application, offline profiling data consisting of MRPI,  
572 utilisation and application performance ( $\frac{1}{Execution\ time}$ ) are  
573 collected at different DVFS settings available on the chosen  
574 platform. The collected utilisation and MRPI for various

575 applications are then grouped into utilisation bins and MRPI  
576 bins, and a corresponding voltage-frequency setting is assigned  
577 to each bin of the second classification layer. At runtime,  
578 the DVFS governor measures the MRPI and utilisation and  
579 uses workload prediction to set an appropriate DVFS level  
580 (lines 3-9). To avoid violation of performance constraints,  
581 the frequency is never scaled down below  $f_{req}$  (lines 11-14).  
582 Workload prediction is based on exponential moving average  
583 filter. Prediction error during previous time epoch for MRPI  
584 ( $e_m$ ) and utilisation ( $e_u$ ) is used as feedback to improve the  
585 workload prediction accuracy (lines 7 & 9). Furthermore,  
586 it can manage both per-core (lines 2-14), i.e., supporting  
587 fine-grained power management [32], and cluster-wide DVFS  
588 platforms (lines 15-20). For more details on binning-based  
589 DVFS approach, readers can refer to [31], [33].

## V. EXPERIMENTAL RESULTS

591 This section presents the details of the experimental setup,  
592 covering the platform, benchmark applications and reported  
593 approaches considered for the comparison. Furthermore, an  
594 evaluation of the performance prediction models and benefits  
595 of the AdaMD approach over the previous approaches are  
596 discussed, including associated overheads.

## A. Experimental Setup

597 *Platform:* We use the Odroid-XU3 [1], containing the ARM  
598 big.LITTLE technology based Samsung Exynos 5422 chip.  
599 This has four ARM Cortex-A15 (big) cores, four ARM Cortex-  
600 A7 (LITTLE) cores. The platform supports per-cluster DVFS,  
601 and all cores within a cluster can only run at the same DVFS  
602 level. The big cores have a range of frequencies between 0.2  
603 GHz and 2.0 GHz with a 0.1 GHz step, whereas the LITTLE  
604 cores can vary their frequencies from 0.2 GHz to 1.4 GHz in  
605 steps of 0.1 GHz. The device firmware automatically adjusts  
606 the voltage for a selected frequency. The platform also contains  
607 four real-time current sensors that facilitate measurement of  
608 power consumption of each CPU cluster, GPU and memory.  
609 We used Ubuntu OS with kernel version 3.10.96. Energy  
610 consumption is computed as the product of average power con-  
611 sumption (dynamic and static) and application execution time.  
612 This includes both the core and memory energy consumption  
613 of all the software components, including our implementation,  
614 OS, applications and other background processes.

615 *Implementation:* The proposed AdaMD approach is imple-  
616 mented as a *user space* application by using the Perfmon2  
617 [34] and cpufrequtils framework. Perfmon2 en-  
618 ables the *user space* access to the performance moni-  
619 toring unit (PMU), and cpufrequtils helps in set-  
620 ting/getting the operating frequencies. Standard Linux API  
621 ( $sched\_setaffinity(2)$ ) is used to control the CPU  
622 affinity of processes, i.e., to bind the applications to specific  
623 cores. The thread-to-core mapping algorithm operates at a  
624 coarser granularity (500 ms) considering its higher migration  
625 overhead. As the workload of application changes randomly,  
626 to capitalize on these changes for energy savings, the DVFS  
627 governor is operated at a finer granularity of 100 ms.  
628

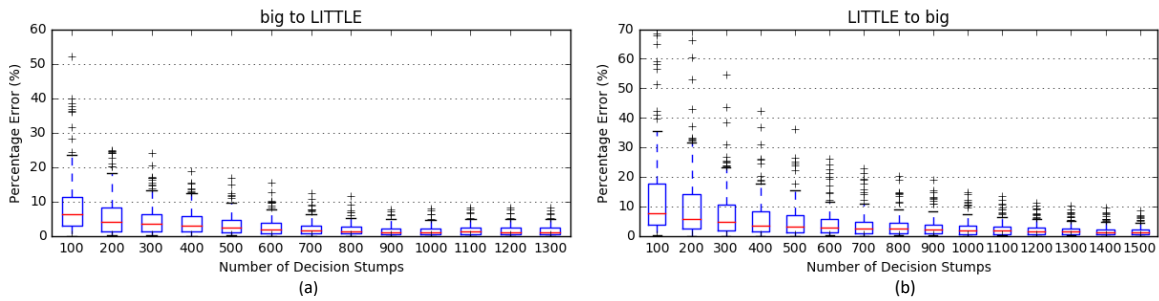


Fig. 4. Box plot of absolute percentage error in IPC prediction by our performance model for different number of decision stumps used in the additive regression, showing the median, lower quartile, upper quartile and outliers – (a) Estimating the performance of LITTLE given the information about the big core (b) Estimating the performance of big given the information about the LITTLE core.

