

This is an electronic reprint of the original article. This reprint may differ from the original in pagination and typographic detail.

Quantifying the Interaction Between Structural Properties of Software and Hardware in the ARM Big.LITTLE Architecture

Stepanovic, Srboľjub; Georgakarakos, Georgios; Holmbacka, Simon; Lilius, Johan

Published in:

2018 26th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP)

DOI:

[10.1109/PDP2018.2018.00027](https://doi.org/10.1109/PDP2018.2018.00027)

Published: 01/01/2018

Document Version

Accepted author manuscript

Document License

Publisher rights policy

[Link to publication](#)

Please cite the original version:

Stepanovic, S., Georgakarakos, G., Holmbacka, S., & Lilius, J. (2018). Quantifying the Interaction Between Structural Properties of Software and Hardware in the ARM Big.LITTLE Architecture. In I. Merelli, P. Lio, & I. Kotenko (Eds.), *2018 26th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP)* (pp. 138–144). IEEE. <https://doi.org/10.1109/PDP2018.2018.00027>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Quantifying the Interaction Between Structural Properties of Software and Hardware in the ARM big.LITTLE Architecture

Srboljub Stepanovic, Georgios Georgakarakos, Simon Holmbacka and Johan Lilius

Faculty of Science and Engineering

Åbo Akademi University

Turku, Finland

Email: firstname.lastname@abo.fi

Abstract—Heterogeneous architectures offer the opportunity to achieve high performance and energy efficiency by selecting appropriate cores for execution of ever changing software applications. Appropriate core selection depends on the interaction between the structural properties of the software and the hardware that influences performance of the software. We propose a model for efficient core selection when executing software on ARM's big.LITTLE heterogeneous architecture. It features a metric based on the correlation between the performance and the number of last level data cache (LLC) misses on a big and a LITTLE core. Additionally our model defines a soft threshold in terms of the number of LLC misses that determines efficient core selection. We verify the model on both a stress benchmark (stress-ng) and a performance and energy demanding application (HEVC decoding) using XMEM and Linux perf dynamic tools.

I. INTRODUCTION

ARM's big.LITTLE is an asymmetric multi-core architecture where cores share the same instruction set but feature significantly different microarchitectures. The Exynos 5422 chip available in the odroid-xu4 board [1] is an instantiation of the big.LITTLE architecture with two sets of cores organized in clusters: four Cortex A7 cores with simple microarchitecture and short pipelines optimized for energy efficiency and four Cortex A15 with complex microarchitecture and deep pipelines optimized for high performance. However, mapping an application on big.LITTLE is not trivial. In order to efficiently use these resources, we need to understand how to optimally map an application on the suitable core to minimize energy consumption while maintaining performance.

In this paper we examine the interaction between structural properties of hardware and software, and evaluate its impact on the achieved performance on both A7 and A15 cores on selected benchmarks. In particular we study how the microarchitecture can exploit the structural properties of the software.

We measure the interaction between structural properties of software and hardware in terms of the number of last level cache (LLC) misses that appear during execution of a program. LLC misses dominantly affect the execution

time because memory access time is an order of magnitude larger than latencies caused by other miss events like instruction cache misses, L1 data cache misses and branch mispredictions. LLC latency will cause stalls in execution of instructions that are dependent on the LLC miss instruction. This in turn prevents a big core to benefit from instruction level parallelism using out-of-order execution of independent instructions. This insight can be exploited in order to optimally utilize heterogeneous big.LITTLE architectures by scheduling applications with a large number of LLC misses on a LITTLE core and the ones with a small number of LLC misses on a big core.

The main contribution of this paper is a model to predict which processor core will achieve better cycles per instruction (CPI) when executing particular application. The model establishes a linear correlation between CPI and the number of LLC misses of an application. Each benchmark from the stress-ng benchmark suite [2] is evaluated for the number of LLC misses and CPIs on both cores. The model is then used to predict CPI for an unknown program after profiling it to get the number of LLC misses. Evaluated CPIs on both cores determine on which core a certain benchmark will be mapped. Furthermore, we establish a soft threshold in terms of LLC misses that divides all benchmarks in the two regions. The first one contains benchmarks that have the number of LLC misses below the threshold, so they will be scheduled on a big core. The second one contains benchmarks that have the number of LLC misses above the threshold, so they will be scheduled on a LITTLE core.

