

CHOAMP: Cost Based Hardware Optimization for Asymmetric Multicore Processors

Jyothi Krishna V S, Shankar Balachandran, and Rupesh Nasre

Abstract—Heterogeneous Multiprocessors (HMPs) are popular due to their energy efficiency over Symmetric Multicore Processors (SMPs). Asymmetric Multicore Processors (AMPs) are a special case of HMPs where different kinds of cores share the same instruction set, but offer different power-performance trade-offs. Due to the computational-power difference between these cores, finding an optimal hardware configuration for executing a given parallel program is quite challenging. An inherent difficulty in this problem stems from the fact that the original program is written for SMPs. This challenge is exacerbated by the interplay of several configuration parameters that are allowed to be changed in AMPs. In this work, we propose a probabilistic method, CHOAMP to choose the *best* available hardware configuration for a given parallel program. Selection of a configuration is guided by a user-provided *run-time* property such as energy-delay-product (EDP), and CHOAMP aspires to optimize the property in choosing a configuration. Core part of our probabilistic method relies on identifying the behavior of various program constructs in different classes of CPU cores in the AMP, and how it influences the cost function of choice. We implement the proposed technique in a compiler which automatically transforms a code optimized for SMP to run efficiently over an AMP, eliding requirement of any user annotations. CHOAMP transforms the same source program for different hardware configurations based on different user requirement. We evaluate the efficiency of our method for three different run-time properties: execution time, energy consumption and EDP, in NAS Parallel Benchmarks for OpenMP. Our experimental evaluation shows that CHOAMP achieves an average of 65%, 28% and 57% improvement over baseline HMP scheduling while optimizing for energy, execution time, and EDP respectively.

Index Terms—Asymmetric Multicore Processors, big.LITTLE, Optimal Hardware configuration, scheduling, Compiler Optimization

1 INTRODUCTION

Advancement in processor technology with symmetric multicore processor (SMP) scaling has resulted in large increase in power consumption per unit of chip area. Increasing importance of embedded systems has further emphasized the need of power-efficient computations while designing the CPUs. The quest to reduce power consumption has led to the fabrication of asymmetric multicore processors (AMPs). big.LITTLE [1] from ARM and Tegra-E1 [2] from NVIDIA are examples of AMPs.

AMPs are characterized by multiple types of cores. However, unlike in CPU-GPU heterogeneous multicore processors, all cores in AMPs generally implement the same instruction set architecture (ISA) and the processes can easily be migrated from one core type to another (similar to SMP).¹ In AMPs, a core type differs from another in terms of frequency bandwidth, micro-architecture, and cache size, among other aspects. These design choices are made based on power and performance constraints. The cores in AMPs can be logically divided in memory-centric cores (low-power cores) and compute-centric cores (high-power cores). The memory-centric cores have lower cache miss penalties (based on CPU cycles), lower static and dynamic power ratings, but at a cost of lower computational efficiency when compared to the compute-centric cores.

big.LITTLE is an example of an AMP system. In

big.LITTLE, the powerful cores are called the big (brawny) cores and the weaker power-efficient cores are called the LITTLE (wimpy) cores. The number of big and LITTLE cores may vary across implementations. A comparison of big and LITTLE cores is given in Table 1.

Table 2 shows the idle and peak power consumption of big.LITTLE cores obtained by running several micro-benchmarks (*cf.* Section 6) on Odroid-XU3 board [6]. The power consumption readings are from the onboard power sensors. We collect the power consumption of a single core by switching off the other cores in the cluster. We observe that even for similar frequencies, the power consumption of big is much higher than that of LITTLE. big cores also perform better compared to LITTLE in case of compute-intensive programs while running at the same frequency [7].

Earlier implementations of big.LITTLE relied on cluster switching [8], which disallowed simultaneous usage of big and LITTLE cores. Global Task Scheduling / Heterogeneous

TABLE 1: Comparison of big and LITTLE cores

	LITTLE core	big core
Core Types	Cortex-A7, Cortex-A35, Cortex-A53	Cortex-A15, Cortex-A17, Cortex-A57, Cortex-A72
Pipeline	simple 8-10 stage in-order	15-24 stage out-of-order
Frequency	400 - 1400 MHz (For A7)	800 - 2000 MHz (For A15)
Speed	1.9 DMIPS [5]	3.5-4.01 DMIPS
L1 cache size	32 KB	32-48 KB
Instruction Set	Thumb-2	

- Jyothi Krishna V S is with IIT Madras, India.
E-mail: jkrishna@cse.iitm.ac.in
- Shankar Balachandran is with Intel Labs, Bangalore.
- Rupesh Nasre is with IIT Madras.

1. Heterogeneous-ISA AMPs also exist [3], [4]. But we do not deal with those in this work.

TABLE 2: Peak and idle power consumption for different frequencies of big and LITTLE in Watts.

Cortex A15 (big) Frequency	Single Core		Four Cores	
	Idle	Peak	Idle	Peak
2.0 GHz	0.95	2.40	1.15	5.28
1.8 GHz	0.70	1.65	0.70	3.50
1.6 GHz	0.51	1.40	0.53	2.80
1.4 GHz	0.38	0.85	0.40	2.30
1.2 GHz	0.30	0.70	0.32	1.80

Cortex A7 (LITTLE) Frequency	Single Core		Four Cores	
	Idle	Peak	Idle	Peak
1.4 GHz	0.20	0.50	0.22	1.40

Multicore Processor (HMP) scheduling [9] addressed this issue, allowing all the cores to be used at the same time. HMP scheduling is now integrated in the Linux kernel.

When the AMPs are designed with considerable frequency difference between the powerful brawny cores and weaker wimpy cores, the selection of optimal type of cores turns out to be a trade-off between execution time and power consumption [10]. One can choose the brawny cores to run the threads faster and the weaker cores to save on energy consumption. However, such a clear cut favorite between these two cores is absent when a composite cost functions such as Energy-Delay-Product (EDP) [11] needs to be optimized for arbitrary programs. The issue gets exacerbated due to variation in the frequencies the cores of a type are clocked at. For instance, as the frequency difference between different types of cores decreases, the weaker (simpler) cores exhibit better execution time than the powerful cores (with complex hardware optimizations) for heavily memory-bound programs [7]. In a similar manner, a powerful core may consume less memory than a weaker core depending upon their operating frequencies and the program characteristics. In a system with more than two types of cores, where the computational power difference between two adjacent classes of cores is not substantial, the choice of an optimal core type becomes obscure. The optimal solution point can easily shift to next core type with slight changes in the input program. Our work in this article deals with identifying such an optimal configuration for a program by looking at the program characteristics, to be executed on a hardware with multiple core-types with varying execution and power behavior. The notion of optimality is also dictated by the user.

The application performance in a multicore environment degrades considerably in the presence of intensive inter-thread communication. In such cases, a power-efficient solution may discard using big cores in favour of a more power-efficient core type (such as LITTLE) with minimal or no performance penalty. On one extreme, for strong CPU-bound programs, performance of big cores is twice as that of LITTLE cores [12]. On the other extreme, for memory-bound programs, it is now known that a large number of cores suffers from poor performance, with the big cores incurring higher penalty. For such programs, an optimal number of cores is present [12] for each of the core types for which we can obtain an optimal execution time.

We anticipate most real world programs to have a complex mix of characteristics (compute, memory, synchroniza-

tion). Therefore, an optimal configuration such programs would need to make a judicious use of cores of varying power (in our experimental setup, big and LITTLE cores). This motivates us to holistically model scheduling of parallel tasks, their assignments to various core-types, and configuration of each core. We illustrate that parallel systems exhibit wide variation in throughput based on the scheduling, and how our model helps achieve near-optimal schedule and configuration. While we illustrate the efficacy of our model using big.LITTLE system, we emphasize that our model is more general and applies to an AMP configuration having capability of several different types of cores.

Former research in this area has focused on scheduling in AMPs to improve execution time [13], [14], energy consumption [15], [16], [17], or a function of both [18], [19]. These works propose novel techniques to exploit the hardware heterogeneity for optimal performance and / or power consumption. However, most of the proposals are tied to optimizing one of the runtime features. The proposed systems are inflexible towards the user-requirements. One of the primary goals of our proposal is to optimize a runtime cost-function provided by the user. Thus, some users may wish to optimize execution time given an energy budget, while some others may wish to optimize peak power subject to a time budget, while yet another group of users may wish to optimize the energy-delay product (EDP).

To cater to such user-requirements, we propose a holistic system named CHOAMP which makes use of a compiler, program analysis and transformation, regression analysis, and a runtime scheduler. A close-knit interaction of these system components helps us overcome the challenges arising due to unknown variable values during compilation, and unknown execution paths in the control-flow graph, and achieve a more precise modeling of the parallel program. Together, the components lead to a scheduling with appropriate core-configurations leading to near-optimal performance (based on the user-driven cost-function).

This article makes the following contributions:

- We design a Predictor that can predict the optimal core-configuration, optimizing a user-provided runtime cost-function for running OpenMP programs. This involves an in-depth learning of all possible hardware configurations and its response to parallel programming elements. We have devised a supervised learning system where the runtime behavior is captured using different regression functions. We then use this knowledge-base to analyze the input program to predict its runtime behavior.
- We build a compiler to transform OpenMP programs to run with the predicted hardware configuration.
- We present an in-depth experimental study of thread-scheduling in asymmetric multiprocessing environment using a suite of NAS Parallel Benchmarks [20]. We observe that the predicted hardware configurations are mostly better than the base case and very often close to the best. We also show CHOAMP's ability to work well with dynamic scheduling using the HMP [9] scheduler for big.LITTLE and Compiler Enhanced Scheduling [21], a custom scheduler for OpenMP for AMPs.

