

Received March 6, 2019, accepted April 3, 2019, date of publication April 11, 2019, date of current version April 30, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2910559

Microarchitecture-Aware Code Generation for Deep Learning on Single-ISA Heterogeneous Multi-Core Mobile Processors

JUNMO PARK^{1,2}, (Student Member, IEEE), YONGIN KWON²,
YONGJUN PARK³, (Member, IEEE), AND
DONGSUK JEON¹, (Member, IEEE)

¹Graduate School of Convergence Science and Technology, Seoul National University, Seoul 08826, South Korea

²System LSI, Samsung Electronics Co., Ltd., Hwaseong 18448, South Korea

³Department of Computer Science, Hanyang University, Seoul 04763, South Korea

Corresponding author: Dongsuk Jeon (djeon1@snu.ac.kr)

This work was supported in part by the National Research Foundation of Korea under Grant NRF-2019R1C1C1004927 and Grant NRF-2016R1C1B2016072, in part by the Information Technology Research Center Support Program under Grant IITP-2019-2018-0-01421, supervised by the Institute for Information and Communications Technology Promotion, and the System LSI Division of Samsung Electronics.

ABSTRACT While single-ISA heterogeneous multi-core processors are widely used in mobile computing, typical code generations optimize the code for a single target core, leaving it less suitable for the other cores in the processor. We present a microarchitecture-aware code generation methodology to mitigate this issue. We first suggest adopting Function-Multi-Versioning (FMV) to execute application codes utilizing a core at full capacity regardless of its microarchitecture. We also propose to add a simple but powerful backend optimization pass in the compiler to further boost the performance of applicable cores. Based on these schemes, we developed an automated flow that analyzes the program and generates multiple versions of hot functions tailored to different microarchitectures. At runtime, the running core chooses an optimal version to maximize computation performance. The measurements confirm that the methodology improves the performance of Cortex-A55 and Cortex-A75 cores in Samsung's next-generation Exynos 9820 processor by 11.2% and 17.9%, respectively, while running TensorFlow Lite.

INDEX TERMS Edge computing, function multi-versioning, single-ISA heterogeneous multi-core.

I. INTRODUCTION

A. DEEP LEARNING ON SINGLE-ISA HETEROGENEOUS MULTI-CORE PROCESSORS

The big.LITTLE architectures consisting of multiple cores with different microarchitectures are widely adopted in mobile environments where the tradeoff between power efficiency and performance is a critical issue. In order to fully exploit the benefit of those architectures, operating systems need to handle each program appropriately. On the Android platforms, tasks are divided into background, foreground, and top-app. The type of each task determines where it should be executed; on most of big.LITTLE systems, top-app, and foreground tasks can run on any core, whereas background tasks can run only on little core to save power consumption. For instance, if a task runs in the background, the scheduler will

assign the task to the little core group. In addition, many optimized schedulers such as an energy aware scheduler (EAS) have been applied to dynamically select an appropriate core in the assigned core group by scrutinizing the usage pattern and computational requirements [1]–[4]. The type of task may be changed from background to top-app or vice versa by user input, making it impossible to predict which core group the given task will be assigned to.

As deep learning algorithms have evolved in the last few years, developers have made many efforts to accelerate an inference process and improve energy efficiency in mobile environments. For instance, Google proposed multiple convolutional neural network (CNN) models optimized for real-time processing in mobile systems, accompanied by quantized models to further reduce overheads [5], [6].

The most performance-critical part of the inference using deep neural networks is a general matrix multiplication (GEMM). To accelerate GEMM operations, developers

The associate editor coordinating the review of this manuscript and approving it for publication was Xinyu Du.

often apply hand-tuned assembly codes or intrinsic functions such as NEON in ARM architectures. For instance, two types of libraries are used for maximizing GEMM performance in the Tensorflow Lite infrastructure: 32-bit floating-point models use the Eigen library [7] based on ARM-NEON intrinsic functions, whereas 8-bit quantized models employ the gemmlowp library [8] consisting of hand-tuned assembly codes.

Compilers typically generate assembly codes with different patterns depending on the microarchitecture characteristics for maximal performance. Since the compiler-generated code pattern is optimized only for a target core, execution of the code may undergo performance degradation on the other cores in the processor with a big.LITTLE configuration.

The GCC compiler currently supports `big.LITTLE -mcpu` option in order to address this issue, which tries to optimize the code considering hardware characteristics of both big and little cores simultaneously. However, the resulting code still exhibits inferior performance to the codes solely optimized for each core due to architectural differences. For big.LITTLE architectures with big out-of-order and little in-order cores, traditional static compilers such as GCC generate assembly codes primarily focusing on little cores that do not have sufficient hardware resources such as a register renaming unit and a re-order buffer. As a result, significant performance degradation may be incurred when the code is executed in big cores.