629 **Applications:** To evaluate AdaMD, applications – Blacksc-  
 630 holes (bl), Bodytrack (bo), Swaptions (sw), Freqmine (fr), Vips  
 631 (vi), Water-Spatial (wa), Raytrace (ra), fmm (fm) – from pop-  
 632 ular benchmark suites, such as PARSEC 3.0 [18] and SPLASH  
 633 [19], are taken. These applications exhibit different memory  
 634 behavior, data partitions, and data sharing patterns. Different  
 635 execution scenarios – single application, concurrent execution  
 636 of multiple applications, dynamic addition of application(s) at  
 637 runtime – are also considered to mimic the real-world behav-  
 638 ior. To ensure the deterministic execution of application and  
 639 to meet its performance constraint, no two applications share  
 640 the same cores. However, the threads of the same application  
 641 share the allocated cores to maximize resource utilisation.  
 642 For each application, performance constraints are defined in  
 643 terms IPC. Such performance requirements can be translated  
 644 to throughput requirements for frame based applications like  
 645 audio/video applications, where throughput is expressed as a  
 646 frame rate to guarantee a good user experience.

647 **Comparison:** To show the benefits of our approach AdaMD  
 648 compared to the state-of-the-art, the selected comparison can-  
 649 didates from the relevant reported works are given below.

- 650 1) HMP+ $x$  [35]: The state-of-the-art solution for  
 651 big.LITTLE multi-processing, the Heterogeneous Multi-  
 652 Processing (HMP) scheduler, with various default Linux  
 653 power governors  $x$  (= Ondemand (O), Conservative (C)  
 654 and Interactive (I)) is considered. For a fair comparison,  
 655 we ran applications with different numbers of threads  
 656 and chose the one meeting the performance constraint.
- 657 2) MIM [14]: This approach maps application threads onto  
 658 only one type of core(s) based on workload memory-  
 659 intensity, called a memory-intensity based mapping  
 660 (MIM). For the single-application execution scenario, a  
 661 memory-intensive application is mapped onto LITTLE  
 662 cores, whereas a compute-intensive one is executed on the  
 663 big cores. In a multiple-application scenario, applications  
 664 are sorted based on their memory-intensity, and the  
 665 one with the highest memory-intensity is mapped onto  
 666 LITTLE cores, and remaining applications are allocated  
 667 onto the big cluster with an equal number of cores.
- 668 3) EAM [15]: An energy-efficient mapping is selected  
 669 through an exhaustive search of voltage-frequency set-  
 670 tings and thread-to-core mappings. For each possible  
 671 thread-to-core mapping, voltage-frequency settings are  
 672 varied from the lowest possible value to the highest and

673 the one that meets performance requirement with the  
 674 lowest energy consumption is chosen. We refer to this  
 675 approach as energy-aware mapping (EAM).

- 676 4) ITMD [6]: This approach uses offline analysis of energy  
 677 and performance for individual applications to decide on  
 678 an energy-efficient mapping when multiple applications  
 679 are run concurrently. Furthermore, it also applies work-  
 680 load classification-based DVFS periodically to minimize  
 681 the power consumption.

### 682 B. Evaluation of Performance Predictor

683 The performance prediction model estimates the perfor-  
 684 mance of the big core given the performance of a LITTLE  
 685 core ( $P_{bl}$ ) and vice versa ( $P_{lb}$ ). The number of base learners  
 686 (decision stumps)  $M$  in Eq. 3 impacts the model accuracy  
 687 and runtime overhead. We tested our model over 148 distinct  
 688 samples to evaluate the model accuracy in IPC estimation and  
 689 the corresponding box plot of percentage error distribution  
 690 for  $P_{bl}$  and  $P_{lb}$  are given in Figures 4a and 4b respectively.  
 691 As shown, the error range gets narrower with the number  
 692 of decision stumps, as it would help in better predicting  
 693 the speedup. Furthermore, increasing the number of decision  
 694 stumps also reduces the outliers, shown as cross in Figures 4a  
 695 and 4b, improving model stability. However, choosing more  
 696 decision stumps could increase the runtime overhead, and  
 697 sometimes accuracy of the prediction may not be improved  
 698 after reaching a certain number of decision stumps. There-  
 699 fore, to balance this, we built additive regression models for  
 700 different numbers of decision stumps. It can be seen from  
 701 Fig. 4a and 4b that the the improvement in model accuracy  
 702 is negligible after 900 and 1100 decision stumps for  $P_{bl}$  and  
 703  $P_{lb}$  respectively. Therefore, we have chosen these numbers  
 704 for our models  $P_{bl}$  (mean absolute percentage error (MAPE)  
 705 = 1.57%; maximum error (ME) = 8.1%) and  $P_{lb}$  (MAPE =  
 706 3.45%; ME = 8.5%). The maximum error of  $P_{bl}$  and  $P_{lb}$  is  
 707 about 7.9% and 5% lower compared to the previous model  
 708 [17], respectively. The prediction accuracy of  $P_{lb}$  is 1.88%  
 709 worse than  $P_{lb}$  and requires 200 extra decision stumps. This is  
 710 because the LITTLE cores support accessing only four PMCs  
 711 simultaneously, compared to six PMCs supported by big cores.

### 712 C. Comparison of Energy Consumption

713 This section presents the energy consumption results for var-  
 714 ious approaches to show the benefits of the proposed AdaMD



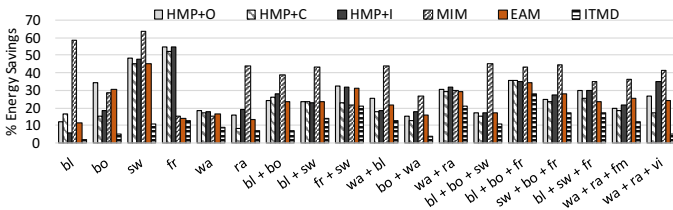


Fig. 5. Percentage improvement in energy consumption achieved by the AdaMD compared to reported approaches for single and concurrent applications.