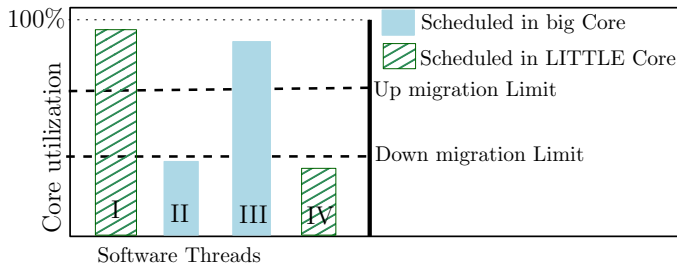


Fig. 1: Example of HMP dynamic scheduling.

The sequel is organized as follows. Section 2.3 briefly introduces HMP scheduling, CES and OpenMP. Section 3 explains our micro-benchmarking framework. Section 4 provides the theory behind CHOAMP. Section 5 provides detailed explanation of our proposed methodology. Section 6 discusses the implementation details including the compiler, regression analysis and the transformation and illustrates the effectiveness of our approach using experimental evidence. Section 7 compares and contrasts with the related work, and Section 8 concludes.

2 BACKGROUND

We provide a brief background on HMP, CES and OpenMP.

2.1 HMP

Heterogeneous multiprocessor (HMP) scheduler is a dynamic scheduler designed for big.LITTLE systems based on the core CPU-cycle utilization by a thread. Thus, threads with high CPU utilization are executed on big cores, while others in LITTLE cores. A thread that is currently scheduled in LITTLE thread can be *up-migrated* to big core if the thread has a higher CPU utilization than the up-migration limit. Similarly, a thread that is currently scheduled in big can be *down-migrated* if the current CPU utilization of the thread is lower than the down-migration limit. The scheduler is invoked periodically to keep high performing threads in big.

Figure 1 shows an example HMP scheduling system depicting percentage utilization of CPU cycles by different threads in big and LITTLE. Thread I is showing high CPU utilization (beyond a threshold) in a LITTLE core and the thread can benefit and run faster if it is scheduled in a more powerful core. Hence, the thread is scheduled for up-migration. On the contrary, thread II is currently scheduled in a big core with a low CPU utilization (below another threshold). Down-migrating to a LITTLE core can reduce energy consumption without much increase in its execution time. The HMP scheduling is integrated into the Completely Fair Scheduler [22] (CFS) in the Linux kernel.

2.2 CES

Compiler Enhanced Scheduling (CES) [21] is a *static* scheduler for OpenMP programs in AMPs. CES optimizes execution time by dividing the workload in the parallel region based on the task workload and the processing power of the core. The compiler estimates the relative speed-up of each core type for executing the threads and divides the parallel tasks accordingly such that the execution time of each thread

in a parallel region is normalized. The big cores usually receive larger chunks of work when compared to the LITTLE cores. There is a dynamic part to CES which attenuates the imbalance in the initial work distribution which might have sneaked in due to the compile-time assumptions.

2.3 OpenMP API

We provide an introduction to a subset of OpenMP parallel program features that are relevant to this paper. OpenMP API is designed for shared memory parallelism in C, C++ and FORTRAN. Our focus is on C and C++ programs. Special *#pragma* directives are used by the programmer to specify the intended OpenMP program behaviour to the compiler. A parallel program written using OpenMP API can have interleaved sequential and parallel regions. A parallel region, written inside *#pragma omp parallel* (*omp_parallel*), creates a team of threads, which may run in parallel to execute the code defined. The number of threads (*N_THREADS*) in a team can be set using environment variables or using OpenMP library calls (*set_omp_num_thread()*). For better performance, a parallel program needs to exhibit more parallelism and lesser synchronization.

The master thread is the initial thread which creates the team on encountering the *omp_parallel* region. In a team, thread id 0 is the master thread, while the remaining threads with ids 1..*N_THREADS*-1 are non-master threads. The code specified under *omp_master* (*#pragma omp master*) is executed only by the master. The remaining threads in the team skip *omp_master* and continue executing the parallel region. The parallel region is executed by master too.

Threads may have private variables, and they can communicate via shared variables. Shared variables can be synced to main memory using *omp_flush* (*#pragma omp flush*). The set of variables to be synced can be optionally specified as parameters to *omp_flush*. If the set is not specified, entire local cache of the thread is synced with main memory before proceeding. A barrier construct (*#pragma omp barrier*) is used to synchronize among all the team threads. Each thread waits at the barrier until all the other team threads reach the barrier. Barriers have implicit flush operations on all the threads. There is implicit barrier at the end of every *omp_parallel* construct.

Apart from these, OpenMP provides *omp_atomic* (*#pragma omp atomic*) and *omp_critical* (*#pragma omp critical*) for updating and accessing shared variables without incurring a data race. *omp_atomic* is used to atomically read, write or update a single shared scalar variable at a time. The *omp_atomic* is only applicable to the immediately following statement and not all binary operations are allowed. A *omp_critical* block is executed by all the team threads in a mutually exclusive manner. All unnamed *omp_criticals* are considered identical and maximum one thread can be in an *omp_critical* at a given point of time. Unlike for atomics, there is no restriction on the operations or the data types allowed inside an OpenMP critical region.

Work-sharing constructs define units of work distributed across team threads. The work sharing constructs have an implicit barrier at the end, which can be removed using *nowait* clause. There are three work-sharing constructs available in OpenMP C/C++. First, the *omp_for* (*#pragma omp*

for) implements a parallel for loop in N_ITRS iterations, where each iteration is executed once by one of the threads. The iterations can be executed in parallel with other iterations. The scheduling pattern of iterations to threads may be specified using `static`, `dynamic`, or `guided` clauses. In static scheduling, each thread is assigned equal-sized chunks of iterations in a round-robin fashion, until there are no iterations left. The size of a chunk, unless specified, is $N_ITRS / N_THREADS$. In dynamic scheduling, each thread is supplied with a chunk of iterations on demand, until there are no iterations left. Guided scheduling is very much similar to dynamic scheduling except that the chunk size (whose initial value can be specified using the optional parameter) starts off large and decreases in later steps to improve load-balancing. Guided and dynamic scheduling are often useful for better load-balancing, but have higher overheads compared to static scheduling. The chunk size plays an important role in the performance of dynamic and guided scheduling. If chunk size is very small, the scheduling overhead increases. On the other hand, if chunk size is very large, the load-imbalance is high. Therefore, a judicious setting of chunk size often affects parallel performance.

`omp_reduction` (`reduction(identifier :list)`), is often used with `omp_for` to update a shared scalar variable on an associative operation. `identifier` is used to identify the reduction operation and `list` contains one or more variables on which the reduction operation might be applied inside the region. We can have multiple `omp_reductions`, one for each operation in a single `omp_for`.

In `omp_sections` (`#pragma omp sections`), we have multiple sections each encapsulated by `#pragma omp section` (`single_omp_section`), executed by one of the threads in the team. Section scheduling is arbitrary. A single `omp_section` once scheduled to a thread, will be executed by the thread.

`omp_single` (`#pragma omp single`) contains a block of code that is executed by only one of the randomly-selected team threads. The rest of the threads jump to the end of the `omp_single` block and wait at the implicit barrier at the end of the construct, unless `nowait` is specified.

Figure 2 provides an OpenMP code that computes the cumulative sum and the maximum element in the array A . Inside the `omp_parallel` region (from Line# 6 to Line# 20) we have an `omp_for` (from Line# 8 to Line# 19) which iterates over A . Location `cumSum[i]` stores the cumulative sum of A 's elements from position 0 to i . Variable MAX stores the maximum element of A . To avoid race conditions, MAX is updated inside a critical region (from Line# 15 to Line# 18). Variable `localSum` stores the sum of A 's elements and is updated using `omp_reduction` in Line# 10.

3 MICRO-BENCHMARKS

We now describe the generation of micro-benchmarks and training of the prediction model. Our micro-benchmark framework is inspired by Eigenbench [23], a work on benchmarking different Transactional Memory (TM) systems.

EigenBench describes a tool (EigenTool) which uses a set of micro-benchmarks to analyze the efficiency of a TM system. The micro-benchmarks are designed to test a set of orthogonal application characteristics that form a basis for transactional behaviour. Precise extraction of the Eigen characteristics is a key step in designing the EigenTool.

L#	
1	#define M 50000
2	int f(int *s, int A[], int cumSum[], int L)
3	{
4	int MAX = 0, localSum = 0, temp = L/128;
5	int N = M - temp;
6	#pragma omp parallel
7	{ int i, j;
8	#pragma omp for reduction(+:localSum)
9	for(i = 0; i<N ; i++) {
10	localSum += A[i];
11	cumSum[i] = 0;
12	for(j=0; j<N; j++) {
13	if(j<=i) cumSum[i] += A[j];
14	}
15	#pragma omp critical
16	{
17	if(MAX < A[i]) MAX = A[i];
18	}
19	}
20	}
21	return MAX;
22	}

Fig. 2: Running example: An OpenMP code to find the maximum element and the cumulative sum of an array.

3.1 Feature Identification

In CHOAMP, we identify a set of parallel program characteristics in OpenMP (from hereon termed prime features), whose existence and intensity variations affect the runtime behaviour of the program in an AMP environment. The choice of prime characteristics depends on both the parallel programming platform (OpenMP) constructs and the hardware characteristics (in general, HMP; in our experiments, big and LITTLE). Similar to EigenTool, we use a collection of set of micro-benchmarks (one for each prime feature selected) to carefully study the behaviour of all the available hardware configurations and how each of the runtime cost functions of choice is affected. This helps us identify the efficiency of different hardware configurations of big.LITTLE. We use the data collected from these micro-benchmark runs to train a predictor, which is used to predict the optimal configuration O_{bl} for a given parallel program.

Each micro-benchmark is a valid OpenMP program generated with a known *intensity vector* (*ivector*) of the prime feature set. Each ivector value represents intensity of an prime feature as a fraction of the total ALU operations. By varying the value of a single feature in ivector over a range we create a set of micro-benchmarks whose run-time values shed light onto the effect of that program feature in different hardware configurations.