B. PROPOSED APPROACH

To solve the aforementioned code inefficiency problem, we propose a novel code generation methodology for single-ISA heterogeneous multi-core mobile processors aimed at deep learning applications. We first suggest adopting Function-Multi-Versioning (FMV) to execute application codes utilizing a core at full capacity regardless of its microarchitecture. The concept of FMV was initially proposed several decades ago [9], but it has not been adopted broadly due to large code space overhead from generating multiple copies of the code. However, for mobile deep learning frameworks such as Tensorflow Lite, it is observed that a program is spending most of its time performing only a few operations such as GEMM. Under this observation, we find that applying FMV to only a set of hot functions enhances performance noticeably while imposing very little code space overhead.

We also propose to add a simple but powerful back-end optimization pass in the compiler to further boost the performance of smaller cores. By modifying the baseline AArch64 load-store optimization pass with the introduction of a load split pass, the code generation scheme achieves significant performance improvement for smaller cores running GEMM.

In order to apply these techniques to general systems, we also develop an automatic microarchitecture-aware code generation flow. It first analyzes a target program and selects candidate functions from the list of functions sorted by execution frequency based on the profiling result. Then the flow

Algorithm 1 Load Split Pass

1. **function** load-split-pass(loop $L \in$ in function)
 2. **for** each instruction $I \in$ basic blocks in L
 3. **if** $I == 128\text{-bit_load}$ **then**
 4. Find available register $r \in$ GPR64RegClass
 5. **if** exist(r)**then**
 6. Create $I_{64\text{-bit_load}}$ with $R_{source} = R_{source,I}$,
 $R_{dest} = R_{dest,I}$
 7. Create $I_{64\text{-bit_load}}$ with $R_{source} = R_{source,I} + 8$,
 $R_{dest} = r$
 8. Create $I_{64\text{-bit_mov}}$ with $R_{source} = r$ and $R_{dest} = R_{dest,I}$
 9. Remove I
 10. **end if**
 11. **end if**
 12. **end for**
-

clones the target functions and inserts a runtime selector, resulting in target-specific assembly codes after compilation. At runtime, the selector checks IDCODE in the Performance Monitoring Unit (PMU) register of the current core and chooses an optimal version among the clones.

The proposed code generation flow was tested on Samsung Exynos 8895 processor as well as next-generation Exynos 9820 processor. Measurements show a performance improvement of 12.7% for Exynos-M2 core in Exynos 8895, and performance boosts of 11.2% and 17.9% for Cortex-A55 and A75 cores in Exynos 9820, respectively, confirming that the proposed methodology is effective for processing deep learning models on single-ISA heterogeneous multi-core processors.

II. RELATED WORKS AND MOTIVATION

A. FUNCTION MULTI VERSIONING (FMV)

FMV has been extensively studied in the last few years in order to maximize performance while suppressing code space increase. A program may run differently depending on the input data pattern due to unexpected control flows from conditional statements and varying call paths in a function. The conventional FMV technique profiles the program and generates multiple copies of the code, each optimized for specific input data pattern, through feedback-driven program optimization. This process may be repeated until an optimal point is reached [10], [11]. Some works also suggest applying the multi-versioning technique to smaller code regions such as a loop or set of basic blocks [12], [13].

Although conventional FMV schemes shorten the execution time of the program through adaptive code selection at runtime, deep learning algorithms usually exhibit very regular data patterns, making input-dependent FMV less attractive. In addition, it still optimizes the code only for a single target core although each microarchitecture may have largely different hardware configurations.

The LLVM compiler [14], which is widely used in mobile environments, currently supports optimizations of target

programs for a specific microarchitecture. If the compiler obtains target microarchitecture information using the *-mcpu* option, the compiler optimizes the code considering hardware resources of the given microarchitecture. For instance, the compiler can enforce different register allocation policies based on the information about the target microarchitecture. For in-order machines, the register numbers cannot be reused due to lack of register renaming unit in the microarchitecture. However, out-of-order machines usually have register renaming units, and hence the compiler can employ an improved instruction scheduling through reusing register numbers.

B. GENERAL MATRIX MULTIPLICATION (GEMM)

GEMM is considered to be a key library in deep learning algorithms based on deep neural networks. For minimizing inference time on mobile devices, some prior works proposed to optimize the neural networks (e.g., Mobilenets [5] and Mnasnet [15]), whereas others try to reduce arithmetic calculation overhead by replacing costly 32-bit floating-point operations with 16-bit operations or even smaller fixed-point operations [6], [16].