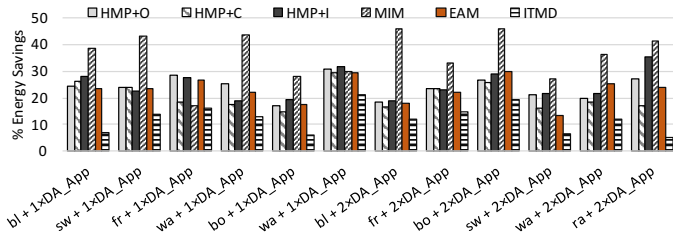


Fig. 6. Energy savings achieved by the AdaMD with respect to different approaches for one and two applications added dynamically to the system while an application is executing.

715 approach. Fig. 5 shows the energy savings achieved by the  
 716 AdaMD with respect to reported approaches for different  
 717 single and concurrently executing applications (launched at the  
 718 same time). We observed substantial energy savings compared  
 719 to the reported approaches for all the application execution  
 720 scenarios. For single application execution (bl, bo, sw, fr, wa,  
 721 and ra), with our approach AdaMD, average energy savings of  
 722 30.7%, 25.8%, 27.3%, 37.4%, 21.8% and 7.8% are observed  
 723 compared to HMP+O, HMP+C, HMP+I, MIM, EAM, and  
 724 ITMD, respectively. Furthermore, for concurrent execution of  
 725 two and three applications, AdaMD shows 25.5%, 22.4%,  
 726 26.5%, 37.5%, 24.8%, and 14.2% lower energy consumption  
 727 than HMP+O, HMP+C, HMP+I, MIM, EAM, and ITMD,  
 728 respectively. In the single application scenario, we observed  
 729 that ITMD, EAM, and AdaMD chooses a similar thread-  
 730 to-core mapping, however, the energy savings observed are  
 731 mainly because of the proposed DVFS technique. Unlike,  
 732 ITMD and EAM, AdaMD takes the thread synchronisation  
 733 overhead into account while selecting a voltage-frequency  
 734 setting. In concurrent execution scenarios, the energy savings  
 735 are due to both DVFS and the utilisation of freed resources of  
 736 a finished application for active applications.

737 Furthermore, to demonstrate the adaptiveness of AdaMD  
 738 to application arrival, the following experimental evaluation is  
 739 performed. The execution starts with one application and later,  
 740 one (+ 1xDA\_Apps) or two (+ 2xDA\_Apps) applications  
 741 are added at runtime. The dynamically added applications,

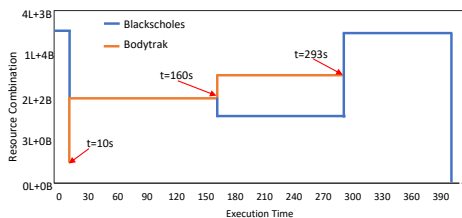


Fig. 7. Resource combination (number of big (B) and LITTLE (L) cores) allocated to Blackscholes and Bodytrak by the proposed AdaMD approach to adapt to application arrival/completion and performance variation.

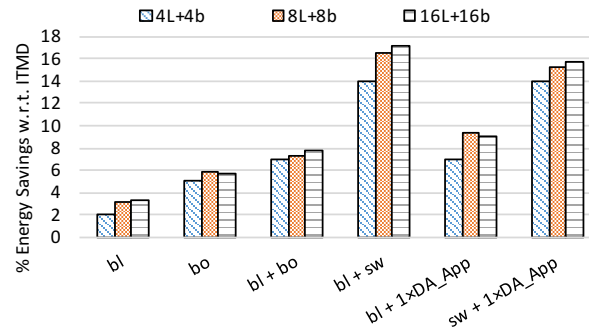


Fig. 8. Scalability of AdaMD for different core configurations of big (b) and LITTLE (L) cores: energy savings achieved by AdaMD with respect to ITMD.

742 abbreviated as DA\_Apps in Fig. 6, are from those mentioned  
 743 in Section V-A. The advantages of the AdaMD with respect  
 744 to other approaches in terms of energy consumption are  
 745 shown in Fig. 6. On an average, AdaMD reduces the energy  
 746 consumption by 23.8%, 20.6%, 24.8%, 35.8%, 12.2%, and  
 747 23.0% compared to HMP+O, HMP+C, HMP+I, MIM, EAM,  
 748 and ITMD, respectively. To illustrate AdaMD's ability to  
 749 adapt to different runtime scenarios, we plotted the resource  
 750 combination (number of active core and their type) versus  
 751 execution time for Blackscholes and Bodytrak in Fig. 7.  
 752 While Blackscholes is executing with four LITTLE and three  
 753 big cores (4L+3B), Bodytrak is added to the system at t=10s.  
 754 Considering the performance constraints of Bodytrak, 2B+2L  
 755 are allocated to Bodytrak by freeing the cores from over-  
 756 performance of Blackscholes. Due to the workload variations,  
 757 Bodytrak experiences performance loss at t=160s, thereby  
 758 triggering the Resource Reallocator to readjust the mappings  
 759 of both Blackscholes and Bodytrak. Upon Bodytrak's com-  
 760 pletion at t=293s, the freed cores are again allocated to  
 761 Blackscholes as it can benefit from faster execution to lower  
 762 the energy consumption.