The selection of the set of program features should be such that it covers all the important factors affecting the runtime behavior of our test programs. A sub-optimal set of features would result in an inaccurate predictor whereas choosing a lot of *non-influential* features would increase the training time of the predictor. The prime features can be divided into language-independent and language-dependent features. The language-independent features mostly revolve around dissimilarity in the hardware features of big and LITTLE. The language-dependent features vary across par-

allel programming models and across the parallel constructs provided by them (in this case OpenMP).

The language-independent features represent the basic set of operations which are executed differently by different core types. They include memory operations, branch operations and data dependencies. The cache latencies and miss penalties of the **big** cores are much higher than those of the **LITTLE** cores. For instance, L2 latency for **big** is 21 cycles while for **LITTLE** it is 10 cycles [7]. Thus, as we increase the amount of memory operations in the program the speed-up advantage of **big** cores over **LITTLE** cores reduces. The amount and density of the memory operations also plays an important role in determining a core’s performance. The **big** cores are usually stacked with hardware optimizations such as out-of-order execution, reorder buffers, and can withstand multiple cache misses. If the interval between two consecutive memory operations is adequate enough, the effect of a stall in **big** can be mitigated to some extent. Further, **big** cores are usually equipped with larger caches. **LITTLE** cores, on other hand, have in-order execution and the pipeline stalls on every non-cached memory access. To capture this behaviour, we create multiple micro-benchmark programs with the same amount of memory operations but with different spreads of the load and store operations.

We also vary the virtual address spacing between the target address for consecutive memory operations from adjacent locations (to mimic array accesses which are cache friendly) to purely random locations in a buffer larger than the L3 cache size (to mimic pointer based accesses which shows low spacial locality) This helps us capture the effect of memory access on different types of cores.

A branch misprediction can affect the ILP, and in turn, a core’s performance. As we move from **LITTLE** to **big**, the complexity and the number of pipeline stages increase. This results in **big** cores incurring a higher penalty for every branch misprediction. Similar to the memory operations, we generate a set of micro-benchmarks with different spreads and intensities of branch operations. Our model also considers the possible effect of false-sharing. The effect of false-sharing becomes prominent when we consider using smaller chunk for scheduling the `omp_for` iterations having store operations on array locations. This increases the memory traffic, and thereby, memory latency.

The OpenMP specific features selected include `barrier_constructs`, `omp_criticals`, `omp_flushes` and `omp_for` for this work. The OpenMP constructs chosen as prime features are based on their impact on inter-thread communication and synchronization. As we increase the `N_THREADS` value, the synchronization cost also increases. So, for high intensities of `barrier_constructs`, larger value of `N_THREADS` would suffer from higher execution time. In case of `omp_criticals`, more work done inside critical sections will be advantageous for **big** cores. This is because we will be able to execute the critical regions faster and thereby reduce the critical section waiting time. On the other hand, larger value of `N_THREADS` increases the waiting time penalty of the system, leading to lower energy-efficiency. `omp_flush` is costlier for **big** than a normal memory operation as the former blocks the pipeline until the flush operation is complete. The effect of `omp_flush` on **LITTLE** is not as severe due to in-order execution. `omp_reductions` have a similar effect on

big and **LITTLE** as the reductions use local flushes at the end of the region to update values. **CHOAMP** also needs to train for different kinds of scheduling and sizes of `omp_for`. Static scheduling would result in equal work-distribution (assuming homogeneous workload in each iteration). **big** cores at high frequencies would execute their workloads faster and reach the barrier. This would lead to hardware under-utilization and thread migrations due to HMP scheduling.

3.2 Using Micro-Benchmarks

The code for micro-benchmark generators is available online [24]. We vary the intensity of the selected feature in the micro-benchmarks in a practical range (based on the analysis of various OpenMP benchmarks, namely, NPB [20], and FSU OpenMP programs [25]), and record various run-time parameters of interest. This is repeated for all the hardware configurations. The micro-benchmarks are semi-automatically generated to stress various performance aspects in the program. **CHOAMP** supports three cost functions: execution time, energy and EDP.

For each prime factor, a benchmark set is executed with the intensities of other prime factors being constant. We then combine the outputs of a benchmark set into the following equation. Since we have kept all other intensities constant and assume that the dependencies between prime factors are negligible we can generate the following equation

$$X * c_i = y_i^T \quad (1)$$

Here, each row in the matrix X corresponds to an intensity of the prime factor represented in different powers starting from zero. The c_i corresponds to the co-efficient vector for hardware configuration i . The y_i vector corresponds to the cost function values recorded using the micro-benchmark set for hardware configuration i . If we consider a linear relationship between the prime factor and the cost function, the matrix X will contain simply two columns ($c_0 * x^0 + c_1 * x = y$).

Using a regression tool we can estimate the co-efficient vector and which are used by our Predictor to find y of an unknown intensity of the prime factor in the input program for each possible hardware configuration.

4 CHOAMP THEORY

CHOAMP employs a clustering-based approach for scheduling iterations to cores. Thus, iterations that are clustered together are scheduled together. The challenging aspect in applying clustering to parallel programs is that a single feature is often not sufficient to identify the cluster an iteration belongs to. Overall performance of an application is a complex interplay of several features such as ALU operations, memory instructions, and synchronization constructs. However, since an iteration needs to be executed by a single core in our setup, it is imperative to identify the most prominent feature which affects a given cost function. The challenge gets exacerbated when multiple features are prominent. Such scenarios demand prioritization over features to identify the final assignment of an iteration to a core. We first discuss our clustering approach in the context of various program features in isolation, and then explain how

CHOAMP computes the final core-assignment combining individual recommendations from each feature.

4.1 Branch Instructions

Arbitrary branch instructions reduce cache effectiveness, and affect execution time. Executing branch instructions generally consume more energy than usual ALU operations as shown in Table 7. Thus, all our cost functions (execution time, energy and EDP) are directly affected by branches. Scheduling branch-heavy iterations to big cores can lead to their under-utilization and idling, and would also consume more energy. Therefore, whenever the percentage of branch instructions is higher than a threshold (determined by the predictor based on microbenchmark results), it is beneficial to cluster all such iterations to LITTLE cores.

4.2 Shared Data

Iterations accessing a lot of shared data should be scheduled together. Such a clustering leads to better cache efficiency and reduces inter-thread communication latency. A key concern here is to choose the type of core to schedule the iterations together. In our experience, it needs to be guided by the frequency of memory operations (and other features such as synchronization). For instance, for low memory traffic, it is better to schedule iterations on big, and so on.

4.3 Memory Traffic Diversity

An interesting feature that affects scheduling is the diversity of memory traffic. We observed that scheduling the iterations with diverse memory demands together improves performance. Also, contrary to the usual belief, scheduling of iterations with high memory demand onto LITTLE leads to inferior performance (when the cost function is execution time). Instead, the performance is better if such iterations are scheduled on big but with lower frequency.

4.4 Synchronization Instructions

Synchronization instructions play a crucial role in the execution time as well as energy consumed by a parallel program. CHOAMP models the behavior of three synchronization constructs: atomics (omp atomic), reductions (omp reduction) and critical (omp critical). The clustering behavior differs across each of them. For atomics, if the percentage (within an iteration) is beyond a threshold, the iterations are clustered onto LITTLE, otherwise big. When a medium type of core is present in an HMP system, the iterations can be clustered on it based on another threshold. The thresholds are suggested by the microbenchmarking (Section 3). On the other hand, reductions involve thread coordination. Hence, CHOAMP clusters the participating iterations onto the same type of core. For execution time, big is preferred. Critical section involves yet another consideration. If iterations involve many critical sections (beyond a threshold), CHOAMP prefers to schedule the corresponding iterations sequentially. Otherwise, depending upon the cost function, the iterations get clustered on different types of cores. In particular, for execution time, CHOAMP prefers big. For energy as the cost function, CHOAMP prefers LITTLE if the critical

section involves shared data. For EDP, CHOAMP uses a combination of the above two rules. In case of barriers, the biggest consideration would be the synchronization cost and the barrier waiting time. When the number of barriers increases and the cost function is execution time, CHOAMP prefers smaller number of cores. For energy as the cost function, CHOAMP prefers LITTLE for executing the barrier instructions. In the case of EDP, CHOAMP prefers a symmetric core system, where all cores have similar processing power, thus reducing the overall barrier waiting time.

4.5 Putting Features Together

Individual feature vectors lead to independent suggestions on the placement of iterations. However, depending upon the program, sometimes, the suggestions may lead to conflicting placements. For instance, if the number of branch instructions is high, the individual suggestion from the rules in Section 4.1 would be to schedule the iteration on LITTLE. However, if in addition, if the memory traffic is high, based on the suggestion from the rules in Section 4.3, the placement would be on big. Therefore, it is imperative to combine individual suggestions by appropriately prioritizing them. This leads to devising a weight-function to reach the final placement prediction. CHOAMP supports three types of weight functions: linear, quadratic and Gaussian, as shown in Table 3, which we discuss next in more detail.

5 CHOAMP IMPLEMENTATION

We now discuss using micro-benchmarks for training the predictor, finding the optimal O_{bl} configuration and the transformation of the program. The workflow of our technique follows the block diagram shown in Figure 3, and consists of primarily two phases. (i) Training Phase, where we generate micro-benchmarks to train the predictor for the selected program features and populate the knowledge base. (ii) Compilation Phase, where we analyze the input program and output the O_{bl} configuration for the input program. We explain the two phases in more detail below.

5.1 Training Phase

During the training phase we run sets of micro-benchmarks in various available hardware configurations and collect the power, energy and runtime values. We use the Java Scientific Library [26] provided by M. Flanagan to estimate the coefficient values represented in equation 1. Once the coefficients are estimated, they are input into our Predictor.

The Predictor combines the different prime planes (listed in Section 3) of the source program using weighted euclidean distance to generate the combined cost of all the prime features of the program.

We can assume different regression techniques which can be fixed into Equation 1. For example, if we consider linear regression, the equation boils down to $y = c_1 * x + c_0$ for each prime factor, x . In CHOAMP, we have used three different regression functions for the prediction model: (i) Linear, (ii) Quadratic and (iii) Gaussian. All three models are trained using the same micro-benchmark set. From our experiments, we observe that different regression functions have different levels of accuracy based on how well they

TABLE 3: Cost calculation for various Regression functions.