The rest of this paper is organized as follows: Section II presents related work done in performance estimation and latency characterization. Section III analyses the features of A7 and A15 cores of big.LITTLE. Section IV describes impact of miss events on performance. Section V describes the details of our performed measurements regarding cache and memory latencies, number of cache misses and overall performance. Also in section V the results of our measurements are presented and discussed. Section VI concludes the paper.

II. RELATED WORK

The work presented in [3] shows that there is a correlation between a small and a big core slowdown and memory level parallelism (MLP) ratio as well instruction level parallelism (ILP) ratio between the cores but it does not quantify this correlation. This work also does not specify what the threshold is in terms of memory intensity that would indicate when either MLP- or ILP-ratio should be applied to predict CPI. The work compares the following scheduling policies in terms of performance on a two-core heterogeneous multi-core: their proposed performance impact estimation including both MLP- and ILP-ratio, MLP-ratio, memory-dominance, random and optimal scheduling. We extend this work by creating a unified linear model for predicting CPI based on the number of LLC misses for all benchmarks. We also propose a threshold in terms of LLC misses that determines on which core an application will be scheduled.

The work presented in [4] introduces X-MEM, a cross-platform and extensible memory characterization tool for the cloud. This tool can be used to characterize the memory hierarchy of a specific platform. The tool is also able to measure cache and main memory unloading latency, main memory load latency and read/write behavior. We enhanced this tool to statistically determine each cache and memory component separately, based on the number of cache and main memory hits.

The design of new performance counters for measuring CPI stacks based on interval analysis is presented in [5]. The work explains the impact of long latency instructions and all relevant miss events such as L1 and LLC cache misses, and branch mispredictions on execution time of a program that is divided into intervals between these miss events. We show that LLC miss latency prevails in relation to latencies of other miss events. We take advantage of this fact to construct a model for CPI prediction based only on the number of LLC misses.

The work presented in [6] proposes a symbiotic core execution mechanism for detecting regions of code with high ILP or MLP and shows how to exploit such features on a heterogeneous architecture. The authors propose coarse-grained scheduling to determine appropriate cores outside of the regions with high MLP and also fine-grained scheduling to determine appropriate cores inside such regions. Contrary to their mechanism that is applied on chunks of instructions, our approach is applied on application level. It enables to determine CPI for an unknown benchmark after only profiling benchmarks from representative testing set such as the stress-ng benchmark suite.

Characterization of branch misprediction penalty is explained in [7] through interval analysis. In this interval analysis, superscalar processor performance is viewed as a sequence of inter-miss intervals. The misses that define

the intervals are branch mispredictions, (L1 and L2) i-cache misses and long (L2) d-cache misses. They show how the penalty for a particular branch misprediction depends on a preceding miss event and hence how this event affects overall performance. In our work we focus on L2 (LLC) d-cache misses as they most dominantly affect performance.

III. ARM BIG.LITTLE ARCHITECTURE

The platform chosen for this study is the ARM big.LITTLE multi-core system. We use the Exynos 5422 implementation with four Cortex-A7 [8] cores and four Cortex-A15 [9] cores. The LITTLE A7 is a 32 bit core using a partial dual-issue in-order microarchitecture with an 8-10 stage pipeline as shown in the upper part of Fig. 1. This means that it enables dual-issue instruction execution only for some integer operations; for all other operations A7 is a single-issue machine. This core does not have reorder buffer. The argument for using this kind of microarchitecture is to improve energy efficiency because of its simple construction with few speculative execution units such as branch prediction, Translation Lookaside Buffer (TLB) lookups and no instruction scheduling. The memory system consists of separate L1 instruction (i-cache) and data (d-cache) caches and shared L2 cache as shown in Table 1.