763 Fig. 8 demonstrates the scalability of AdaMD, showing  
 764 energy savings with respect to ITMD for two different core  
 765 configurations (8L+8b, 16L+16b). The reported values in  
 766 the figure have been obtained through analytical analysis of  
 767 experimental results (performance and energy) collected on  
 768 the Odroid-XU3 (4L+4b) and extrapolating for the considered  
 769 application execution scenarios. We used linear extrapolation  
 770 that takes runtime overheads associated with each application  
 771 as it varies depending upon workload characteristics (e.g.,  
 772 frequent workload variations may incur DVFS transition la-  
 773 tencies/thread migration overheads). As can be seen, AdaMD  
 774 is able to adapt to increased design space and achieve energy  
 775 savings. The increase in energy savings is mainly due to  
 776 proposed DVFS which exploits the synchronisation overheads  
 777 and workload variations to lower power consumption of more  
 778 number of active cores.

D. Performance

780 The proposed approach outperforms all reported approaches  
 781 in meeting application performance constraints, as shown in  
 782 Fig. 9. We evaluated the percentage of performance constraint  
 783 misses for all the application scenarios presented in Fig. 5  
 784 (Without Application Addition) and Fig. 6 (With Application

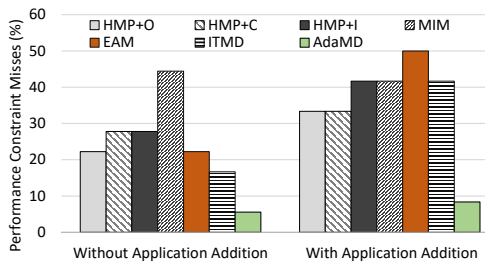


Fig. 9. Evaluation of various approaches in meeting application performance constraints.

815 Addition). For the without application addition case, AdaMD  
 816 meets application performance constraint for 95% of the con-  
 817 sidered application scenarios, i.e., 17 out of 18 cases, shown  
 818 on the horizontal axis in Fig. 5. The only scenario where  
 819 the AdaMD fails to satisfy the performance requirements is  
 820  $bl+sw+fr$ . Even in this case, except for  $sw$ , the performance  
 821 constraints of  $bl$  and  $fr$  are met. This is mainly because  
 822 of the diverse workload profiles of the three applications  
 823 and relatively higher performance requirements chosen for  
 824  $sw$  (approximately  $2\times$  compared to  $bl$  and  $fr$ ). In case  
 825 of application addition at runtime, AdaMD is able to satisfy  
 826 performance constraints for 92% of the evaluated scenarios,  
 827 i.e., out of 12 scenarios shown in the Fig. 6, except for  $sw +$   
 828  $2\times DA\_Apps$ , the performance constraints are met.

829 Compared to the recently reported approach ITMD [6],  
 830 AdaMD achieves energy savings of up to 28% (for  
 831  $bl+bo+fr$ ). Further, AdaMD satisfies performance con-  
 832 straints for up to 95% of the application scenarios (80% better  
 833 than ITMD).  
 834

### 804 E. Runtime Overheads

805 To compute the runtime overheads of the AdaMD, we  
 806 measured the amount of time that the algorithm takes to  
 807 complete various steps (A to E) explained in Section IV. Steps  
 808 A (Runtime Data Collector), B (Performance Predictor), C  
 809 (Resource Combination Enumerator) and D (Resource Selec-  
 810 tor) are triggered when an application arrives into the system,  
 811 whereas step E (Resource Manager) operates periodically. The  
 812 runtime overheads can be analytically represented as follows  
 813 for each time epoch (500 ms):

$$T_o = T_{AdaMap} + \eta \times T_{DVFS} \quad (7)$$

814 where,

$$T_{AdaMap} = T_{pmc_m} + T_{pm} + T_{th} + T_{rar} \quad (8)$$

$$T_{DVFS} = T_{pmc_{vf}} + T_{metrics} + T_{wp} + T_{classify} + T_{vfs} \quad (9)$$

815 where,  $T_{pmc_m}$ ,  $T_{pm}$ ,  $T_{th}$ ,  $T_{rar}$ ,  $T_{pmc_{vf}}$ ,  $T_{metrics}$ ,  $T_{wp}$ ,  
 816  $T_{classify}$ , and  $T_{vfs}$  represent time taken for PMC data col-  
 817 lection for mapping; performance prediction; identification of  
 818 resource combination; resource allocation/reallocation; PMC  
 819 data collection for DVFS; computation of MRPI and util-  
 820 isation; workload prediction; finding DVFS setting through

821 classification bins; and DVFS transition latency, respectively.  
 822 Note that performance prediction happens only when a new  
 823 application is launched, therefore the overhead  $T_{pm}$  may  
 824 not be present in every time epoch. Moreover, the runtime  
 825 overhead  $T_{DVFS}$  is multiplied by a factor of 2.5 ( $\eta$ ), as it  
 826 operates at a finer granularity of 100 ms compared to the  
 827 mapping time interval of 500 ms.