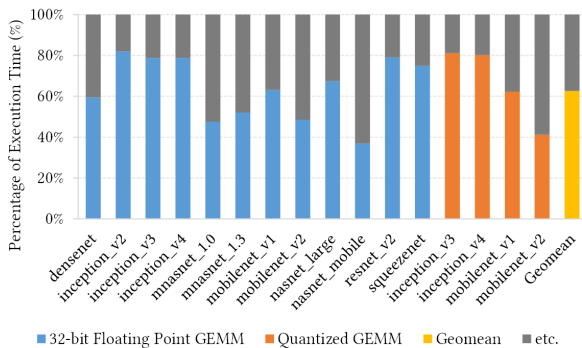


FIGURE 1. Proportion of GEMM libraries in execution time.

In this work, for all analyses we use the Tensorflow Lite which is a lightweight version of TensorFlow aimed at running machine learning models on mobile and embedded devices. Tensorflow Lite supports both 32-bit floating-point and 8-bit quantized models. Figure 1 shows the percentage of GEMM libraries in execution time when different CNN models are processed using Tensorflow Lite on Exynos 9820 processor, which implies that the main bottleneck can be only a few functions, and they consume most of the computation time. Specifically, the hottest function alone takes 64% and 61% of the runtime in 32-bit floating-point and 8-bit quantized models, respectively. This result suggests that if FMV is applied only to those functions, we could obtain significant performance improvement without code bloat. Note that some 32-bit floating-point models such as Mobilenet [5] and Mnasnet [15] have the second hottest function due to their specific neural network architectures that require depth-wise convolutions [5]. Table 1 displays representative models used for analyses throughout the paper, which are extracted from the hosted models and have varying neural network

Table 1. Selected test models of TensorFlow Lite.

	Model	Top-1 Accuracy	Top-5 Accuracy
32-bit Floating-point	densenet	64.2%	85.6%
	inception_resnet_v2	77.5%	94.0%
	inception_v3	77.9%	93.8%
	inception_v4	80.1%	95.1%
	mnasnet_1.0	74.1%	91.8%
	mnasnet_1.3	75.2%	92.6%
	mobilenet_v1	71.0%	89.9%
	mobilenet_v2	71.8%	90.6%
	nasnet_large	82.6%	96.1%
	nasnet_mobile	73.9%	91.5%
	resnet_v2	76.8%	93.6%
	squeezenet	49.0%	72.9%
8-bit Quantized	inception_v3	77.5%	93.7%
	inception_v4	79.5%	93.9%
	mobilenet_v1	70.0%	89.0%
	mobilenet_v2	70.8%	89.9%

architectures and input sizes of $224 \times 224 \times 1$ or larger. Other hosted models are mostly derivative of the selected models.

C. MICROARCHITECTURE-AWARE CODE GENERATION

Both GCC and LLVM currently support *-mcpu* option for microarchitecture-aware code generation. Without such an option, the compiler generates a generic code for generalized architecture such as AArch64 and ARMv7-a. This is equivalent to using *-mcpu = generic* option at compile time.

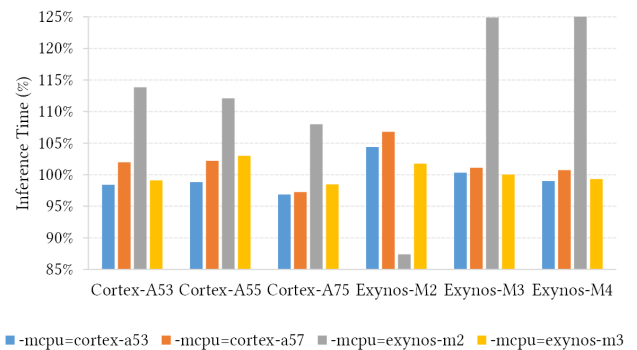


FIGURE 2. Measured inference time on various cores with different -mcpu options used for compilation.

Figure 2 shows the measurement results when we compile Tensorflow Lite with different *-mcpu* options and run the models from Table 1 on various cores. Note that the baseline is *-mcpu = generic* which is being used by most Android developers. Google provides a toolchain for building Android applications called Native Development Kit (NDK). Although NDK does support *-mcpu* option, in most cases developers use a generic option since there exists a wide range of SoCs with different microarchitectures. In Figure 2, if we use *-mcpu = exynos-m2* option during compilation, a significant improvement of 12.6% is observed for Exynos-M2

core in Exynos 8895 processor. However, that option results in 13.8% performance degradation for Cortex-A53 core in the same Exynos 8895 processor, confirming that a simple microarchitecture-aware code generation is not effective in heterogeneous multi-core processors.

III. PROPOSED FMV SCHEME FOR HETEROGENEOUS MULTI-CORE PROCESSORS

As discussed earlier, the conventional FMV scheme generates multiple versions of a given code but does not allow dynamic version change when a task migrates to a different core. In other words, if a code is optimized for one core, the code may suffer from performance degradation on the other cores.