Regression	Cost Estimation
Linear	$\sum_{e=0}^N w_e c_e [1] * wl(e) + (\sum_{e=0}^N w_e * c_e [0]) / N$
Quadratic	$\sum_{e=0}^N w_e * (c_e [1] * wl(e) + c_e [2] * wl(e)^2 + c_e [0]) / N$
Gaussian	$\sum_{e=0}^N w_e Y_e / (\sqrt{2 * \pi * sd_e * (e^{(wl(e)-X_e)/sd_e})})$

N: Number of prime Features
c[*]: coefficients from regression model using micro-benchmarks
 w_e : weight of prime feature
 $wl(e)$: intensity of the prime feature in the program
 Y_e : Y scale prime feature in the program
 X_e : Mean of prime feature in the program

are able to fit each feature and the runtime intensity of the feature. CHOAMP with predictor trained with Linear (Quadratic, Gaussian) regression is termed linear-CHOAMP (quadratic-CHOAMP, gaussian-CHOAMP respectively).

The cost function calculation for different regression training functions is given in Table 3. The cost function is computed as a weighted sum of values for each prime vector. The wl vector contains the intensities of the prime factors. N denotes the number of trained prime features. The regression functions estimate the relationship between the prime factors and the selected cost function. w_e provides differential weight for different prime features. Linear regression (c_l) uses a linear combination of ivector. The quadratic regression (c_q) is more powerful, as the cost function is modelled as a quadratic function of ivector. Though the function is more powerful, c_q is more susceptible to errors such as overfitting. The gaussian regression (c_g) assumes all the prime factors to be in a continuous domain. Using weighted sum, we can vary the importance of an individual prime feature by varying the weights. The choice of the predictor is guided by a global flag.

Table 4 shows the different prime features trained for CHOAMP. We show the study of memory micro-benchmark set in Table 5 for a subset of the hardware configurations trained. For instance, configuration 0L4b(2) contains zero LITTLE and four big, with big’s frequency set to 2.0 GHz. The table presents the average values of the runtime cost function obtained, for the highest (1.0) and the lowest ivector (0.1) values for memory operations. The ivector values for the remaining prime vectors are kept constant. For memory operations, as we can see for execution time, big is advantageous over LITTLE; for energy, LITTLE is preferred; and for EDP, LITTLE is advantageous beyond a threshold.

Table 6 shows the average values of the runtime cost function obtained, for the highest (3.5E-5) and the lowest (4.77E-10) ivector values for the barrier micro-benchmark set. We can see that for higher intensities of barriers, using more cores is expensive on all fronts.

Table 7 shows the average values of the runtime cost function obtained, for the highest (0.5) and the lowest (5E-4) ivector values for the branch micro-benchmark set. We see that for very high percentage of branch instructions, LITTLE outperforms big in terms of execution time.

5.2 Compilation Phase

The first step in the compilation phase is the analysis of the input program to extract the intensities of prime features. A normal OpenMP program can have multiple parallel regions interleaved with sequential processing. In

TABLE 4: Prime features chosen

Prime Feature	Description	ivector value
Branch operations	Percentage of Branch operations	varied from 0 to 30%.
Memory operations	Trained for multiple densities of LOAD and STORE operations for same number of memory operations.	varied from 50 - 100%.
Atomic operations	percentage of omp_atomic operations	varied from 0 to 15%.
Barriers	Number of barriers in a parallel region. Trained for multiple densities of barriers for same number of barriers.	from 0 to 3.05E-03%
Critical sections	Percentage of operation inside critical section when compared to entire workload of a parallel section	from 0 to 6.10E-03%
False sharing	Taken as a percentage of all STORE operations in the parallel region.	Trained for different strides with potential false sharing.
Flush operations	Percentage of memory locations flushed using omp_flush.	varied from 0 to 5%.
omp_for	Percentage of entire parallel region inside a omp_for.	varied from 80% to 100%.
Reduction operations	Percentage of reductions in a parallel region as a factor of all ALU operations	varied for 0 to 25%.

TABLE 5: Runtime values for the highest intensity (1.0) to the lowest intensity (0.1) in memory set micro-benchmark for a subset of configurations. The remaining prime values in the ivector are kept constant.

Config.	Exec. time in s		Energy in J		EDP in Js	
	high	low	high	low	high	low
4L0b	5.70	2.01	5.82	1.82	33.31	3.67
0L4b(2)	2.71	0.62	15.71	2.88	42.57	1.81
4L2b(2)	2.33	0.67	13.43	1.46	38.56	2.55
4L4b(2)	1.80	0.47	11.39	2.43	20.93	1.15

TABLE 6: Runtime values for the highest intensity (3.05E-05) to the lowest intensity (4.77E-10) in barrier set micro-benchmark for a subset of configurations. The remaining prime values in the ivector are kept constant.

Config.	Exec. time in s		Energy in J		EDP in Js	
	high	low	high	low	high	low
4L0b	5.45	5.27	5.26	4.86	28.68	25.65
0L4b(2)	1.91	1.78	8.55	8.09	15.35	21.09
4L2b(2)	3.87	2.26	17.49	10.20	67.75	23.02
4L4b(2)	43.66	1.80	231.42	11.99	1.0E4	22.59

TABLE 7: Runtime values for the highest intensity (0.5) to the lowest intensity (5E-4) in branches set micro-benchmark for a subset of configurations. The remaining prime values in the ivector are kept constant.

Config.	Exec. time in s		Energy in J		EDP in Js	
	high	low	high	low	high	low
4L0b	0.835	0.98	1.23	1.05	1.03	1.03
0L4b(2)	1.13	0.63	4.83	4.62	5.24	3.05
4L2b(2)	0.68	0.52	4.65	4.35	3.12	2.8
4L4b(2)	0.61	0.42	5.88	4.45	2.49	3.30

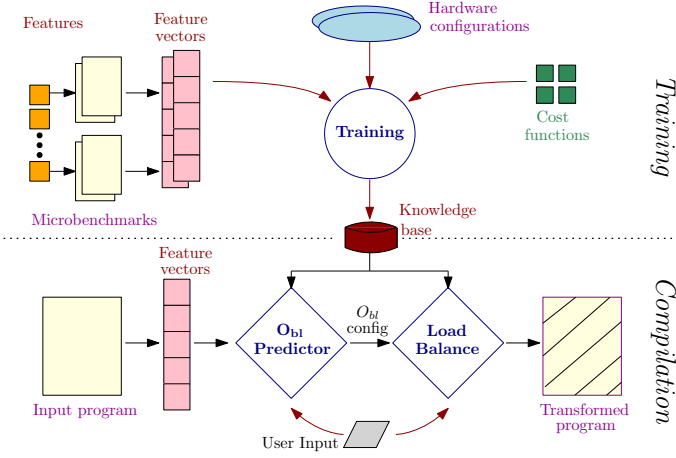


Fig. 3: Block diagram depicting the workflow of our framework. The top half shows the training phase and the bottom half shows various compiler stages in code-scheduling.

this work we focus on optimizing only the parallel regions, which allows us to analyze each parallel region separately. CHOAMP creates a concurrent control-flow graph (CCFG) for the input program. All the basic-blocks in the CCFG are associated with a workload (wl) vector. The wl vector is assigned the intensity of each program feature (prime) of interest in the block. Each basic block appends its wl to the wl of its parent block for computing the combined effect.

CHOAMP traverses the CCFG to compute the intensity of each of the prime features in the input source program. For precise and effective modeling of OpenMP programs, we need to have an inter-procedural workload accumulation. CHOAMP creates a call-graph to store all the function calls with hooks to the call-sites. The call-graph analysis is a two-step process. The first step is intra-procedural wherein individual functions are analyzed in isolation. The second step is inter-procedural which takes into account the calling context of a function. To handle the context, CHOAMP topologically sorts the call-graph and processes functions in that order. Use of recursion generates potentially infinite number of contexts as recursion-depth is unknown at compile-time. CHOAMP currently does not handle recursion.

The intensity of each prime characteristic is stored as the fraction of the total number of prime characteristic operations performed in the code. After feature extraction is over, CHOAMP passes the generated wl vector to the trained prediction model, and calculates a cost function value for each selected hardware configuration for the runtime cost function of interest. The model then predicts a hardware configuration having the least cost function value to run the parallel region.

In presence of multiple parallel regions in the program, CHOAMP computes the O_{bl} configuration for each `omp_parallel` region separately. However, it is difficult in `big.LITTLE` to apply different hardware configurations to different parallel regions of the same program. There are multiple ways in which a single configuration can be chosen. CHOAMP prioritizes the configuration associated with the most intensive parallel region (dictated by wl vectors), as that is likely to dominate the cost function.

TABLE 8: Value ranges with prominent binary operations

Operation	Result range
$i + j$	$\text{range}(\text{MAX}(R(i)) + \text{MAX}(R(j)))$
$i * j$	$\text{range}(\text{MAX}(R(i)) * \text{MAX}(R(j)))$
$i - j$	$\text{range}(\text{MAX}(R(i)) - \text{MIN}(R(j)))$
i / j	$\text{range}(\text{MAX}(R(i)) / \text{MIN}(R(j)))$
$i >> j$	$i - \text{MAX}(R(j))$
$i \% j$	$R(j)$
$i = j$	$\text{MAX}(R(i), R(j))$
$\text{MAX}(R(i))$: Known maximum value in range i	
$\text{range}(k)$: Range to which the value k belongs	

Consider our OpenMP example from Figure 2. The example contains a single parallel region with an embedded `omp_for` with N iterations. The wl vector values are collected as a function of N for some of the prime factors (like memory operations and `omp_critical`) and are independent of N for some other primefactors (barriers). As a result, the accuracy of the Predictor would depend on the estimation of the value of N . If the runtime value of N is unknown, CHOAMP equates it to the largest known constant in the program.