828 We observed an average runtime overhead of  $600\mu s$  and  
 829 1.4 ms for A to D when executed at 2 GHz and 1 GHz  
 830 on a big core of Odroid-XU3, respectively. The DVFS part  
 831 of step E incurs  $320\mu s$  and others parts take up to  $15\mu s$   
 832 when the overhead is measured at the maximum frequency  
 833 (2 GHz). The DVFS algorithm operates at a granularity of  
 834 100 ms, so the overhead is less than 0.5%. Performance and  
 835 Resource manager part of E is invoked for every 500 ms. The  
 836 overhead associated with this part depends on the number of  
 837 times the application misses its performance constraint and  
 838 thread migrations across the cores. Here, we observed an  
 839 overhead between 0.15% to 0.75%. Our results show that the  
 840 total runtime overhead is very minimal and moreover, they  
 841 have been included when computing energy consumption and  
 842 performance.

## 843 VI. RELATED WORK

844 To achieve energy savings and/or to meet performance  
 845 constraints in multi-core platforms, various approaches for  
 846 DVFS and/or task mapping have been proposed [3]–[17], [20],  
 847 [31], [36]–[41]. These works perform offline, online or hybrid  
 848 (offline & online) optimization for resource management.

849 Approaches based on offline optimization utilize extensive  
 850 design space exploration of the underlying hardware and  
 851 target application(s). The techniques proposed in [7], [40]  
 852 are used for DVFS and/or task mapping. However, they  
 853 consider execution of a single application at a time, and thus  
 854 are not suitable for the concurrent execution of applications.  
 855 The approach presented in [40] generates multiple mappings  
 856 for each application offering a tradeoff between resource  
 857 requirements and throughput, while Quan and Pimentel [8]  
 858 proposed scenario-based online mapping approaches targeting  
 859 homogeneous multi-core platforms in which mappings derived  
 860 from design-time DSE are stored for runtime mapping deci-  
 861 sions. Evidently, these techniques consume more time, and  
 862 cannot cope with dynamic application behavior, especially  
 863 when multiple applications are run concurrently.

864 To adapt to dynamic application workloads, pure online  
 865 optimization based approaches, performing all processing at  
 866 runtime, have also been investigated [4], [9]–[11]. In [4], an  
 867 online reinforcement learning based adaptive DVFS approach  
 868 targeting frame-based applications is presented to improve  
 869 energy efficiency. In [9], an online spatial mapping technique  
 870 to map streaming applications onto a multi-core system is  
 871 discussed. Brião *et al.* [10] present dynamic task allocation  
 872 strategies based on bin-packing algorithms for soft real-time  
 873 applications. An online task allocator using the adaptive task  
 874 allocation algorithm and clustering approach for minimizing  
 875 the communication load is described in [11]. All of these  
 876 approaches perform well for unknown applications to be exe-  
 877 cuted at runtime, but lead to inefficient results as optimization

878 decisions need to be taken quickly without offline analysis  
879 results [3].

880 Hybrid approaches using results of offline analysis in  
881 making online decisions have been widely proposed to im-  
882 prove energy efficiency/performance in homogeneous multi-  
883 core platforms [8], [12]–[17]. Such approaches usually achieve  
884 better performance/energy savings compared to pure online  
885 optimizations as they take advantage of both offline and online  
886 computation. In [12], task mapping and DVFS under power  
887 constraints are discussed. Similarly, in [13], first thread-to-core  
888 mapping is obtained based on utilization, and then DVFS is  
889 applied depending upon the power budget. When considering  
890 the power-performance tradeoffs, recent research focus has  
891 shifted to heterogeneous architectures [3], [6], [14]–[17]. For  
892 multi-threaded applications, most approaches tend to map an  
893 application completely onto one type of processing core(s)  
894 [14], [16], [17]. This simplifies the thread-to-core mapping  
895 problem, but cannot benefit from the power-performance trade-  
896 offs offered by simultaneously mapping application threads  
897 onto multiple types of cores. Van Craeynest *et al.* [14] pre-  
898 sented a performance impact estimation technique to predict  
899 which application-to-core mapping is likely to provide the best  
900 performance to map the application onto the most appropriate  
901 core type. In a similar direction, some proposals have used  
902 workload memory-intensity as an indicator to guide task  
903 mapping [38], [39]. A domain-specific hybrid task mapping is  
904 presented in [3], which relies heavily on offline DSE. However,  
905 approaches reported in [3], [14] do not consider DVFS which  
906 can help to improve energy savings.

907 On the other hand, techniques proposed in [5], [6], [15]–[17]  
908 use DVFS, but they have several shortcomings. For example,  
909 in [16], the design space is explored for a single application,  
910 which increases exponentially for concurrent execution of ap-  
911 plications. Donyanavard *et al.* [17] consider applications with  
912 only one thread and thus use only one type of core for each ap-  
913 plication. The approach presented in [15] considers concurrent  
914 execution and mapping of application threads onto more than  
915 one type of cores. However, it requires extensive offline and/or  
916 online exploration for building regression models for perfor-  
917 mance and energy for all possible thread-to-core mappings and  
918 voltage-frequency settings, which is non-scalable. Moreover,  
919 online periodic adjustment of  $V$ - $f$  setting is not explored, which  
920 is essential for adapting to workload variations and achieving  
921 better energy savings. This has been addressed in [5], [6],  
922 however, they also require extensive offline characterisation,  
923 and in particular, [5] requires application instrumentation to  
924 guide the runtime selection. Moreover, all these approaches  
925 do not perform adaptive mapping at application arrival/exit,  
926 and thus they are not efficient if a new/unknown application  
927 arrives/existing application finishes. The approach (AdaMD)  
928 presented in this paper addresses the above limitations by  
929 removing dependency on the application-dependent offline  
930 results, and adapting to application arrival/completion times.

