

# The Importance of Worst-Case Memory Contention Analysis for Heterogeneous SoCs

Lorenzo Carletti, Gianluca Brilli, Alessandro Capotondi, Paolo Valente, Andrea Marongiu

**Abstract**—Memory interference may heavily inflate task execution times in Heterogeneous Systems-on-Chips (HeSoCs). Knowing worst-case interference is consequently fundamental for supporting the correct execution of time-sensitive applications. In most of the literature, worst-case interference is assumed to be generated by, and therefore is estimated through read-intensive synthetic workloads with no caching.

Yet these workloads do not *always* generate worst-case interference. This is the consequence of the general results reported in this work. By testing on multiple architectures, we determined that the highest interference generation traffic pattern is actually hardware dependant, and that making assumptions could lead to a severe underestimation of the worst-case (in our case, of more than 9×).

## I. INTRODUCTION

Heterogeneous on-chip Systems (HeSoC) combine the benefits of low power consumption Systems on Chip (SoC) with the ability to use accelerators to execute specialized workloads efficiently. Commercial-off-the-shelf (COTS) HeSoCs often make use of GP-GPU or FPGA as accelerators, combined with multi-core general-purpose *host* CPUs. This provides both specialization and flexibility.

These systems typically rely on a shared-memory organization, where the aforementioned compute units are interconnected through a shared bus to the main system DRAM, as shown in Fig. 1. As the number of compute engines grows in the chip, the main memory is subject to increasing contention.

This potentially causes the tasks executing on the various units to experience decreased bandwidth and, as a consequence, an increased execution time due to mutual interference [1]. This is particularly problematic for time-sensitive applications. The problem has been extensively studied before [2]–[7], and several different approaches to mitigate the effects of memory interference have been developed [8]–[12].

In order to prove the validity of an interference mitigation solution, it's important to have a proper understanding of the timing effects (i.e., slowdown) a workload under test experiences in the worst case, and showing that a particular approach can handle such a scenario. We define as worst-case the program which is slowed down the most by other tasks causing DRAM interference.

Some previous works [4], [6], [12] make the assumption that read-only synthetic benchmarks, which are programs executing memory operations at the highest speed possible, have to be the ones which either:

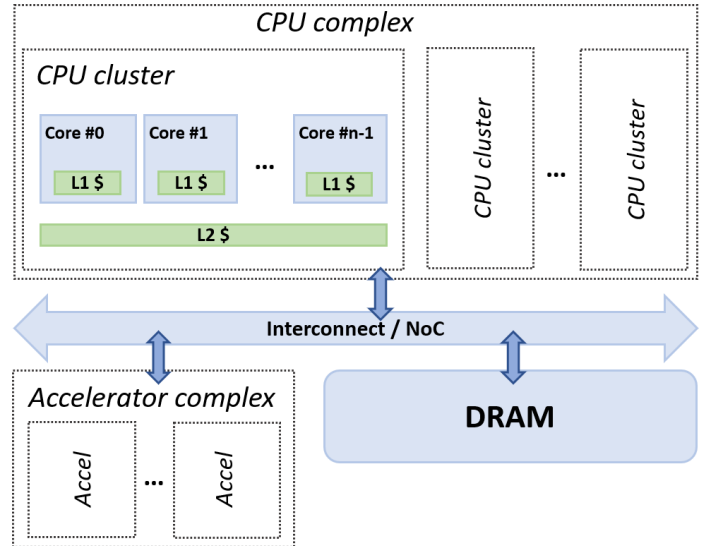


Fig. 1: Shared memory architectural template.

- 1) Cause the highest amount of DRAM interference for other concurrently running tasks.
- 2) Are the most slowed down by the effects of DRAM interference.

Using such concepts without a proper study of the hardware is particularly problematic, since they could lead to a wrong characterization of the worst-case, which, in turn, could mean that certain conditions might cause slowdowns greater than expected for important time-sensitive applications.