5.3 Estimating Values of Unknown Variables

Precision of the estimated wl vector is critically affected by the unknown variables. As with most static analyses, in general, the values of several variables are unknown at compile time. Especially in the context of our analysis in CHOAMP, knowing range of variable values can considerably improve prediction accuracy, leading to improved prediction. There are primarily two kinds of unknown variables that plague effectiveness of CHOAMP: loop bounds, and variables used in conditional statements. We handle these as follows. We first preprocess the source code to collect the program constants and the variables of interest; that is, loop bounds and variables used in conditional statements. To keep the method effective and practical, we choose to bucketize their values in exponential ranges. Thus, we define a set of variable ranges (R) bounded by powers of a selected base (such as 10) $R(i)$ representing the variable range to which the variable i belongs. A variable is expected to be pushed into the range which would bound the largest value the variable can attain during execution. The preprocessing happens as below. First, the program constants are pushed into a corresponding R with a flag to identify the constant. All the statements updating the variables of interest (loop bounds and those used in conditionals) are processed, and the unknown variables are pushed into the variable ranges based on their update values and the expressions they are used in. CHOAMP currently models the binary expressions (which are more frequent) as shown in Table 8. In the end, the remaining (unclassified) variables of interest are pushed into the largest non-empty variable range. After such a bucketing, for a non-constant variable i , it assumes the largest known value in the $R(i)$ as the value of i .

For example, consider Figure 2. We have N , the variable of interest which depends on variable L , which is unknown. The initial unknown value L is put into the current highest non empty bucket ($R_{1E4-1E5}$ containing constant M with value 50,000). The range of `temp` is estimated to be $R_{1E2-1E3}$ ($\text{range}(\text{MAX}(R(L)))$) i.e 50,000/128). Once the range of `temp` is estimated we estimate the range of N

TABLE 9: The runtime values of O_{bl} configuration predicted for different cost functions for the example given in Figure 2.

	Base case 4L4b(2)	Manual Best case	CHOAMP predictions		
			Linear	Quadratic	Gaussian
Time in s	6.34	5.18 4L4b(1.6)	5.36 4L4b(1.8)	6.34 4L4b(2)	5.36 4L4b(1.8)
Energy in J	30.16	12.52 4L0b	12.52 4L0b	12.52 4L0b	12.52 4L0b
EDP in Js	191.13	116.91 4L4b(1.4)	116.91 4L4b(1.4)	121.85 4L4b(1.2)	166.33 4L2b(1.2)

($range(M - MIN(R_{1E2-1e3}))$ i.e 50,000-1000) as $R_{1E4-1E5}$. Hence the value of N is estimated as 50,000 which is the maximum known value in the range of N. Once the variable values are resolved, CHOAMP proceeds to find the cost function value for each of the hardware configuration trained, using the functions described in Table 3.

5.4 CHOAMP Example

Now we revisit our example in Figure 2 to demonstrate the efficiency of our Predictor. Figure 4 shows the scaled runtime characteristics with the scaled cost function values calculated by each of the regression functions. We use $iL_jb(k)$ notation to represent a hardware configuration, where i stands for the number of LITTLE cores, j stands for the number of big cores and k stands for the frequency of big in GHz. The LITTLE core frequency is kept constant at 1.4 GHz. All values are scaled with respect to values of 4L4b(2) configuration. The CF values predicted by CHOAMP mostly follows the same pattern as the scaled runtime values helping us to effectively identify the best configuration.

Table 9 shows the runtime values of the example for the O_{bl} configuration predicted by each Predictor along with the O_{bl} configuration predicted by each Predictor and the best hardware configuration. We can see that the Predictors are able to predict hardware configurations that are close to the optimum configuration.

5.5 CHOAMP Scope

CHOAMP continues to be under active development. It currently supports almost all of the OpenMP features including various parallel sections, scheduling strategies, and synchronization constructs. However, there are a few non trivial aspects which CHOAMP does not support: nested parallelism and recursion and OpenMP tasks.

OpenMP allows nested parallelism where each member in the team can create a new set of threads, where the creator thread assumes the role of the master in the new team. The newly created team will have its own local barriers and pragma constructs. Even though nested parallelism theoretically increases the effective number of threads, it is going to be limited by the hardware configuration.

As discussed earlier, recursion poses issues with respect to potentially infinite number of calling contexts. In our study with NPB, none of the nine benchmarks uses nested parallelism or recursion. Therefore, we leave recursion handling as a future work.

We restrict the discussion of CHOAMP to the subset of synchronization constructs that are prominent in the benchmarks we have chosen for initial learning [25], [27]

and NPB [20] benchmark we use for testing our proposed method. Constructs like OpenMP tasks will introduce more challenges. With tasks, all the work in a parallel region is not available at once. Also, the tasks can switch threads (untied tasks). All these are dependent on runtime which makes it difficult to estimate and model the ivector for prediction.

While considering a different AMP hardware, we need to update the language-independent prime features. For example, consider a distributed AMP, with NUMA architecture, the shared-data (including false-sharing) induced stalls can be costlier than say, branch mis-prediction stalls. Also barrier synchronization can be differentially costly for same base hardware configuration but with different core chosen based on the selected cores' locations. We can address these issues by differentiating each core with an id and thereby the selected hardware configuration based on the participating cores. This will exponentially increase the number of possible hardware configurations to be trained for and increase the span of the training phase.

We use regression model as our learning algorithm as a simple model that can capture the effects of the prime factors, with low modelling and training costs. We are able to capture the essence of the program features with the regression model, though more complex and costly ML kernels like SVMs and Neural Networks can replace the regression model with possibly greater efficiency while learning for more complex prime features.

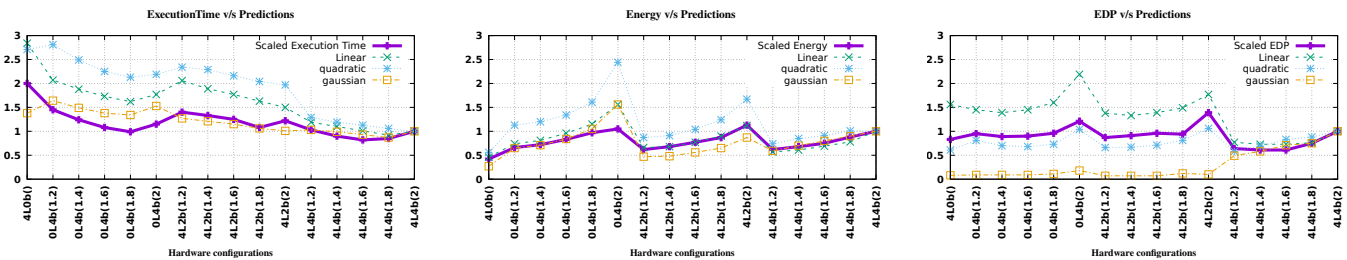
Even though the model has been explained using OpenMP, the underlying CHOAMP methodology of (i) identifying the prime factors, (ii) learning a Predictor over the prime factors based on the parallel programming model of choices and AMP and (iii) finally using the trained Predictor for predicting the best possible configuration available with the AMP to a good effect. For example, for considering pthreads, as the possible sets of language dependent prime features could include pthread_create, pthread_join and pthread_exit. These features pose similar challenges as OpenMP tasks and would require more exhaustive sets of micro-benchmarks to train a good Predictor.

5.6 Handling Imbalanced Workload

The inherent power imbalance in Asymmetric Multicore Processors makes it difficult to balance running times of threads in parallel programming platforms such as OpenMP. Current compilers are tuned for optimizing the code for SMP hardware with equivalent amount of workload given to each thread between two barriers.

Since the big cores can usually execute threads much faster than the LITTLE cores, often the threads that are initially scheduled in big cores will reach the barrier earlier than the threads that are scheduled in LITTLE core. Once a thread running on a big core reaches the barrier, its CPU utilization reduces and the HMP scheduler would eventually down-migrate it to a LITTLE core. On the other hand, a thread running on a LITTLE core that is yet to reach the barrier may be up-migrated. Such a scenario can lead to overall hardware under-utilization. The associated performance penalty is high if there is an imbalance in the number of different kinds of cores.

In omp_for, static scheduling is the most preferred because of its low overhead. In AMPs, however, static



(a) Scaled execution time compared to different regression functions. (b) Scaled Energy compared to different regression functions. (c) Scaled EDP compared to different regression functions.

Fig. 4: Runtime characteristics of the example in Figure 2 for different hardware configurations.

scheduling often results in load imbalance and frequent thread-migrations, which can considerably hamper performance [28]. CHOAMP employs multiple load-balancing mechanisms to reduce this inefficiency, as we discuss below.

- Using dynamic or guided scheduling for a `omp_for`: dynamic and guided scheduling mechanisms offer relatively better overall load-balancing. A naïve usage of these mechanisms incurs additional scheduling overhead, but it can be mitigated by choosing an appropriate chunk-size with the help of the predictor, to reduce the number of scheduling points.
- Using custom scheduling: Here, CHOAMP first fixes a core-to-thread mapping and then divides the iterations based on the speed-ratios of various types (such as big-to-medium and medium-to-LITTLE). Note that these ratios are not fixed, but vary based on the core frequencies chosen by the predictor. For our experiments, we use CES [21] for the custom scheduling.

6 EXPERIMENTAL EVALUATION

6.1 Setup

We use IMOP [29], a source-to-source transformation and analysis framework for OpenMP. The benchmark programs are pre-compiled using gcc 4.8. Precompilation evaluates the preprocessor directives to increase accuracy and reduce compilation time. For regression training, we use Java Scientific Library provided by Flanagan [26]. The evaluation system configuration details are presented in Table 10. For our experimentation purposes we have trained for a subset of the numerous hardware variations achievable with big.LITTLE: (i) vary the number of big cores from 0 to 4 (step size of 2), (ii) vary the number of LITTLE cores from 0 to 4 (step size of 4), and (iii) vary the frequency of big cores from 1.2 GHz to 2 GHz (step size of 0.2 GHz). For all the experiments we fix the frequency of LITTLE cores to 1.4 GHz, as found empirically to provide benefits.