## 931 VII. CONCLUSIONS

932 The increasing demand for performance and energy effi-  
933 ciency has forced mobile systems to employ heterogeneous

934 multiprocessor system-on-chips. These systems offer a diverse  
935 set of core and frequency configurations to runtime manage-  
936 ment systems for online tuning. This paper has presented an  
937 adaptive thread-to-core mapping and DVFS technique, called  
938 AdaMD, for choosing a configuration for each performance-  
939 constrained application that minimises energy consumption.  
940 By using runtime information while applications are executing  
941 and eliminating the need for application-dependent offline  
942 results, AdaMD is capable of managing even unknown appli-  
943 cations efficiently. Proposed algorithm first selects a resource  
944 combination (number of cores and their type) that meets the  
945 application performance requirement using an accurate perfor-  
946 mance prediction model and resource enumerator/selector. It  
947 then monitors application performance, workload and its status  
948 (finished or newly arrived) for tuning voltage-frequency set-  
949 tings and adjusting thread-to-core mappings. Our experiments  
950 show an improvement of up to 28% in energy consumption  
951 compared to the most promising existing approaches. The  
952 proposed approach also outperforms previous approaches in  
953 meeting application performance constraints. Our future work  
954 includes validation with more number of cores and types  
955 having different ISA (e.g., CPU, GPU, etc.) to show the  
956 scalability and adaptability of the approach.

## 957 ACKNOWLEDGEMENT

958 This work was supported in parts by the EPSRC  
959 Grant EP/L000563/1 and the PRiME Programme Grant  
960 EP/K034448/1 ([www.prime-project.org](http://www.prime-project.org)). Experimental data  
961 used in this paper can be found at [https://doi.org/10.5258/](https://doi.org/10.5258/SOTON/D1041)  
962 SOTON/D1041.

## 963 REFERENCES

- 964 [1] “Odroid-XU3,” [www.hardkernel.com/main/products](http://www.hardkernel.com/main/products).
- 965 [2] “Mediatek X20,” <http://www.96boards.org/product/mediatek-x20/>.
- 966 [3] W. Quan and A. D. Pimentel, “A hybrid task mapping algorithm for  
967 heterogeneous MPSoCs,” *ACM Transactions on Embedded Computing  
968 Systems*, vol. 14, no. 1, p. 14, 2015.
- 969 [4] R. A. Shafik, A. K. Das, L. A. Maeda-Nunez, S. Yang, G. V. Merrett, and  
970 B. Al-Hashimi, “Learning transfer-based adaptive energy minimization  
971 in embedded systems,” *IEEE Transactions on Computer-Aided Design  
972 of Integrated Circuits and Systems*, vol. 35, no. 6, pp. 877–890, 2016.
- 973 [5] U. Gupta, C. A. Patil, G. Bhat, P. Mishra, and U. Y. Ogras, “DyPO:  
974 Dynamic pareto-optimal configuration selection for heterogeneous MP-  
975 SoCs,” *ACM Transactions on Embedded Computing Systems*, vol. 16,  
976 no. 5s, p. 123, 2017.
- 977 [6] B. K. Reddy, A. K. Singh, D. Biswas, G. V. Merrett, and B. M.  
978 Al-Hashimi, “Inter-cluster thread-to-core mapping and dvfs on hetero-  
979 geneous multi-cores,” *IEEE Transactions on Multi-Scale Computing  
980 Systems*, vol. 4, no. 3, pp. 369–382, 2018.
- 981 [7] M. Goraczko, J. Liu, D. Lymberopoulos, S. Matic, B. Priyantha,  
982 and F. Zhao, “Energy-optimal software partitioning in heterogeneous  
983 multiprocessor embedded systems,” in *Proc. of the Design Automation  
984 Conference*. ACM, 2008, pp. 191–196.
- 985 [8] W. Quan and A. D. Pimentel, “A scenario-based run-time task mapping  
986 algorithm for MPSoCs,” in *Proc. of the Design Automation Conference*.  
987 ACM, 2013, p. 131.
- 988 [9] P. K. Hölzenspies, J. L. Hurink, J. Kuper, and G. J. Smit, “Run-time  
989 spatial mapping of streaming applications to a heterogeneous multi-  
990 processor system-on-chip (MPSoC),” in *Design, Automation and Test  
991 in Europe*. ACM, 2008, pp. 212–217.
- 992 [10] E. W. Brião, D. Barcelos, and F. R. Wagner, “Dynamic task allocation  
993 strategies in MPSoC for soft real-time applications,” in *Design, Automa-  
994 tion and Test in Europe*. ACM, 2008, pp. 1386–1389.
- 995 [11] J. Huang, A. Raabe, C. Buckl, and A. Knoll, “A workflow for run-  
996 time adaptive task allocation on heterogeneous MPSoCs,” in *Design,  
997 Automation and Test in Europe*, 2011.