The big A15 cores consist of a more complex 32 bit architecture with a 15-25 stage pipeline with two separate integer and Neon units as shown in the lower part of Fig. 1. The core is also using separate load/store units and a separate unit for branch instructions. The A15 core executes instructions out-of-order to gain high throughput that enables more concurrency in the hardware. Three instructions can be processed per clock cycle. Out-of-order execution requires register renaming to dynamically alter register names in order to resolve false dependences and enable parallel hardware to operate on data without dependences. A reorder buffer is used to commit the instructions in correct order

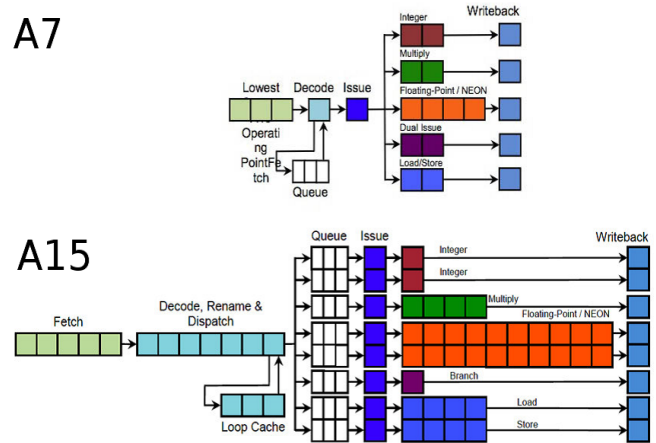


Figure 1. Microarchitecture of the A7 and the A15 [10].

after out-of-order execution. It also makes possible that the correct process state can be restored in case of a branch misprediction [7]. Reservation stations store instructions which await operands to be broadcasted without first written back to memory. The memory system consists of separate L1 i-cache and d-cache and shared L2 cache as shown in Table 1. The L1 i-cache and d-cache have Least Recently Used (LRU) cache replacement policy.

IV. IMPACT OF MISS EVENTS ON PERFORMANCE

In this section we will show how a particular miss event affects performance and what the parameters are that characterize the impact of each miss event on performance. Our approach is based on an observation from [5] that pipeline miss events can be split into frontend and backend pipeline miss events.

Indeed, miss events depend on the interaction between structural properties of software and hardware. Structural properties of software related to cache miss events are the frequency and the order of accessing particular memory locations by a program. Structural properties of hardware related to cache miss events are cache size, associativity, fixed line length and replacement policy. L1 cache may contain a subset of the data in L2 (LLC) cache. L2 cache contains a subset of the data in main memory. This means that more than one of the memory blocks can be mapped to the same cache line. Therefore, every access to a location of the memory blocks mapped to the same cache line as the memory block currently stored into this cache line will cause a cache miss event. Also, the cache is too small and cannot hold all of the memory blocks referenced in the program. A cache replacement policy defines how old data in the cache is replaced with new data. In this way two different cache replacement policies applied to the same program will cause different number of cache miss events. The parameters that characterize each miss event are penalty, rate and overlapping effect. Penalty is the time needed for particular miss event to be resolved. Rate is the number of particular miss events relative to the total number of instructions. An overlapping effect appears when two miss events happen at the same time.

Table I
SMALL/BIG CORE CACHE STRUCTURE

Parameter	L1 i-cache	L1 d-cache	L2 cache
Size[KB]	32	32	512/2048
Associativity[way]	2	4/2	8/16
Fixed line length[B]	32/64	64	64
Cache replacement policy	pseudo random/LRU	pseudo random/LRU	pseudo

The frontend pipeline miss events are L1 instruction cache misses and branch mispredictions. Penalty of these miss events consists of two components on a big core. The first one is caused by discharging instructions that follow a miss event from the reorder buffer. The second one is caused by refilling the pipeline with new instructions. On a LITTLE core refilling the pipeline is the only penalty. Because the frontend pipeline length of a big core is larger than the one of a LITTLE core and because of the size of a big core reorder buffer, the frontend pipeline miss event penalty is much larger on a big core. Instructions are executed only serially in the frontend pipeline. Therefore, the frontend pipeline miss events do not overlap at all.

The backend pipeline miss events are L1 and LLC data cache misses as well as long latency instructions. L1 data cache misses and long latency instructions are hidden through out-of-order execution on a big core. When they occur, some other instructions present in the reorder buffer will be executed until data accessed from LLC data cache become available or the long latency instruction is finished. The compulsory condition that must be fulfilled in order to avoid stalls due to these miss events is that the LLC access time and the long latency instruction execution time are smaller than the execution time of instructions present in the reorder buffer. We show that this condition is always fulfilled in section V. This shows clear advantage of a big core when considering these events.