We use NPB OpenMP benchmarks [20] for evaluating CHOAMP. Characteristics of NPB benchmarks relevant to this work are presented in Table 11. The number of barriers listed corresponds to the explicit barriers in the code. The number of lines of code includes the code only from the main program, and not from libraries.

Our experiments are carried out on a big.LITTLE system where the HMP scheduler is enabled. We use the default hardware configuration as the baseline, which is 4L4b(2) (4 LITTLE cores running at 1.4 GHz and 4 big at 2 GHz).

TABLE 10: Test environment description

Hardware	Odroid XU3 [6]
Processor	Samsung Exynos 5422
Microbenchmark scripts	Python version 2.7
Parallel platform	OpenMP [30]
Frontend compiler	IMOP [29]
Backend compiler	gcc 4.8.2
Regression tool	Java Scientific Library [26]
OS	LUbuntu 16.04
big frequencies tested	2GHz, 1.8GHz, 1.6GHz, 1.4GHz, 1.2GHz
Configurations tested	4 big & 4 LITTLE, 4 big & 0 LITTLE, 2 big & 4 LITTLE, 0 big & 4 LITTLE, 2 big & 2 LITTLE

TABLE 11: Characteristics of NPB Benchmarks

Name	Lines	parallel	omp_for	barriers	omp_critical
BT	2653	2	70	38	0
CG	550	5	24	23	0
DC	2478	1	0	0	1
EP	169	1	1	1	1
FT	775	2	7	11	1
IS	285	2	2	2	1
LU	2658	3	36	27	1
MG	828	6	11	22	1
SP	2206	2	78	65	0

6.2 Predicted Configurations and their Performance

Table 12 shows the O_{bl} configurations predicted by CHOAMP when trained using the linear regression. For execution time as the cost function, all big cores are chosen in all but one benchmark (LU), which follows our expectation.

CG and SP have high percentage of barriers; this prompts CHOAMP to choose a lower number of cores (no LITTLE

TABLE 12: O_{bl} selected for NPB benchmarks computed by linear regression. All LITTLE cores run at 1.4 GHz. The value in brackets shows the frequency of big in GHz.

Benchmark	Configuration chosen for a cost function		
	Execution Time	Energy	EDP
BT	4L4b(1.8)	4L0b	4L4b(1.4)
CG	0L4b(1.8)	4L0b	0L4b(1.4)
DC	4L4b(1.8)	4L0b	4L4b(1.4)
EP	4L4b(1.8)	0L4b(1.2)	4L4b(1.6)
FT	4L4b(1.8)	4L0b	4L4b(1.6)
IS	4L4b(1.8)	4L0b	4L4b(1.4)
LU	4L2b(2.0)	4L0b	0L4b(1.4)
MG	4L4b(1.8)	4L0b	4L4b(1.6)
SP	0L4b(1.8)	4L0b	0L4b(1.4)

cores). While trying to optimize for energy consumption, CHOAMP chooses a core configuration which is closer to the lower end of the spectrum. Thus, when the cost function is energy, CHOAMP ends up choosing mostly the LITTLE cores, except for EP. EP is a highly parallel benchmark with very little synchronization cost and high ILP, allowing the big cores to execute instructions much faster than the LITTLE cores, prompting CHOAMP to choose big over LITTLE for energy consumption. While training for a more complex cost function like EDP, a more balanced configuration is chosen for most of the benchmarks. The Predictor always tend to choose a hardware configuration which lie in-between the one chosen for Execution time and Energy in the spectrum. Most of the configurations chosen have the maximum number of cores possible to exploit maximum parallelism. For CG and SP, CHOAMP selects less number of cores because of the high synchronization present in the benchmarks. The big cores frequency chosen are similar to that of LITTLE. This helps in balancing the workload (most `omp_fors` in NPB have static scheduling).

Table 13 compares the execution time, energy consumption and EDP values obtained for baseline with the corresponding values for the program optimized for the cost function by CHOAMP (with configurations mentioned in Table 12). For CG and SP, we get high improvement in execution time as CHOAMP chose a lower number of cores for executing the programs with high synchronization costs. SP, in fact, shows even better runtime for fewer threads (2 and 1) outside the trained spectrum. Since the baseline has high power consumption, we see huge improvement in energy consumption of CHOAMP-optimized programs. A similar trend is observed when the cost function is EDP. This clearly indicates the efficacy of CHOAMP.

A closer look at some benchmarks: Now, we look into a few of these benchmarks in detail. Figures 5-7 show scaled execution time, energy consumption and EDP values obtained for EP, CG and FT for different hardware configurations we trained for. The hardware configurations are inversely sorted first in terms of the total number of cores and then in terms of big core frequency. The first entry in the graphs represents the value corresponding to the base case (4L4b(2)) to which the rest of the values are scaled to. In all the cases, CHOAMP predicts the configuration with the lowest cost function value as per the model. In the plots we have marked the configuration selected by CHOAMP.

First, we look into EP (Figure 5), a highly parallel benchmark with very little synchronization inside a parallel region. Due to its inherently parallel nature, the execution time reduces with increasing number of cores and core frequency. Linear-CHOAMP is able to accurately predict the best configuration for running EP while optimizing for execution time. High parallelism in EP helps higher hardware configurations to get better EDP values. Both linear-CHOAMP and gaussian-CHOAMP are able to predict hardware configuration close to the optimal.

As shown in Table 11, CG has a large number of barriers which increase the synchronization costs considerably (Figure 6). This restricts the ability of the program to scale well with the number of processors. As a result, CG performs better in terms of execution time and EDP for a lower num-

ber of cores. Linear-CHOAMP is able to predict hardware configuration that is very close to the optimal value.

In terms of program characteristics, FT (Figure 7) lies in between EP and CG. The benchmark is fairly scalable. The best configuration for execution time lies at the high end of the spectrum with best configuration for EDP in the middle of the spectrum. Linear-CHOAMP and gaussian-CHOAMP predict a similar trend with all O_{bl} predictions showing substantial improvement in the cost-function values.

6.3 Cost Functions

Figure 8 shows the overall percentage improvement (reduction) in execution time for CHOAMP-compiled programs over base case compilation. We see an average improvement of 28% in execution time with linear-CHOAMP over the nine benchmarks in NPB. DC contains very few program-specific prime features we have trained. SP and CG (for linear-CHOAMP) show high improvement in execution time for lower number of cores (4), due to its high number of barrier instructions which induces high synchronization costs for more cores. Gaussian-CHOAMP predicted all LITTLE cores for LU (4L0b), resulting in a higher execution time. The improvement while using quadratic-CHOAMP (2.7%) and gaussian-CHOAMP (3%) is much lower.

Figure 9 shows the overall percentage improvement (reduction) in energy consumption for CHOAMP-compiled programs over the program running in base case. We see an average improvement of 65% in energy consumption while using linear-CHOAMP over the nine benchmarks in NPB. Both quadratic-CHOAMP (average of 50%) and gaussian-CHOAMP (67%) are also able to provide very good improvements in energy consumption values. The high level parallelism of EP prompted quadratic-CHOAMP to choose all big cores (0L4b(2)) for energy consumption. But high power consumption of big cores has resulted in negative improvement in energy consumption.

Figure 10 shows the overall percentage improvement (reduction) in EDP values for CHOAMP-compiled programs over the program running in base case. We see an average improvement of 58% in EDP for linear-CHOAMP over the nine benchmarks in NPB. For SP, the EDP value obtained for CHOAMP-compiled program is negligible when compared to the base case and hence reported as 100% improvement (corrected to two decimal places). Quadratic-CHOAMP predicted the base configuration for MG for optimizing EDP resulting in zero improvement. EP was given all big (0L4b(2)) by quadratic-CHOAMP. The high power consumption of big cores gives us a higher EDP than base. However, linear-CHOAMP provides considerable benefits.

6.4 CHOAMP and CES

In this subsection, we analyze the compatibility of CHOAMP with a dynamic scheduler (CES [21]). Figure 11 shows the average percentage improvement in execution time of (CHOAMP + CES)-compiled code over CES-compiled code. If CHOAMP outputs the program for a homogeneous hardware configuration (for example CG and SP while using linear-CHOAMP), the CES compiler is left ineffective and will not change the code. We see an average of 14% improvement in execution time while using Linear-CHOAMP when

TABLE 13: Runtime cost function values of base case (4L4b(2)) compared with O_{bl} runtime cost function

Benchmark	Execution time in s		Energy in J		EDP in J.s	
	baseline	O_{bl}	baseline	O_{bl}	baseline	O_{bl}
BT	295.10	258.64	1446.92	486.95	426987.66	240872.73
CG	7.29	3.46	34.74	10.67	253.29	44.54
DC	113.59	122.42	363.78	124.38	41320.86	32228.11
EP	18.52	15.62	89.61	53.00	1659.82	1159.98
FT	12.60	8.33	61.61	26.23	776.19	214.03
IS	1.83	1.73	9.49	4.12	17.36	8.65
LU	235.95	169.89	1178.02	280.41	277954.89	32355.37
MG	4.71	4.08	23.63	10.68	111.35	73.06
SP	74.36	0.18	358.73	0.22	26676.79	0.06
Mean	84.88	64.93	396.28	110.74	86195.36	34106.28

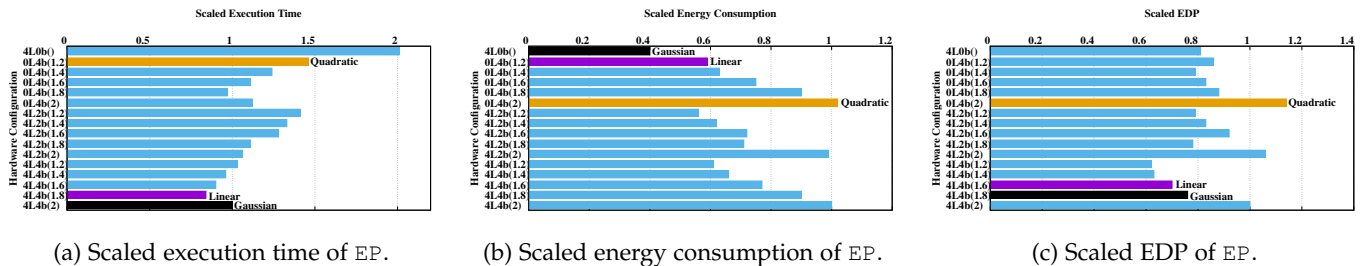


Fig. 5: Runtime characteristics of EP for different hardware configurations. Smaller the better.