- 998 [12] R. Cochran, C. Hankendi, A. K. Coskun, and S. Reda, "Pack & cap:  
999 adaptive DVFS and thread packing under power caps," in *Proc. of the*  
1000 *IEEE/ACM Intl. symposium on microarchitecture*, 2011, pp. 175–185.
- 1001 [13] H. Sasaki, S. Imamura, and K. Inoue, "Coordinated power-performance  
1002 optimization in manycores," in *Proc. of the Intl. Conf. on Parallel*  
1003 *architectures and compilation techniques*. IEEE, 2013, pp. 51–61.
- 1004 [14] K. Van Craeynest, A. Jaleel, L. Eeckhout, P. Narvaez, and J. Emer,  
1005 "Scheduling heterogeneous multi-cores through performance impact  
1006 estimation (PIE)," in *ACM SIGARCH Computer Architecture News*,  
1007 vol. 40, no. 3, 2012, pp. 213–224.
- 1008 [15] A. Aalsaud, R. Shafik, A. Rafiev, F. Xia, S. Yang, and A. Yakovlev,  
1009 "Power-aware performance adaptation of concurrent applications in  
1010 heterogeneous many-core systems," in *Intl. Symp. on Low Power Elec-*  
1011 *tronics and Design*. ACM, 2016, pp. 368–373.
- 1012 [16] E. D. Sozzo, G. C. Durelli, E. Trainiti, A. Miele, M. D. Santambrogio,  
1013 and C. Bolchini, "Workload-aware power optimization strategy for  
1014 asymmetric multiprocessors," in *Design, Automation & Test in Europe*  
1015 *Conference & Exhibition*. IEEE, 2016, pp. 531–534.
- 1016 [17] B. Donyanavard, T. Mück, S. Sarma, and N. Dutt, "SPARTA: runtime  
1017 task allocation for energy efficient heterogeneous many-cores," in *Proc.*  
1018 *of the Intl. Conf. on Hardware/Software Codesign and System Synthesis*.  
1019 ACM, 2016, p. 27.
- 1020 [18] C. Bienia, "Benchmarking modern multiprocessors," Ph.D. dissertation,  
1021 Princeton University, January 2011.
- 1022 [19] S. C. Woo, M. Ohara, E. Torrie, J. P. Singh, and A. Gupta, "The  
1023 SPLASH-2 programs: Characterization and methodological considera-
- 1024 tions," in *ACM SIGARCH Computer Architecture News*, vol. 23, no. 2,  
1025 1995, pp. 24–36.
- 1026 [20] B. K. Reddy, G. V. Merrett, B. M. Al-Hashimi, and A. K. Singh, "Online  
1027 concurrent workload classification for multi-core energy management,"  
1028 in *Design, Automation Test in Europe Conference & Exhibition*, March  
1029 2018, pp. 621–624.
- 1030 [21] M. Pricopi, T. S. Muthukaruppan, V. Venkataramani, T. Mitra, and  
1031 S. Vishin, "Power-performance modeling on asymmetric multi-cores,"  
1032 in *Intl. Conf. on Compilers, Architecture and Synthesis for Embedded*  
1033 *Systems*. IEEE, 2013, pp. 1–10.
- 1034 [22] J. C. Saez, A. Fedorova, D. Koufaty, and M. Prieto, "Leveraging core  
1035 specialization via OS scheduling to improve performance on asymmetric  
1036 multicore systems," *ACM Transactions on Computer Systems*, vol. 30,  
1037 no. 2, p. 6, 2012.
- 1038 [23] J. L. Henning, "Spec cpu2006 benchmark descriptions," *ACM SIGARCH*  
1039 *Computer Architecture News*, vol. 34, no. 4, pp. 1–17, 2006.
- 1040 [24] L. McVoy and C. Staelin, "Lmbench: Portable tools for performance  
1041 analysis," in *USENIX Annual Technical Conference*, 1996, pp. 23–23.
- 1042 [25] "Roy Longbottom's PC Benchmark Collection," <http://www.roylongbottom.org.uk>. [Online; accessed 10-Oct-2018].
- 1043 [26] M. R. Guthaus, J. S. Ringenberg, D. Ernst, T. M. Austin, T. Mudge,  
1044 and R. B. Brown, "MiBench: A free, commercially representative  
1045 embedded benchmark suite," in *IEEE Intl. Workshop on Workload*  
1046 *Characterization*. IEEE, 2001, pp. 3–14.
- 1047 [27] F. Eibe, M. Hall, and I. Witten, "The WEKA workbench. online  
1048 appendix for" data mining: Practical machine learning tools and  
1049 techniques," *Morgan Kaufmann*, 2016.
- 1050 [28] J. H. Friedman, "Stochastic gradient boosting," *Computational Statistics*  
1051 *& Data Analysis*, vol. 38, no. 4, pp. 367–378, 2002.
- 1052 [29] G. Southern and J. Renaud, "Analysis of parsec workload scalability,"  
1053 in *2016 IEEE International Symposium on Performance Analysis of*  
1054 *Systems and Software (ISPASS)*. IEEE, 2016, pp. 133–142.
- 1055 [30] B. K. Reddy, M. J. Walker, D. Balsamo, S. Diestelhorst, B. M. Al-  
1056 Hashimi, and G. V. Merrett, "Empirical CPU power modelling and  
1057 estimation in the gem5 simulator," in *IEEE Intl. Symp. on Power and*  
1058 *Timing Modeling, Optimization and Simulation*, 2017, pp. 1–8.
- 1059 [31] K. R. Basireddy, E. W. Wachter, B. M. Al-Hashimi, and G. V. Merrett,  
1060 "Workload-aware runtime energy management for HPC systems," in *Intl.*  
1061 *Conf. on High Performance Computing & Simulation*, 2018, p. 8.
- 1062 [32] U. Y. Ogras, R. Marculescu, D. Marculescu, and E. G. Jung, "Design and  
1063 management of voltage-frequency island partitioned networks-on-chip,"  
1064 *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*,  
1065 vol. 17, no. 3, pp. 330–341, 2009.
- 1066 [33] G. Liu, J. Park, and D. Marculescu, "Dynamic thread mapping for high-  
1067 performance, power-efficient heterogeneous many-core systems," in *Intl.*  
1068 *Conf. on Computer Design*. IEEE, 2013, pp. 54–61.
- 1069 [34] S. Eranian, "Perfmon2: a flexible performance monitoring interface for  
1070 linux," in *Proc. of Ottawa Linux Symp.*, 2006, pp. 269–288.
- 1071 [35] K. Yu, D. Han, C. Youn, S. Hwang, and J. Lee, "Power-aware task  
1072 scheduling for big.LITTLE mobile processor," in *Intl. SoC Design Conf.*  
1073 IEEE, 2013, pp. 208–212.
- 1074 [36] A. K. Singh, C. Leech, B. K. Reddy, B. M. Al-Hashimi, and G. V.  
1075 Merrett, "Learning-based run-time power and energy management of  
1076 multi/many-core systems: current and future trends," *Journal of Low*  
1077 *Power Electronics*, vol. 13, no. 3, pp. 310–325, 2017.
- 1078 [37] A. K. Singh, A. Prakash, K. R. Basireddy, G. V. Merrett, and B. M.  
1079 Al-Hashimi, "Energy-efficient run-time mapping and thread partitioning  
1080 of concurrent applications on CPU-GPU MPSoCs," *ACM Trans-*  
1081 *actions on Embedded Computing Systems*, vol. 16, no. 5s, p. 147, 2017.
- 1082 [38] T. Li, P. Brett, R. Knauerhase, D. Koufaty, D. Reddy, and S. Hahn,  
1083 "Operating system support for overlapping-isa heterogeneous multi-  
1084 core architectures," in *Intl. Symp. on High Performance Computer*  
1085 *Architecture*. IEEE, 2010, pp. 1–12.
- 1086 [39] V. Petrucci, O. Loques, D. Mossé, R. Melhem, N. A. Gazala, and S. Gob-  
1087 riel, "Energy-efficient thread assignment optimization for heterogeneous  
1088 multicore systems," *ACM Transactions on Embedded Computing Sys-*  
1089 *tems*, vol. 14, no. 1, p. 15, 2015.
- 1090 [40] A. Schranzhofer, J.-J. Chen, and L. Thiele, "Dynamic power-aware  
1091 mapping of applications onto heterogeneous MPSoC platforms," *IEEE*  
1092 *Transactions on Industrial Informatics*, vol. 6, no. 4, pp. 692–707, 2010.
- 1093 [41] D. Stamoulis and D. Marculescu, "Can we guarantee performance  
1094 requirements under workload and process variations?" in *Intl. Symp.*  
1095 *on Low Power Electronics and Design*. ACM, 2016, pp. 308–313.
- 1096