LLC misses are events that cannot be hidden through out-of-order execution because the memory access time is always larger than the time needed to execute all instructions present in the reorder buffer. This means that all LLC misses that are handled in parallel can be considered as one miss. But because instructions may be executed in parallel in the backend pipeline also, LLC misses might overlap if they are present in the reorder buffer at the same time as was observed in [3]. When a LLC cache miss occurs instruction level parallelism cannot be exploited by out-of-order execution of a big core and therefore it is meaningful to schedule benchmarks with high number of LLC misses on a LITTLE core. The main claim of our work is summarized in Table II.

Table II
THE MAIN CLAIM OF OUR WORK

The number of LLC misses of a program	The core type a program to be scheduled on
Low number of LLC misses (below the threshold)	LITTLE core (A7)
High number of LLC misses (above the threshold)	big core (A15)

V. EXPERIMENT AND ANALYSIS

We use the X-Mem tool [4] to measure latency of the ARM big.LITTLE memory system. X-Mem measures the average aggregate L1 cache, L2 cache and main memory latency (L). We measure this latency using one worker thread across working sets from 4KB to 1GB. A working set is a private region of memory the worker thread operates on. X-Mem is doing this in such a way that for working sets up to L1 cache size it measures L1 cache hit latency (L1L), up to L2 cache size it measures average aggregate L1 and L2 hit latency and for working sets over L2 size it measures average aggregate L1, L2 and main memory hit latency. We modify X-Mem to measure the number of L1 cache (L1H), L2 cache (L2H) and main memory (MEMH) hits to be able to statistically determine L2 cache (L2L) and main memory (MEML) latencies. We use the equation shown below to calculate L2L and MEML:

$$L = L1H \times L1L + L2H \times L2L + MEMH \times MEML \quad (1)$$

At the end, we average out latency across all working set sizes belonging to corresponding L1 cache, L2 cache or main memory region respectively. The corresponding latencies for LITTLE core are shown in Fig. 2. We measure them at the highest frequency for a LITTLE core of 1.4GHz. The same set of cache and memory latencies is applied to different applications on a particular core type. The X-Mem tool introduces only negligible overhead. The tool does not specify logical/physical CPU core affinity. It is also not possible to turn off the first LITTLE core on the odroid-XU4 board and measure latencies of a big core separately. But because the same technology process is applied in implementing memory system on both cores, we assume that a big core has the same latencies.

We measure the following latencies for L1 data cache, LLC data cache and memory (Fig. 2): 2.49ns (4.98 cycles),

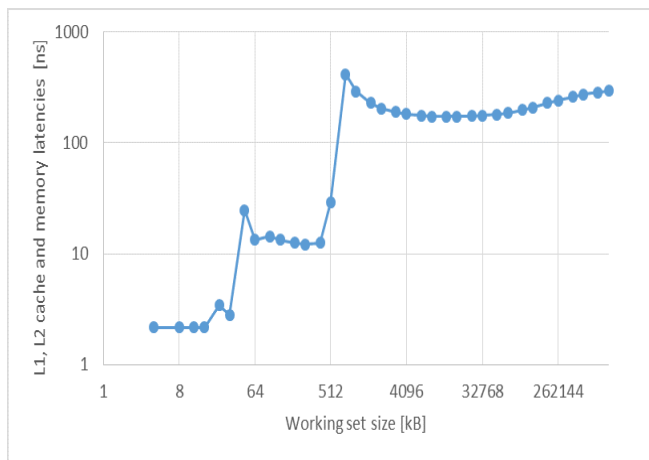


Figure 2. Determining cache and memory latencies.

16.46ns (32.92 cycles) and 223.29ns (446.58 cycles) respectively where the number in brackets is the number of cycles relative to 2GHz. When we compare LLC latency to the size of reorder buffer (128) multiplied by only one cycle needed for an instruction execution, we can notice that L1 data cache misses cannot cause stalls as stated in section IV. When we compare memory latency with size of reorder buffer multiplied by an average latency of 3 cycles, there will be stalls even if the reorder buffer is almost full and all present instructions in the reorder buffer are dependent from each other and from the LLC miss instruction.