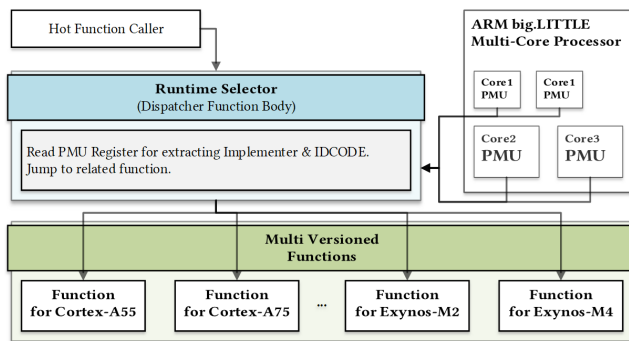


FIGURE 3. Proposed FMV scheme.

To resolve this issue, we propose a new FMV scheme which enables dynamic version change of target code regions at runtime so that each core can execute the version solely optimized for that core. Figure 3 depicts the proposed FMV scheme. Multiple versions of a target function are generated where each version is optimized differently by the compiler considering the target microarchitecture, and a runtime selector is inserted into the program.

For the proposed FMV scheme, the program must know which core it is currently located in at runtime to select a proper version. Hence, the program needs access to PMU to retrieve such information. PMU is an optional feature for some architectures such as ARMv8-A, but the feature is still strongly recommended by ARM [17] and most SoC manufacturers integrate the PMU block in the processor.

A. RUNTIME SELECTOR

After a scheduler assigns a program to an appropriate core, the program must identify the running core for FMV. We insert a runtime selector in the program so that the runtime selector dynamically chooses one of the versions that fits the current microarchitecture during execution. The selector can be realized with a few conditional statements followed by additional function calls. Conventional input-aware FMV schemes also employ a runtime selector in the flow, but the selector must analyze input data in detail to determine an optimal version, which translates to large processing overheads [18], [19]. However, our runtime selector chooses a

version by simply looking at the implementer (IMP) and identification (IDCODE) information retrieved from the PMU. Since the runtime selector always chooses the same version unless the program migrates to a different core, the pattern can be well predicted by the branch predictor in the core, and pipeline stalls due to branch prediction misses can be avoided. The only overhead of the runtime selector is reading PMU register and comparing it against predefined values, which does not impose noticeable performance degradation in real-world experiments as demonstrated in Section 5. The runtime selector is inserted as an inline assembly and Code 1 shows an example runtime selector.

CODE 1. Runtime selector example.

```

1.  function runtime_selector()
2.    pmu ← get_pmu_register()
3.    imp ← get_implementor_from_pmu_register(pmu)
4.    idcode ← get_idcode_from_pmu_register(pmu)
5.    case imp of
6.      ARM:
7.        case idcode of
8.          Cortex-A55:
9.            function_for_cortex-a55()
10.         Cortex-A75:
11.           function_for_cortex-a75()
12.         end case
13.      SAMSUNG:
14.        case idcode of
15.          Exynos-M2:
16.            function_for_exynos-m2()
17.          Exynos-M3:
18.            function_for_exynos-m3()
19.          end case
20.      Default:
21.        function_for_generic()
22.    end case

```

B. MULTI-VERSIONED FUNCTIONS

In Tensorflow Lite, GEMM algorithm is implemented using two libraries: gemmlowp [8] and Eigen library [7]. Contrary to the gemmlowp library which is implemented in hand-tuned assembly, main functions of the Eigen library are written in ARM NEON intrinsics and hence can be optimized by the compiler. For those functions, the compiler selects NEON instructions and performs back-end optimizations including register allocation, instruction scheduling, and peephole optimizations such as load-store optimization. For FMV, target functions are cloned and additional function attributes stating target microarchitecture are inserted as shown in Code 2.

CODE 2. Multi-versioned function example.

```

1.  __attribute__((target("arch=cortex-a55")))
2.  function function_for_cortex-a55 ( )
3.  __attribute__((target("arch=exynos-m2")))
4.  function function_for_exynos-m2 ( )

```

With the added attributes, the compiler can recognize the target microarchitecture of the function even if we use `-mcpu = generic` option or leave it empty. This information is also transferred to the compiler backend for applying target specific optimizations.

IV. LOAD SPLIT OPTIMIZATION AND LLVM-BASED UNIFIED AUTOMATIC CODE GENERATION

While the FMV scheme allows conventional compilers to optimize functions for multiple target microarchitectures, in this section we propose an additional optimization technique for the backend of the LLVM compiler to further enhance performance. The *neon-gemm-kernel-benchmark* for the *gemmlowp* library provides an example hand-tuned assembly code for ARM Cortex-A55 core, where a single 128-bit load instruction is replaced with three separate instructions as shown in Code 3.