**Karunakar R. Basireddy** received his M.Tech. degree in Microelectronics and VLSI from Indian Institute of Technology (IIT), Hyderabad, India in 2015. He is a Ph.D. student in Electronic and Electrical Engineering at the University of Southampton, UK. His current research interests include design-time and run-time optimization of performance and energy in multi-core heterogeneous systems.



**Amit Kumar Singh** (M'09) is a lecturer at University of Essex, UK. He received the B.Tech. degree in Electronics Engineering from Indian School of Mines (IIT), Dhanbad, India, in 2006, and the Ph.D. degree from the School of Computer Engineering, Nanyang Technological University (NTU), Singapore, in 2013. He was with HCL Technologies, India for a year and half until 2008. He has a post-doctoral research experience for over five years at several reputed universities. His current research interests are system level design-time and run-time optimization of 2D/3D multi-core systems for performance, energy, temperature, reliability and security. He has published over 70 papers in reputed journals/conferences, and received several best paper awards.



**Bashir M. Al-Hashimi** (M'99-SM'01-F'09) is an ARM Professor of Computer Engineering, Dean of the Faculty of Physical Sciences and Engineering, and the Co-Director of the ARM-ECS Research Centre, University of Southampton, Southampton, U.K. He has published over 380 technical papers. His current research interests include methods, algorithms, and design automation tools for low-power design and test of embedded computing systems. He has authored or co-authored five books and has graduated 35 Ph.D. students.



**Geoff Merrett** (GSM'06-M'09) is an Associate Professor in the School of Electronics and Computer Science at the University of Southampton, UK, and Head of its Centre for IoT and Pervasive systems. He received the B.Eng. and Ph.D. degrees from Southampton in 2004 and 2009, respectively. His research interests are in energy management of mobile and embedded systems, and has published over 175 articles in journals/conferences in these areas.