A. Stress-ng benchmark suite

We conduct miss events and performance measurements for all benchmarks from stress-ng package on odroid-XU4 board with big.LITTLE processor architecture described in section III. We run these sequential benchmarks either on an A15 core or an A7 core in order to verify the correlation between the number of LLC misses and CPI. We use an A15 core with performance governor on the highest frequency of 2GHz and an A7 core with powersave governor on the lowest frequency of 200MHz. We use perf tool that enables measuring the number of instructions, the number of clock cycles and the number of different miss events.

We use benchmarks from stress-ng benchmark suite which contain different instruction types such as ALU, memory, branch, integer multiply and SIMD instructions to stress corresponding execution units separately and some of them together. Because one LITTLE core is always turned on, we use taskset option of the benchmark suite to specify only one big core. All benchmarks are executed at 10 bogo operations to provide the same number of instructions executed on both cores for a particular benchmark.

We show that the number of LLC cache misses per 10K instructions correlates well with CPI on both cores, by evaluating the accuracy of the CPI predictor model. Because the model is a linear function, model accuracy is a good indicator for correlation: if the model is accurate, CPI correlates well with the number of LLC misses per 10K instructions, and vice versa.

The model that transforms numbers of LLC misses per 10K instructions into CPI is shown in Fig. 3. We build this model using the numbers of LLC misses and CPI values from the benchmarks. We find that a linear relationship between the numbers of LLC misses and CPI fits best. We use the following equation to calculate the CPI from the number of LLC misses:

$$CPI(MPI) = a \times MPI + b \quad (2)$$

where MPI is the number of LLC misses per 10K instructions. Parameters a and b are determined using a least-squares fit. Every point represents the CPI (y-axis) and the corresponding number of LLC misses per 10K instructions (x-axis). There are 58 benchmarks per core. The fitted

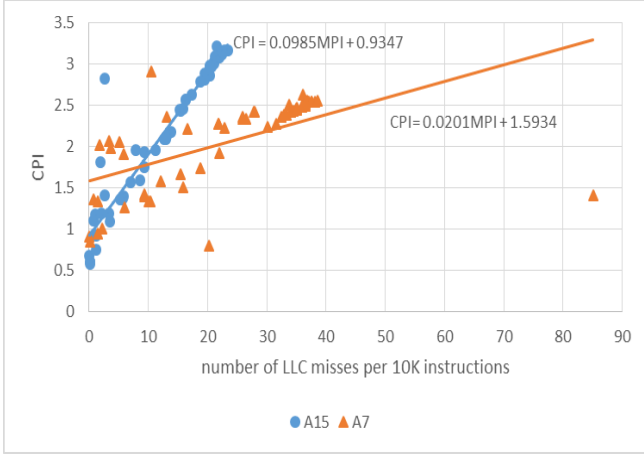


Figure 3. The CPI model for A15 and A7 cores.

model is indicated by the lines. It is clear that a linear model is appropriate for this data, and that fit is relatively good. The most important we get from this model is the crossover point in the intersection of the two lines. This point indicates appropriate core selection for each benchmark. From the equations for both lines we get 8.40 misses per 10K instructions as crossover point. Thus, we can notice that benchmarks having 8.40 or less LLC misses per 10K instructions should be scheduled on a big core, otherwise on a LITTLE one.

Fig. 4 shows the relative error for the proposed model. The average absolute relative error between real measurements and the model for A15 is 9.71% while for A7 is 21.67%. We can notice that the error continuously decreases as the number of LLC misses increases for almost all benchmarks because the memory component dominantly affects CPI in relation to other miss events for the large number of LLC misses per 10K instructions.

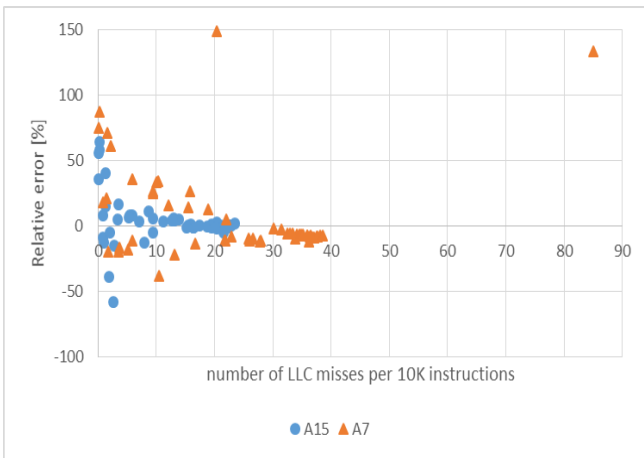


Figure 4. The relative error of CPI model.