### A. Contributions

Our research analyzed the effects of different synthetic memory-intensive benchmarks running alongside other tasks (Polybench) on two separate HeSoCs: the FPGA-based Xilinx ZU9EG and the GPU-based NVIDIA TX2. Our results prove that:

- 1) The type of program causing the highest amount of interference is actually hardware-dependant.
- 2) The type of program most slowed down by the effects of DRAM interference is not guaranteed to be read-only synthetic benchmarks.

## II. SYNTHETIC BENCHMARKS

Synthetic benchmarks are programs which can be fine-tuned to emulate different kinds of programs based on

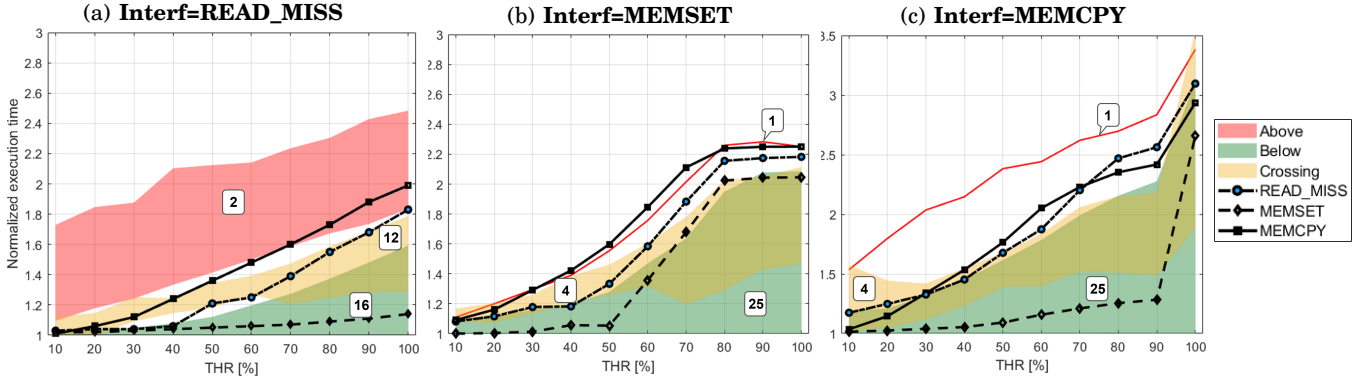


Fig. 2: **NVIDIA TX2**. Execution time increase of the three synthetic benchmarks (curves) and the Polybench benchmarks (solored areas) running on the core under test with increasing interference from the other cores (THR%). The workload executed by the interference cores is indicated above the plot.

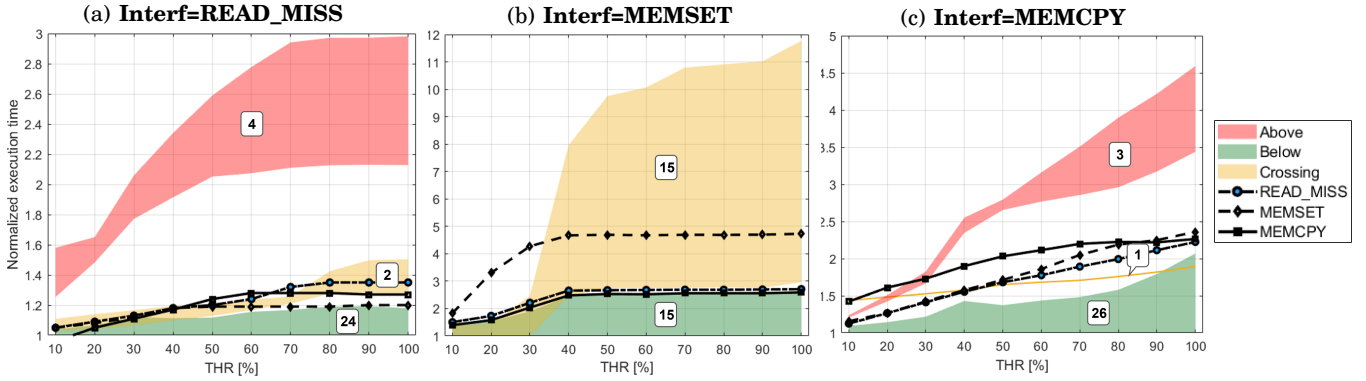


Fig. 3: **Xilinx ZU9EG**. Execution time increase of the three synthetic benchmarks (curves) and the Polybench benchmarks (solored areas) running on the core under test with increasing interference from the other cores (THR%). The workload executed by the interference cores is indicated above the plot.