CODE 3. Load split assembly example.

```

1.  "ldr  q0, [x0]" // Split this as follows
2.  "ldr  d0, [x0]"
3.  "ldr  x18, [x0, #8]"
4.  "ins  v0.d[1], x18"
    
```

Cortex-A55 microarchitecture does not allow loading 128 bits at once. Hence, the core internally realizes the 128-bit load operation by loading 64 bits twice, which cannot be executed in the same cycle although Cortex-A55 supports dual-issue. On the other hand, manually splitting load instructions as shown in Code 3 and placing other arithmetic instructions between them can mitigate pipeline stalls by processing a load and an arithmetic operation in the same cycle through dual-issue. We refer to this scheme as load split. Inspired by this observation, we propose a backend optimization pass for the LLVM compiler.

A. LLVM BACKEND PASSES

The LLVM compiler backend has more than 100 optimization passes, where the instruction selection and load-store optimization passes play an important role in boosting GEMM performance. Since GEMM operations require frequent matrix data load, the overall performance is largely affected by which load instructions are selected [20]. The load-store optimization pass attempts to combine instructions, searching for contiguous loads or stores that can be combined into a single instruction as depicted in Code 4.

CODE 4. Load pairing example.

```

1.  "ldr  q0, [x0]"
2.  "ldr  q1, [x0, #16]" // Combine these as follows
3.  "ldp  q0, q1, [x0]"
    
```

Combining instructions reduces code size and lowers instruction fetch and decode overhead on the core. However,

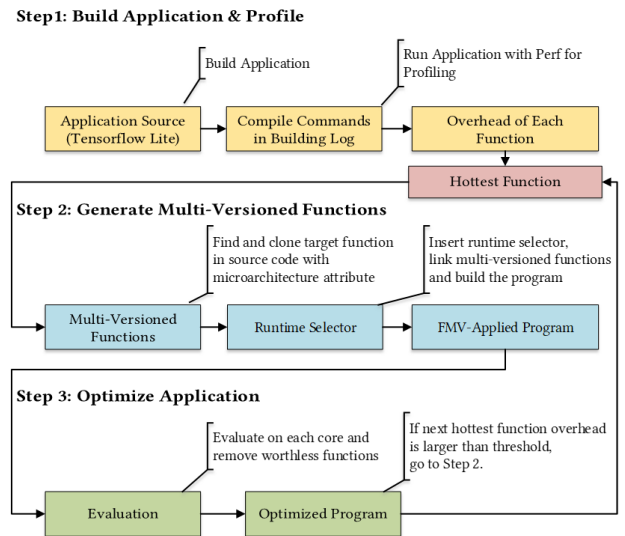


FIGURE 4. Proposed automatic code generation flow.

this optimization reverses the effect of load split and, therefore, this pass is turned off for Cortex-A55 and other similar microarchitectures in our flow.

B. LOAD SPLIT OPTIMIZATION PASS

We implement an additional optimization pass for load split after the load-store optimization pass. The load split optimization can be applied to any load instructions in a function, but we enforce the pass only for loops to minimize code size increase. A typical workflow is described in Algorithm 1.

This optimization provides maximal benefit when enough number of other independent arithmetic instructions can be placed between newly created load instructions. For example, the hottest function of Tensorflow Lite consists of 32 128-bit loads and 96 SIMD-FP-multiply-accumulate instructions, providing a good ratio between load and arithmetic operations. However, the 2nd hottest function consists of 12 128-bit loads, 4 SIMD-FP-multiply and 4 SIMD-FP-add, suggesting less performance boosting due to load split optimization.

After the load split optimization pass, we need to run instruction scheduling again in order to move arithmetic instructions into the space between the split instructions. LLVM already has the *PostRAScheduler* feature that schedules again after register allocation. Currently, the feature is not activated for Cortex-A55/75 in LLVM 7.0, which is the most up-to-date version, and we reactivate this feature in the compilation step of the proposed flow.

C. UNIFIED AUTOMATIC MICROARCHITECTURE-AWARE CODE GENERATION FLOW

The proposed code generation methodology should be applied to only a subset of functions in order to suppress code size increase. We present a unified automatic code generation flow depicted in Figure 4. The flow first builds the target program while collecting build logs which are later used for

tracking source files to be modified. Then the flow profiles the program using Linux Perf [21], which calculates the overhead of each function and sorts the functions by their hotness. In the next step, the function at the top of the list is selected and cloned through FMV. Each clone is compiled with an optimal target microarchitecture directive as well as goes through the load split pass in the backend depending on the target microarchitecture. Finally, the resulting program is profiled on each core and the flow removes the versions that do not exhibit performance improvement. The flow continues onto the next function in the list if its share in runtime is higher than a threshold.