We also run the stress-ng benchmarks on 600MHz and 1GHz either on A15 or A7 core with userspace governor. Parameter a in (2) increases while parameter b is approximately constant as frequency increases on both cores. This means that the highest observed frequencies 1GHz on A7 core and 2GHz on A15 core cause most stalls in execution i.e. the highest CPI for the same benchmark. These frequencies also cause the highest core temperature. Stalls happen in order to decrease the core temperature because userspace governor keeps fixed operating frequency i.e. there is no automatic frequency scaling to scale frequency down and decrease the temperature.

B. HEVC decoding application

In order to evaluate performance behavior of big.LITTLE cores we also use High Efficiency Video Coding (HEVC) [11] as our target test application. The reason to use HEVC comes from the growing adaptation of big.LITTLE in SoCs targeting mobile devices, where video based media applications such as content streaming is increasingly popular [12]. HEVC algorithm also manifests significant complexity as well as diversity of operations ranging from memory intensive operations (motion estimation), compute intensive operations (prediction, rate-distortion optimization, transformation) and hybrid combinations. We benchmark the A7 and A15 cores of the big.LITTLE architecture using an open source implementation of the HEVC decoder. For the decoding measurements, we create input streams with two different structures using the x265 implementation of the HEVC encoder [13] in order to examine performance differentiation with respect to different workloads. We use streams with two different Group of Picture (GOP) structures: one with frames encoded using only Intra prediction (I frames) and one with frames encoded using motion estimation as well (I, P and B frames). The difference is that P and B frames require memory accesses in order to perform motion compensation during decoding. This creates a more memory-intensive workload compared to an all-Intra frame GOP where prediction is solely compute intensive and no memory accesses are needed for prediction except when accessing cached neighbor pixels of each block. We are also varying input streams resolution from the lowest (QCIF), middle (CIF) to the highest (full HD). We used nine QCIF, nine CIF [14] and seven full HD [15] input streams. The rest of the parameters for our streams are kept in typical values. Table III summarizes the input streams configuration.

We obtain that CPI value is approximately constant for almost all input stream configurations of HEVC decoding benchmark running on a particular core type (Fig. 5). We also obtain that CPI value is always lower on a big core (lower is better). The last observation is in line with the proposed threshold in the previous subsection. We get for 49 out of 50 input stream configurations executed on a big core that the number of LLC misses per 10K instructions is

below the threshold and also that CPI on a big core is lower than CPI on a LITTLE core. Therefore, these configurations will be correctly scheduled on a big core. We get for only one input stream with full HD resolution and IBBBP frame GOP structure that the number of LLC misses is above the threshold and CPI on a big core is better. Therefore, only this input steam configuration will be wrongly scheduled on a LITTLE core.

We notice small difference in performance behavior for input streams that have either all Intra or IBBBP frame GOP structure when considering a particular core type (Fig. 6). Decoding B and P frames from IBBBP frame GOP requires large amount of memory accesses when fetching pixels of blocks within previously decoded frames while decoding I frames from all Intra frame GOP requires small amount of memory accesses because it fetches pixels of blocks only from the current frame. This causes that all input streams with all-Intra frame GOP will have lower CPI than the same input streams with IBBBP frame GOP on a big core and around 68% of them on a LITTLE core. We show that the first configuration has around 13% lower CPI in average for all input streams than the second one on a big core and around 2.3% lower CPI on a LITTLE core.

We also notice small difference in performance behavior for input streams that have either QCIF, CIF or full HD resolution when considering a particular core type (Fig. 7). The highest resolution (full HD) causes the most memory accesses on both core types while the middle resolution (CIF) causes the least memory accesses on a big core and the smallest resolution causes the least memory accesses on a LITTLE core in average for all input streams. When considering CPI, the lowest resolution causes highest CPI on both core types while the middle resolution causes the lowest CPI on a big core and the highest resolution causes the lowest CPI on a LITTLE core also in average for all input streams. We show that the input streams with middle resolution has around 16% lower CPI in average than the input streams with lowest resolution on a big core. We also