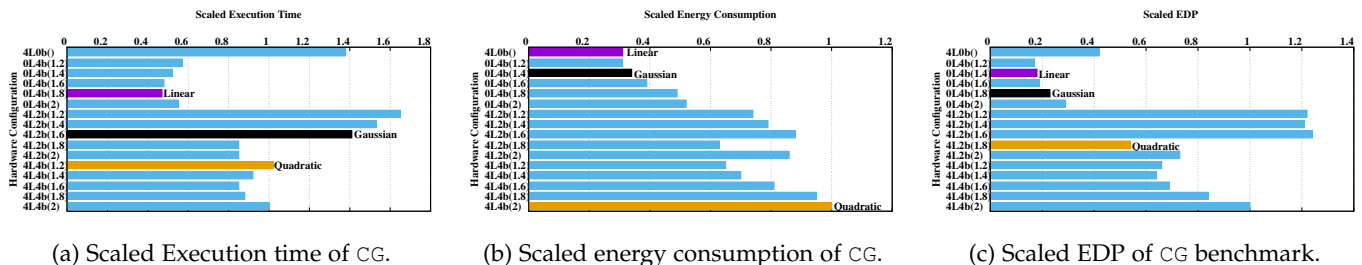


Fig. 6: Runtime characteristics of CG for different hardware configurations. Smaller the better.

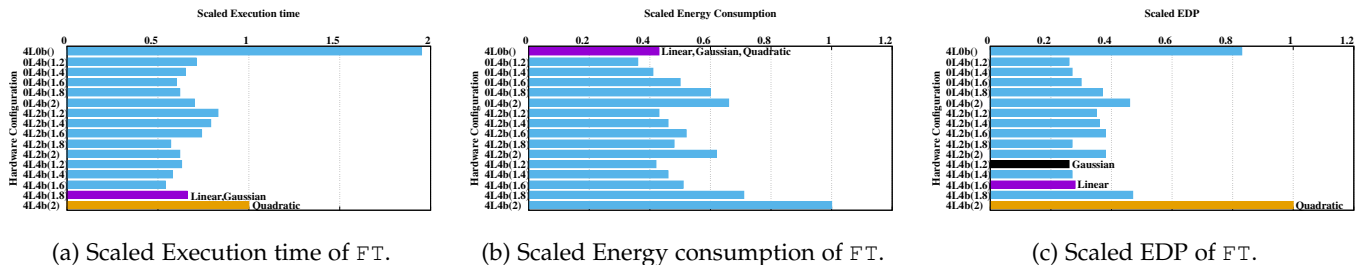


Fig. 7: Runtime characteristics of FT for different hardware configurations. Smaller the better.

compared to CES compiled code. CES removes the implicit barriers in `omp_fors` and balances the thread workload with its core power. As a result the barrier waiting time is reduced and the threads are able to make optimal use of the processing power of both big and LITTLE cores. This was one of the reasons of high improvement in execution time for SP and CG with CHOAMP. We see a huge reduction in CHOAMP's improvement for SP and CG. For EP with quadratic-CHOAMP (0L4b(2)) and for LU with gaussian-CHOAMP (4L0b(2)), there are considerable increases in the execution times.

Overall, we illustrate that CHOAMP offers considerable benefits over the baseline, and allows a user to seamlessly

optimize parallel programs for AMP.

7 RELATED WORK

Our scheme completely rely on compile time analysis and is tested on real hardware and does not require any extra hardware support. A new hardware configuration can be easily integrated into the existing model. We just have to train the Predictor learn for the new hardware configurations. Here we discuss some of the relevant related works on selecting optimal hardware configuration and scheduling in heterogeneous multicore processors.

Majority of the work on scheduling in AMP relies on dynamic analysis as precise information of runtime pa-

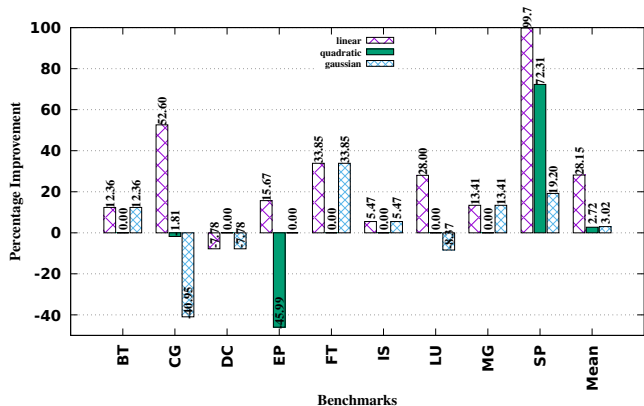


Fig. 8: Percentage improvement in execution time

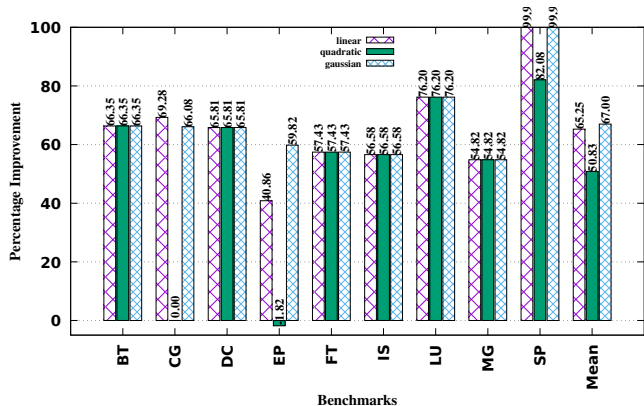


Fig. 9: Percentage improvement in Energy consumption

rameters is available [13] [31] [14] [32] [33]. There have been a few works that rely on compile-time profilers for making scheduling decisions [34], [35]. Compared to the dynamic schedulers (which operate at execution time), static schedulers scale better with the number of cores as well as program complexity [34]. With dynamic scheduling, the schedulers are often required to be fast and light-weighted.

Some of the dynamic scheduling techniques proposed for normalizing the execution time of threads [33], HASS [34], [35] can easily be integrated to CHOAMP in place of HMP and CES. Chen and John [35] proposed a

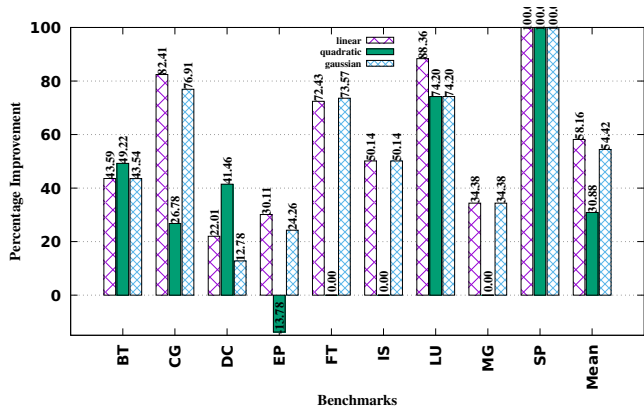


Fig. 10: Percentage improvement in EDP

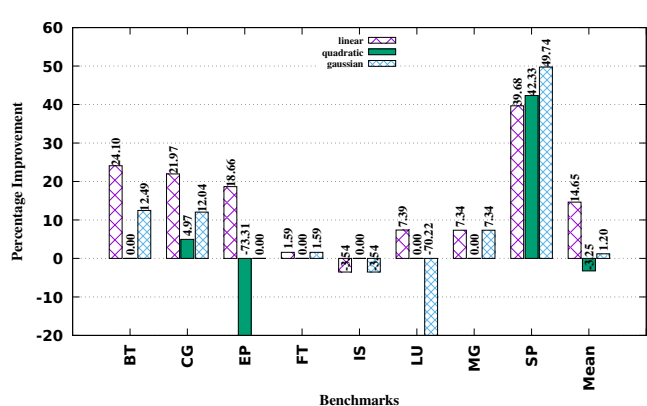


Fig. 11: Percentage execution time gain while executing various NPB Benchmarks (CES + CHOAMP). The base case is CES compiled program for 4L4b(2.0) configuration.

dynamic scheduling technique to map the inherent program characteristics and their corresponding resource demands to the characteristics of different core types in AMPs. The proposed method combines the program characteristics and available core configurations into a unified solution space using weighted Euclidean distance. We also use Euclidean distance to combine the effect of individual prime factors to estimate the overall cost of the program.

A few works have looked into dynamic frequency scaling along with thread scheduling [18] [36] for optimal performance measures. Annamalai et al. [18] proposed a scheduler which dynamically chose between swapping two threads and changing the frequency to optimize for throughput/Watt.

Multiple works have relied on regression for estimating execution time and power [10] with good amount of accuracy. A few of the dynamic algorithms designed for AMPs which uses similar basic concepts are discussed below. Bias scheduling [14] presents a dynamic scheduling technique where a thread is given a big or LITTLE core bias based on the speedup ratio. PIE (Performance Impact Estimation) [13] shows taking the ratio of ILP (Instruction Level Parallelism) and MLP (Memory Level Parallelism) gives a very good estimate of the performance of big core and LITTLE core.

Taylor et al. [37] has proposed to use SVMs to map OpenCL kernels to heterogeneous multicores. The proposed method is implemented on LLVM and trained for multiple CPU-GPU configurations. The CPU configurations are all homogeneous in nature which might not provide the best hardware configuration as demonstrated by Figures 5-7. Memti and Pillana [32] have proposed to use simulated annealing to reduce the solution space exploration to find the optimal system configuration in a heterogeneous system.