certain parameters. For our tests on DRAM interference generation, we decided on using the following configurations:

- *READ\_MISS*: Loads only, at the maximum possible speed. What was believed to be the worst-case in Section I. It produces a 100 % read traffic.
- *MEMSET*: A series of consecutive stores, executed at the maximum possible speed. It produces a 100 % write traffic.
- *MEMCPY*: Consecutive loads and stores. It produces a 50% read - 50 % write traffic.

These synthetic benchmarks were then set to run on three CPU cores on the two platform as *Background Interference Generators*. On the remaining core, another task (either a Polybench, or one of the Synthetic Benchmarks) was set to run, and the amount of slowdown it experienced was measured.

### III. EXPERIMENTAL EVALUATION

We ran our experiments on both the NVIDIA TX2 (Fig. 2) and the Xilinx ZU9EG (Fig. 3) for varying levels of interference intensity (THR%). The *Above* region (red) identifies Polybench which always experience a worse slowdown than *READ\_MISS*. The *Crossing* region (yellow) is for the Polybench which are subject to a worse slowdown than *READ\_MISS* only at certain points. Finally, the *Below* region (green) is for Polybench which are never more slowed down than *READ\_MISS*, which is what previous literature actually accounted for. The results definitively prove that read only synthetic benchmarks (*READ\_MISS*) are not:

- 1) The ones which cause the highest amount of interference, as can be seen by *MEMSET*'s highest slowdown being around 12x on the ZU9EG, compared to *READ\_MISS*'s 3x. On the TX2, *MEMCPY* causes the highest amount of interference, instead. It's important

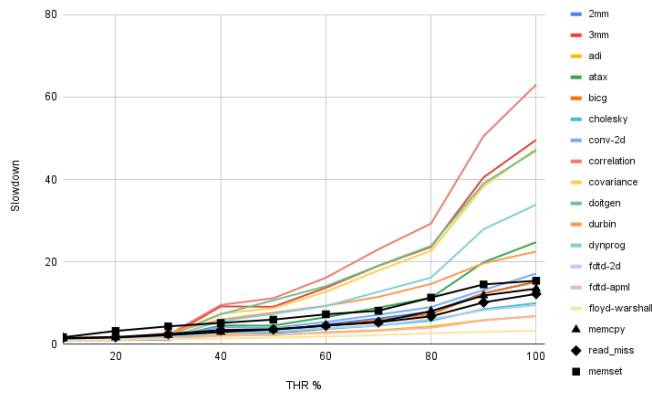


Fig. 4: Slowdown for certain Polybench tasks on the Xilinx ZU9EG, with 3 other CPU cores running MEMSET, and 3 FPGA cores producing RW traffic

to note that there may be even worse interference patterns, that we have not yet found.

2) The programs which are slowed down the most. Depending on the memory access pattern of the interference, different synthetic benchmarks become the worst case. Not only that, but there are certain Polybench which reach even higher degrees of slowdown for all the different interference causing traffic patterns. This, however, is due to Cache events, and not exclusively DRAM interference.

For the ZCU102, the combined effect of these factors makes it so the real worst-case is actually subject to a 12x slowdown factor, instead of the 1.3x which could be observed when using *READ\_MISS* as both the benchmark under test, and the program causing the interference.

#### IV. CONCLUSION

Our research makes it clear that proper worst-case characterization is important when developing memory contention mitigation techniques. While being focused on DRAM interference is important, Cache events must be also accounted for, as they can cause regular tasks to be subject to more interference than the most memory intensive ones (Red/Yellow regions in Fig. 2 and 3). On the Xilinx ZU9EG (Fig. 3), the situation is particularly egregious, as the worst-case for our tests is actually more than  $9\times$  worse than *READ\_MISS*. All of these evaluations have been done with the accelerators turned off, in order to create a valid comparison of the effects of bad worst-case characterization on different HeSoCs. However, we have observed that programs can experience higher degrees of slowdown when the platform-specific hardware is used. For example, on the Xilinx ZU9EG, when the FPGA cores are used to generate traffic, the total slowdown a program can experience can be greater than  $60\times$  (Fig. 4) if the CPU cores are also executing *MEMSET*. Any proposal which fails to account for these kind of scenarios

cannot realistically state to have covered the worst-case slowdown which a program can experience on these kinds of platforms.