V. EXPERIMENTAL EVALUATION

A. EXPERIMENTAL SETUP

In this work, we evaluate the flow using Tensorflow Lite framework which is a strong candidate for deep learning applications on mobile platforms. We run 12 floating point models and 4 quantized models which have different structures and use various sizes of memory. The flow was tested on two mobile processors: Samsung Exynos 8895 and next-generation Exynos 9820 processors. The Exynos 9820 features three processor clusters to realize a big.LITTLE architecture, each consisting of four Cortex-A55 cores, two Cortex-A75 cores, and two Exynos-M4 cores, respectively. The Exynos 8895 has a typical ARM big.LITTLE microarchitecture and includes two clusters in total, each with four Cortex-A53 cores and four Exynos-M2 cores, respectively.

For evaluation, tasks are processed on a specific core using *taskset* command. The frequencies of processors and memory interface are all fixed during experiments for reliable measurements. Since the two processors require different versions of Android OS, Ubuntu is used instead for experiments to align the experiment environments. We use *chroot* tool for running Ubuntu on each device and building Tensorflow Lite for Linux.

B. FLOATING-POINT MODELS

We first evaluate the proposed flow using the floating-point models in the Tensorflow Lite benchmark. Those models employ the Eigen library based on intrinsic functions. Figure 5 and 6 represent the relative runtime of each model using our code generation flow for Exynos 8895 and 9820 processors, respectively, compared to the baseline in which a generic option is used for compilation. In our flow, the FMV feature is used for all experiments in order to generate multiple versions solely optimized for each microarchitecture, and the backend optimization is enabled for Cortex-A55 and A75 cores in Exynos 9820.

For Exynos 8895, the big core (Exynos-M2) exhibits noticeable performance improvements across all models, with an average of 13.0%. For Exynos 9820, the little and middle cores (Cortex-A55/A75) show apparent performance boosting of 11.2% and 17.9%, respectively, on average. Since the load split technique is enabled for those cores, the

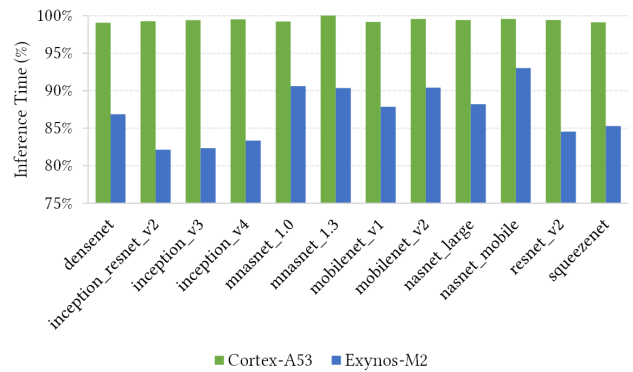


FIGURE 5. Performance improvement of cores in Exynos 8895 processor.

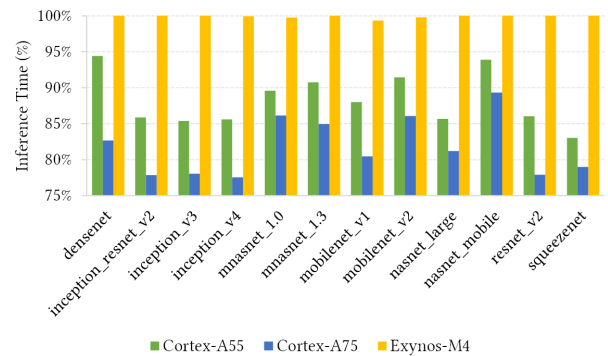


FIGURE 6. Performance improvement of cores in Exynos 9820 processor.

improvements are the combined effects of using an optimal compilation directive for each clone of a target function and applying additional backend optimization. Figure 7 displays the runtime reduction after applying compilation directives only, and after using both techniques. With compilation directives, the throughput of Cortex-A55 and A75 are improved by 3.0% and 2.9%, respectively, whereas the compiler backend optimization enhances the performance further by 8.2% and 15.0%, respectively.

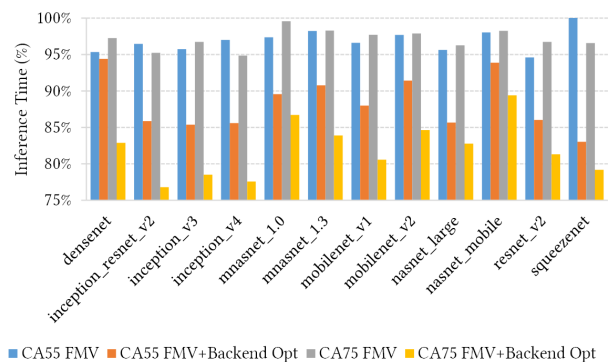


FIGURE 7. Effects of optimal compilation directive and load split optimization.

In the aforementioned experiments, the flow creates 6 clones of the target functions in order to support Cortex-A53/A55/A75 and Exynos-M2/M3/M4 cores altogether. However, the size of the generated binary increases by only 3.66%.