8 CONCLUSIONS

In this work we propose CHOAMP probabilistic model, with a micro-benchmark trained Predictor which predicts the best asymmetric hardware configuration for a given cost function. This eliminates the manual effort involved in trial-and-error method for tuning for optimal configuration. CHOAMP can be flexible to different user needs as the

same program can be transformed for optimizing execution time, energy consumption, and EDP. Experimental results with NAS Parallel Benchmarks show 28% improvement in execution time, 65% improvement in energy consumption and 58% improvement with LTO over the baseline.

In future we plan to add support for more cost functions like execution time given an energy or power limit, energy consumption given a limit on execution time, as a factor of the base case values. Also, we would like to introduce more hardware configurations permitted by AMPs like variable cache sizes and memory bandwidths. We would also like to consider the inter-dependencies of the prime factors which would require using the same training set to for all prime features by using more sophisticated learning models.

REFERENCES

- [1] J. Brian, "Big.LITTLE system architecture from ARM: saving power through heterogeneous multiprocessing and task context migration." in *DAC*, P. Groeneveld, D. Sciuto, and S. Hassoun, Eds., 2012, pp. 1143–1146.
- [2] NVIDIA, "NVIDIA Tegra K1 A New Era in Mobile Computing," NVIDIA, Tech. Rep., 2014.
- [3] A. Venkat and D. M. Tullsen, "Harnessing isa diversity: Design of a heterogeneous-isa chip multiprocessor," in *41st ISCA*. IEEE Press, 2014, pp. 121–132.
- [4] M. DeVuyt, A. Venkat, and D. M. Tullsen, "Execution migration in a heterogeneous-isa chip multiprocessor," *SIGPLAN Not.*, vol. 47, no. 4, pp. 261–272, Mar. 2012.
- [5] R P Weicker, "Dhrystone: A Synthetic Systems Programming Benchmark," *Commun. ACM*, vol. 27, no. 10, pp. 1013–1030, Oct. 1984.
- [6] "Odroid-XU3," http://www.hardkernel.com/main/products/prdt_info.php?g_code=G140448267127.
- [7] J. M. Kim, S. K. Seo, and S. W. Chung, "Looking into heterogeneity: when simple is faster," in *2nd IWPMMP*, 2014.
- [8] "Big.LITTLE Switchers, Evaluation on Exynos BL Processor," 2012. [Online]. Available: http://events.linuxfoundation.org/images/stories/pdf/kf2012_yu.pdf
- [9] H. C. H. Chung, M. Kang, "Heterogeneous Multi-Processing Solution of Exynos 5 Octa with ARM big.LITTLE Technology," 2013. [Online]. Available: http://www.arm.com/files/pdf/Heterogeneous_Multi_Processing_Solution_of_Exynos_5_Octa_with_ARM_bigLITTLE_Technology.pdf
- [10] U. Hölzle, "Brawny cores still beat wimpy cores, most of the time," *IEEE Micro*, vol. 30, no. 4, pp. 23–24, 2010.
- [11] M. Pricopi, T. S. Muthukaruppan, V. Venkataramani, T. Mitra, and S. Vishin, "Power-performance modeling on asymmetric multi-cores," in *CASES '13*. IEEE Press, 2013, pp. 15:1–15:10.
- [12] X. Liang, M. Nguyen, and H. Che, "Wimpy or brawny cores: A throughput perspective," *JPDC*, vol. 73, no. 10, pp. 1351–1361, 2013.
- [13] K. Van Craeynest, A. Jaleel, L. Eeckhout, P. Narvaez, and J. Emer, "Scheduling Heterogeneous Multi-cores Through Performance Impact Estimation (PIE)," in *39th ISCA*, ser. ISCA '12, 2012, pp. 213–224.
- [14] D. Koufaty, D. Reddy, and S. Hahn, "Bias Scheduling in Heterogeneous Multi-core Architectures," in *5th EuroSys '10*. ACM, 2010, pp. 125–138.
- [15] R. Kumar, D. M. Tullsen, P. Ranganathan, N. P. Jouppi, and K. I. Farkas, "Single-isa heterogeneous multi-core architectures for multithreaded workload performance," *SIGARCH Comput. Archit. News*, vol. 32, no. 2, pp. 64–, Mar. 2004.
- [16] R. Nishtala, D. Mossé, and V. Petrucci, "Energy-aware thread collocation in heterogeneous multicore processors," in *EMSOFT '13*. IEEE Press, 2013, pp. 21:1–21:9.
- [17] L. Sawalha and R. D. Barnes, "Energy-efficient phase-aware scheduling for heterogeneous multicore processors," in *2012 IEEE Green Technologies Conference*, April 2012, pp. 1–6.
- [18] A. Annamalai, R. Rodrigues, I. Koren, and S. Kundu, "An opportunistic prediction-based thread scheduling to maximize throughput/watt in amps," in *22nd PACT*, Sept 2013, pp. 63–72.
- [19] R. Rodrigues, A. Annamalai, I. Koren, S. Kundu, and O. Khan, "Performance per watt benefits of dynamic core morphing in asymmetric multicores," in *20th PACT*, Oct 2011, pp. 121–130.
- [20] D. H. Bailey, E. Barszcz, J. T. Barton, D. S. Browning, R. L. Carter, L. Dagum, R. A. Fatoohi, P. O. Frederickson, T. A. Lasinski, R. S. Schreiber, H. D. Simon, V. Venkatakrishnan, and S. K. Weeratunga, "The NAS Parallel Benchmarks-Summary and Preliminary Results," in *SC*. New York, NY, USA: ACM, 1991, pp. 158–165.
- [21] V. S. Jyothi Krishna and S. Balachandran, "Compiler Enhanced Scheduling for OpenMP for Heterogeneous Multiprocessors," in *HiPEAC EEHCO '16*, 2016, pp. 692:1–692:6.
- [22] "CFS Scheduler," <https://www.kernel.org/doc/Documentation/scheduler/sched-design-CFS.txt>, linux, Tech. Rep.
- [23] S. Hong, T. Oguntebi, J. Casper, N. Bronson, C. Kozyrakis, and K. Olukotun, "EigenBench: A simple exploration tool for orthogonal TM characteristics," *IISWC '10*, 2010.
- [24] V. S. Jyothi Krishna, "OpenMP EigenBench," <https://bitbucket.org/jkrishnavs/openmp-eigenbench>, 2016.
- [25] J. Burkardt, "C Examples of Parallel Programming with OpenMP," [urlhttps://people.sc.fsu.edu/~jburkardt/c_src/openmp/openmp.html](https://people.sc.fsu.edu/~jburkardt/c_src/openmp/openmp.html), 2010. [Online]. Available: https://people.sc.fsu.edu/~jburkardt/c_src/openmp/openmp.html
- [26] M. T. Flanagan, "Michael Thomas Flanagan's Java Scientific Library," 2008. [Online]. Available: <http://www.ee.ucl.ac.uk/~mflanaga/java>
- [27] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S. Lee, and K. Skadron, "Rodinia: A benchmark suite for heterogeneous computing," in *IISWC '09*. IEEE Computer Society, 2009, pp. 44–54.
- [28] A. Gutierrez, R. G. Dreslinski, and T. Mudge, "Evaluating private vs. shared last-level caches for energy efficiency in asymmetric multi-cores," in *SAMOS '14*. IEEE, 2014.
- [29] A. Nougrihiya and K. Nandivada, "IIT Madras OpenMP(IMOP) Framework," <http://www.cse.iitm.ac.in/~amannoug/imop>.
- [30] L. Dagum and R. Menon, "OpenMP: An industry standard API for shared-memory programming," *Computational Science & Engineering*, IEEE, vol. 5, no. 1, pp. 46–55, 1998.
- [31] O. Khan and S. Kundu, "A self-adaptive scheduler for asymmetric multi-cores," in *GLSVLSI '10*. ACM, 2010, pp. 397–400.
- [32] S. Memeti and S. Pillana, "Combinatorial optimization of work distribution on heterogeneous systems," *CoRR*, vol. abs/1606.05134, 2016.
- [33] K. Van Craeynest, S. Akram, W. Heirman, A. Jaleel, and L. Eeckhout, "Fairness-aware scheduling on single-isa heterogeneous multi-cores," in *22nd PACT*. IEEE Press, 2013, pp. 177–188.
- [34] D. Shelepov, J. C. Saez Alcaide, S. Jeffery, A. Fedorova, N. Perez, Z. F. Huang, S. Blagodurov, and V. Kumar, "Hass: A scheduler for heterogeneous multicore systems," *SIGOPS Oper. Syst. Rev.*, vol. 43, no. 2, pp. 66–75, Apr. 2009.
- [35] J. Chen and L. K. John, "Efficient program scheduling for heterogeneous multi-core processors," in *DAC '09*. ACM, 2009, pp. 927–930.
- [36] T. S. Muthukaruppan, M. Pricopi, V. Venkataramani, T. Mitra, and S. Vishin, "Hierarchical power management for asymmetric multi-core in dark silicon era," in *DAC '13*. ACM, 2013, pp. 174:1–174:9.
- [37] B. Taylor, V. S. Marco, and Z. Wang, "Adaptive optimization for opencl programs on embedded heterogeneous systems," in *LCTES 2017*. ACM, 2017, pp. 11–20.

Jyothi Krishna V. S. Jyothi Krishna is currently pursuing PhD in the CSE department at IIT Madras. His research interests are Compilers and Heterogeneous Multicore processors.

Shankar Balachandran. Shankar works at Intel Labs, Bangalore. Formerly, he was an Associate Professor in the CSE department at IIT Madras. He received his PhD from the University of Texas at Dallas. His research interests are in Computer Architecture, VLSI Design Automation, and Parallel Algorithms.

Rupesh Nasre. Rupesh is an Assistant Professor in the CSE department at IIT Madras. He completed PhD from IISc Bangalore and Post-Doctoral Fellowship from the University of Texas at Austin. His research focus is in Compilers and Parallelization.