#### REFERENCES

- [1] G. Brilli, A. Capotondi, P. Burgio, and A. Marongiu, "Understanding and Mitigating Memory Interference in FPGA-based HeSoCs," in *2022 Design, Automation Test in Europe Conference Exhibition (DATE)*, 2022, pp. 1335–1340.
- [2] A. Bansal, R. Tabish, G. Gracioli, R. Mancuso, R. Pellizzoni, and M. Caccamo, "Evaluating the Memory Subsystem of a Configurable Heterogeneous MPSoC," in *Workshop on Operating Systems Platforms for Embedded Real-Time Applications (OSPERT)*, 07 2018, p. 55.
- [3] K. Manev, A. Vaishnav, and D. Koch, "Unexpected Diversity: Quantitative Memory Analysis for Zynq UltraScale+ Systems," in *2019 International Conference on Field-Programmable Technology (ICFPT)*, 2019, pp. 179–187.
- [4] P. Radojkovi, S. Girbal, A. Grasset, E. Qui Nones, S. Yehia, and F. J. Cazorla, "On the Evaluation of the Impact of Shared Resources in Multithreaded COTS Processors in Time-Critical Environments."
- [5] J. Nowotsch and M. Paulitsch, "Leveraging Multi-core Computing Architectures in Avionics," in *2012 Ninth European Dependable Computing Conference*, 2012, pp. 132–143.
- [6] N. Capodieci, R. Cavicchioli, I. S. Olmedo, M. Solieri, and M. Bertogna, "Contending memory in heterogeneous SoCs: Evolution in NVIDIA Tegra embedded platforms," in *2020 IEEE 26th International Conference on Embedded and Real-Time Computing Systems and Applications (RTCSA)*, 2020, pp. 1–10.
- [7] M. Bechtel and H. Yun, "Memory-Aware Denial-of-Service Attacks on Shared Cache in Multicore Real-Time Systems," *IEEE Transactions on Computers*, vol. 71, no. 9, pp. 2351–2357, 2022.
- [8] R. Pellizzoni, E. Betti, S. Bak, G. Yao, J. Criswell, M. Caccamo, and R. Kegley, "A Predictable Execution Model for COTS-Based Embedded Systems," in *2011 17th IEEE Real-Time and Embedded Technology and Applications Symposium*, 2011, pp. 269–279.
- [9] H. Yun, G. Yao, R. Pellizzoni, M. Caccamo, and L. Sha, "MemGuard: Memory bandwidth reservation system for efficient performance isolation in multi-core platforms," in *2013 IEEE 19th Real-Time and Embedded Technology and Applications Symposium (RTAS)*, 2013, pp. 55–64.
- [10] R. Cavicchioli, N. Capodieci, M. Solieri, M. Bertogna, P. Valente, and A. Marongiu, *Evaluating Controlled Memory Request Injection to Counter PREM Memory Underutilization*, 11 2020, pp. 85–105.
- [11] M. Mattheeuws, B. Forsberg, A. Kurth, and L. Benini, "Analyzing Memory Interference of FPGA Accelerators on Multicore Hosts in Heterogeneous Reconfigurable SoCs," in *2021 Design, Automation Test in Europe Conference Exhibition (DATE)*, 2021, pp. 1152–1155.
- [12] G. Brilli, R. Cavicchioli, M. Solieri, P. Valente, and A. Marongiu, "Evaluating Controlled Memory Request Injection for Efficient Bandwidth Utilization and Predictable Execution in Heterogeneous SoCs," *ACM Trans. Embed. Comput. Syst.*, sep 2022, just Accepted. [Online]. Available: <https://doi.org/10.1145/3548773>