Table III
INPUT STREAMS CONFIGURATION

Parameter	Configuration
Encoder	x265
Profile, Level	Main
Resolution	QCIF, CIF, full HD
Throughput	25 fps, 100 frames
GOP	IBBBP or full Intra
Motion search method	Exhaustive search
Motion search range	57
Filters	Deblocking and Sample Adaptive Offset

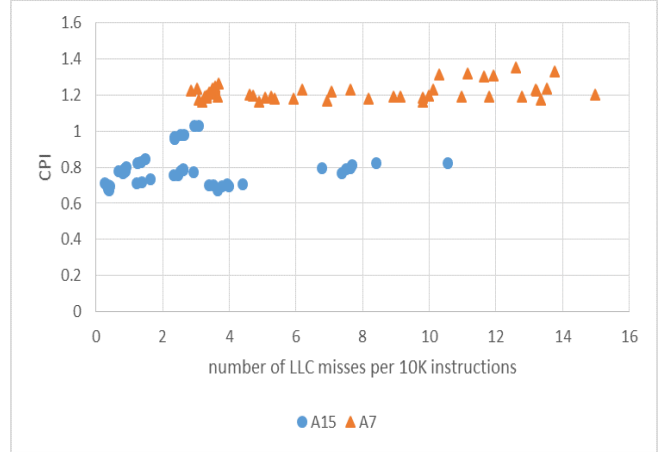


Figure 5. The CPI values for HEVC decoding executed on A15 and A7 cores.

show that the input streams with the highest resolution has around 5.8% lower CPI than the input streams with lowest resolution on a LITTLE core.

VI. CONCLUSION

In this paper we investigated how miss events affect overall performance of the ARM big.LITTLE architecture. We explained what performance impacts of particular miss events are relative to each other. We concluded that LLC misses have the greatest impact on performance. Thus, we constructed the model that establishes correlation between the number of LLC misses and CPI representing overall performance. This model enables us to propose soft threshold in terms of the number of LLC misses to determine on which core to schedule a certain application in order to maximize performance. All applications with the number of LLC misses above this threshold will be scheduled on

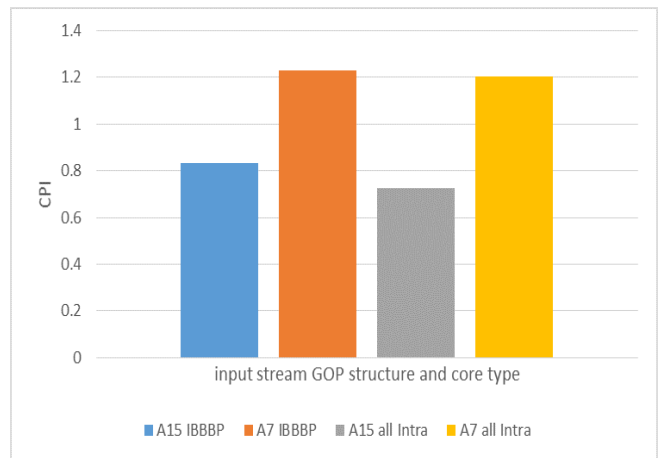


Figure 6. The average CPI values for HEVC decoding input streams with IBBBP and all Intra GOP structures executed on A15 and A7 cores.

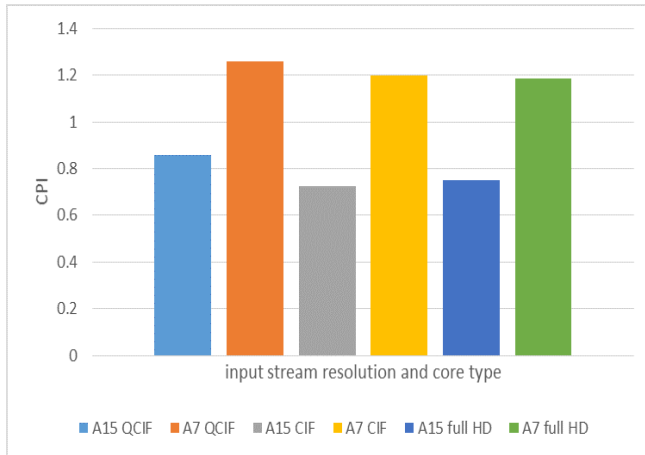


Figure 7. The average CPI values for HEVC decoding input streams with QCIF, CIF and full HD resolutions executed on A15 and A7 cores.

a LITTLE core because large number of LLC misses does not allow a big core to exploit ILP through out-of-order execution. Finally, we validated the proposed threshold with the popular HEVC decoding application. This approach could be completely applied on the other hardware platforms supporting performance monitoring counters.