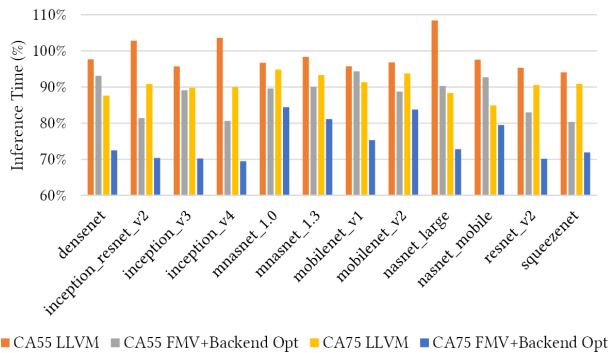


FIGURE 8. Performance improvements over GCC.

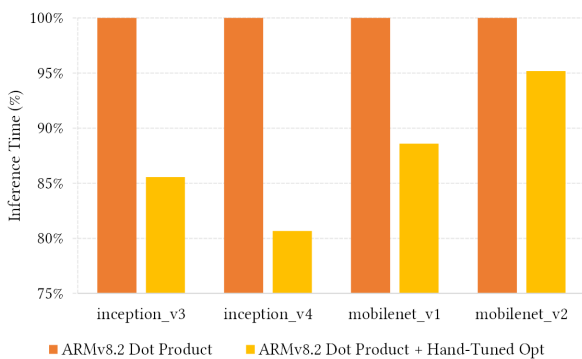


FIGURE 9. Assembly code optimization for gemmlowp library using load split technique.

We also compiled Tensorflow Lite using GCC 8.2 and measured the performance on Cortex-A55 and A75 cores in Exynos 9820 (Figure 8). On Cortex-A55, LLVM and LLVM with load-split optimization show 1.47% and 12.2% better performance than GCC, respectively, on average. On Cortex-A75, LLVM and LLVM with load-split optimization exhibit 9.49% and 24.88% better performance than GCC. We also tested GCC 7.2 for completeness, but it showed even lower performance than GCC 8.2. Google decided to deprecate GCC in their toolchain in October 2016 [22] and hence most of the mobile applications are expected to be built by LLVM.

C. QUANTIZED MODELS

The quantized models use the gemmlowp library, which relies on hand-tuned assembly codes. Therefore, automated code generation flow cannot be directly applied. However, we can still manually adopt the proposed FMV and load split optimization techniques to improve computation performance. Figure 9 shows that the techniques raise the performance by 12.5% when applied to Cortex-A55 core. If the same code is used for Exynos-M4 accompanied in Exynos 9820, the performance drops by 3.5%, and this is remedied by applying the proposed FMV in the assembly, confirming the effectiveness of our flow. Note that we modified the Tensorflow Lite build option to enable ARMv8.2 Dot Product feature.

VI. CONCLUSION

In this work, we present a microarchitecture-aware code generation technique for single-ISA heterogeneous multi-core processors. By applying FMV and introducing additional

optimization pass in the compiler backend, both Exynos 8895 and the next-generation Exynos 9820 processors exhibit significant performance improvements for deep learning applications on all of 16 famous models. Although the flow was tested for two processors, the Cortex-A55 and A75 cores are expected to be used in most mobile SoCs in the near future. In order to adopt the flow in the current Android environments, it is essential to access PMU registers from user level. This is not allowed in the current OS version, but we hope this change is made in the next releases so that deep learning applications can be efficiently accelerated on mobile processors.