More optimal scheduling would be achieved by dividing an application into phases and re-mapping these phases dynamically to different cores according to the number of LLC misses during them [6]. We can distinguish three different groups of phases. The phases with a small number of LLC misses belonging to the first group should be executed on a big core while the ones with a large number of LLC misses belonging to the second group should be executed on a LITTLE core. The phases with very large number of LLC misses belonging to the third group should be divided into chunks. The chunks with very large number of LLC misses should be executed on a big core to exploit MLP because independent LLC misses can access memory in parallel. Eventually, the remaining chunks should be executed on a LITTLE core. Almost all stress-ng benchmarks and all HEVC configurations have only phases that belong to the first two groups.

In future work we plan to further explore behavior of the big.LITTLE architecture with parallel threads and to quantify performance impacts of all mentioned miss events. We also plan to cluster hardware platforms according to particular miss events to maximize performance.

REFERENCES

[1] Odroid-xu4. (2013) Odroid-xu4 board. [Online]. Available: http://www.hardkernel.com/main/products/prdt_info.php

[2] Stress-ng. (2016) Stress-ng benchmark. [Online]. Available: <http://kernel.ubuntu.com/~cking/stress-ng/>

[3] K. Van Craeynest, A. Jaleel, L. Eeckhout, P. Narvaez, and J. Emer, "Scheduling heterogeneous multi-cores through performance impact estimation (pie)," in *ACM SIGARCH Computer Architecture News*, vol. 40, no. 3. IEEE Computer Society, 2012, pp. 213–224.

[4] M. Gottscho, S. Govindan, B. Sharma, M. Shoaib, and P. Gupta, "X-mem: A cross-platform and extensible memory characterization tool for the cloud," in *Performance Analysis of Systems and Software (ISPASS), 2016 IEEE International Symposium on*. IEEE, 2016, pp. 263–273.

[5] S. Eyerhan, L. Eeckhout, T. Karkhanis, and J. E. Smith, "A performance counter architecture for computing accurate cpi components," in *ACM SIGOPS Operating Systems Review*, vol. 40, no. 5. ACM, 2006, pp. 175–184.

[6] G. Patsilaras, N. K. Choudhary, and J. Tuck, "Efficiently exploiting memory level parallelism on asymmetric coupled cores in the dark silicon era," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 8, no. 4, p. 28, 2012.

[7] S. Eyerhan, J. E. Smith, and L. Eeckhout, "Characterizing the branch misprediction penalty," in *2006 IEEE International Symposium on Performance Analysis of Systems and Software*. IEEE, 2006, pp. 48–58.

[8] Cortex-A7. (2011) Cortex-a7 processor. [Online]. Available: <https://developer.arm.com/products/processors/cortex-a/cortex-a7>

[9] Cortex-A15. (2011) Cortex-a15 processor. [Online]. Available: <https://developer.arm.com/products/processors/cortex-a/cortex-a15>

[10] M. Rik. (2011) Deep inside arm's new intel killer. [Online]. Available: http://www.theregister.co.uk/2011/10/20/details_on_big_little_processing/

[11] G. J. Sullivan, J. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (hevc) standard," *IEEE Transactions on circuits and systems for video technology*, vol. 22, no. 12, pp. 1649–1668, 2012.

[12] R. Rodríguez-Sánchez and E. S. Quintana-Ortí, "Architecture-aware optimization of an HEVC decoder on asymmetric multicore processors," *CoRR*, vol. abs/1601.05313, 2016. [Online]. Available: <http://arxiv.org/abs/1601.05313>

[13] (2013) x265 hevc encoder. [Online]. Available: <https://bitbucket.org/multicoreware/x265/wiki/home>

[14] (2017) Yuv video sequences. [Online]. Available: <http://trace.eas.asu.edu/yuv/>

[15] (2013) Ultra video group test sequences. [Online]. Available: <http://ultravideo.cs.tut.fi/#testsequences>