REFERENCES

- [1] ARM and Linaro. *Energy Aware Scheduling (EAS)*. Accessed: Dec. 2018. [Online]. Available: <https://developer.arm.com/open-source/energy-aware-scheduling>
- [2] R. Kumar, D. M. Tullsen, P. Ranganathan, N. P. Jouppi, and K. I. Farkas, "Single-ISA heterogeneous multi-core architectures for multithreaded workload performance," *ACM SIGARCH Comput. Archit. News*, vol. 32, no. 2, p. 64, 2004.
- [3] K. Van Craeynest, S. Akram, W. Heirman, A. Jaleel, and L. Eeckhout, "Fairness-aware scheduling on single-ISA heterogeneous multi-cores," in *Proc. 22nd Int. Conf. Parallel Archit. Compilation Techn. (PACT)*, Sep. 2013, pp. 177–187.
- [4] J. Corbet. (2016). *Scheduling for Android Devices*. Linux Plumbers Conference. [Online]. Available: <https://lwn.net/Articles/706374/>
- [5] A. G. Howard et al., "MobileNets: Efficient convolutional neural networks for mobile vision applications," *CoRR*, vol. abs/1704.04861, 2017.
- [6] B. Jacob et al., "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018.
- [7] G. Guennebaud et al. (2010). *Eigen V3*. [Online]. Available: <http://eigen.tuxfamily.org>
- [8] B. Jacob. (2018). *GEMMLOWP: A Small Self-Contained Low-Precision GEMM Library*. [Online]. Available: <https://github.com/google/gemmlowp>
- [9] K. D. Cooper, M. W. Hall, and K. Kennedy, "Procedure cloning," in *Proc. Int. Conf. Comput. Lang.*, Apr. 1992, pp. 96–105.
- [10] D. Chen et al., "Taming hardware event samples for FDO compilation," in *Proc. 8th Annu. IEEE/ACM Int. Symp. Code Gener. Optim.*, vol. 10, 2010, pp. 42–52.
- [11] D. Chen and D. X. Li, "AutoFDO: Automatic feedback-directed optimization for warehouse-scale applications," in *Proc. Int. Symp. Code Gener. Optim.*, 2016, vol. 1, no. 212, pp. 12–23.
- [12] X. Chen and S. Long, "Adaptive multi-versioning for OpenMP parallelization via machine learning," in *Proc. 15th Int. Conf. Parallel Distrib. Syst. (ICPADS)*, vol. 12, 2009, pp. 907–912.
- [13] K. Koukos, P. Ekemark, G. Zacharopoulos, V. Spiliopoulos, S. Kaxiras, and A. Jimborean, "Multiversioned decoupled access-execute: The key to energy-efficient compilation of general-purpose programs," in *Proc. 25th Int. Conf. Compiler Construct.*, 2016, pp. 121–131.
- [14] C. Lattner and V. Adve, "LLVM: A compilation framework for lifelong program analysis & transformation," in *Proc. Int. Symp. Code Gener. Optim. (CGO)*, Mar. 2004, pp. 75–86.
- [15] M. Tan, B. Chen, R. Pang, V. Vasudevan, and Q. V. Le, "MnasNet: Platform-aware neural architecture search for mobile," *CoRR*, vol. abs/1807.1, 2018.
- [16] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks for mobile devices," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016.
- [17] *ARM Architecture Reference Manual ARMv8, for ARMv8-A Architecture Profile*, ARM, 2018.
- [18] K. Tian, Y. Jiang, E. Z. Zhang, and X. Shen, "An input-centric paradigm for program dynamic optimizations," in *Proc. ACM Int. Conf. Object Oriented Program. Syst. Lang. Appl. (OOPSLA)*, 2010, vol. 45, no. 10, pp. 125–139.
- [19] P.-F. Chuang, H. Chen, G. F. Hoflehner, D. M. Lavery, and W.-C. Hsu, "Dynamic Profile Driven Code Version Selection," in *Proc. 11th Annu. Workshop Interact. Between Compil. Comput. Archit.*, 2007, pp. 74–81.

- [20] X. Su, X. Liao, and J. Xue, "Automatic generation of fast BLAS3-GEMM: A portable compiler approach," in *Proc. IEEE/ACM Int. Symp. Code Gener. Optim. (CGO)*, Feb. 2017, pp. 122–133.
- [21] *Linux Perf*. Accessed: Dec. 2018. [Online]. Available: <https://perf.wiki.kernel.org/>
- [22] Google. (2019). *NDK Revision History*. Accessed: Jan. 2019. [Online]. Available: https://developer.android.com/ndk/downloads/revision_history



JUNMO PARK received the B.S. degree in computer science from Kwangwoon University in 2012. After that, he worked at Samsung Electronics for around 7 years, working on Compiler Optimization and Development. He is currently pursuing the M.S. degree with the Graduate School of Convergence Science and Technology at Seoul National University, South Korea. His research interests include deep learning, compiler, embedded systems, HW/SW co-design and optimizations.



YONGIN KWON received the B.Sc. degree in electrical and electronic engineering from the Korea Advanced Institute of Science and Technology, South Korea, in 2008, and the M.S. and Ph.D. degrees in electrical and computer engineering from Seoul National University, South Korea, in 2010 and 2015, respectively. He is currently a Software Engineer with Samsung Electronics. His research interests include mobile cloud computing, compiler, deep learning, and embedded systems.



YONGJUN PARK received the Ph.D. degree in electrical engineering from the University of Michigan, Ann Arbor, MI, USA, in 2013. He is currently an Assistant Professor with the Department of Computer Science, Hanyang University, Seoul, South Korea. His research interests include compilers and computer architectures for various computer systems.



DONGSUK JEON (S'10–M'15) received the B.S. degree in electrical engineering from Seoul National University, Seoul, South Korea, in 2009, and the Ph.D. degree in electrical engineering from the University of Michigan, Ann Arbor, MI, USA, in 2014. From 2014 to 2015, he was a Post-doctoral Associate with the Massachusetts Institute of Technology, Cambridge, MA, USA. He is currently an Assistant Professor with the Graduate School of Convergence Science and Technology, Seoul National University. His current research interests include energy-efficient signal processing, low-power circuit, and SoC for mobile applications.

...