



UNIMORE
UNIVERSITÀ DEGLI STUDI DI
MODENA E REGGIO EMILIA



The importance of memory in the next generation of real-time systems

Paolo Burgio
paolo.burgio@unimore.it



Future embedded systems

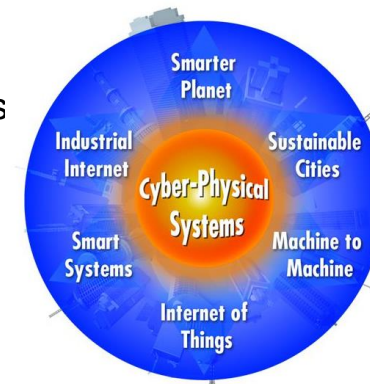


Industry 4.0

The four horsemen

1. **Heavy workloads**
 - Sensor-fusion and image-processing
2. **Reduced power consumption**
 - Smaller batteries and renewable power sources
3. **Quickly interact with the environment**
 - Prompt elaboration of sensor data
4. **Run highest criticality workloads**
 - Replacing safety-critical human activities

Internet-of-Things



Cyber-physical systems

Artificial intelligence



Autonomous driving



Health and medicine



Real-Time multi-core systems?

Multi- and many-core platforms are the solution for 1-2(-3)

- ✓ Climbing "the power wall"
- ✓ High Performance @ poor Watts



Real-Time system: produce result in a guaranteed/bounded amount of time

- ✓ **By construction**
- ✓ Application fields: automotive, avionics, industry, medical...

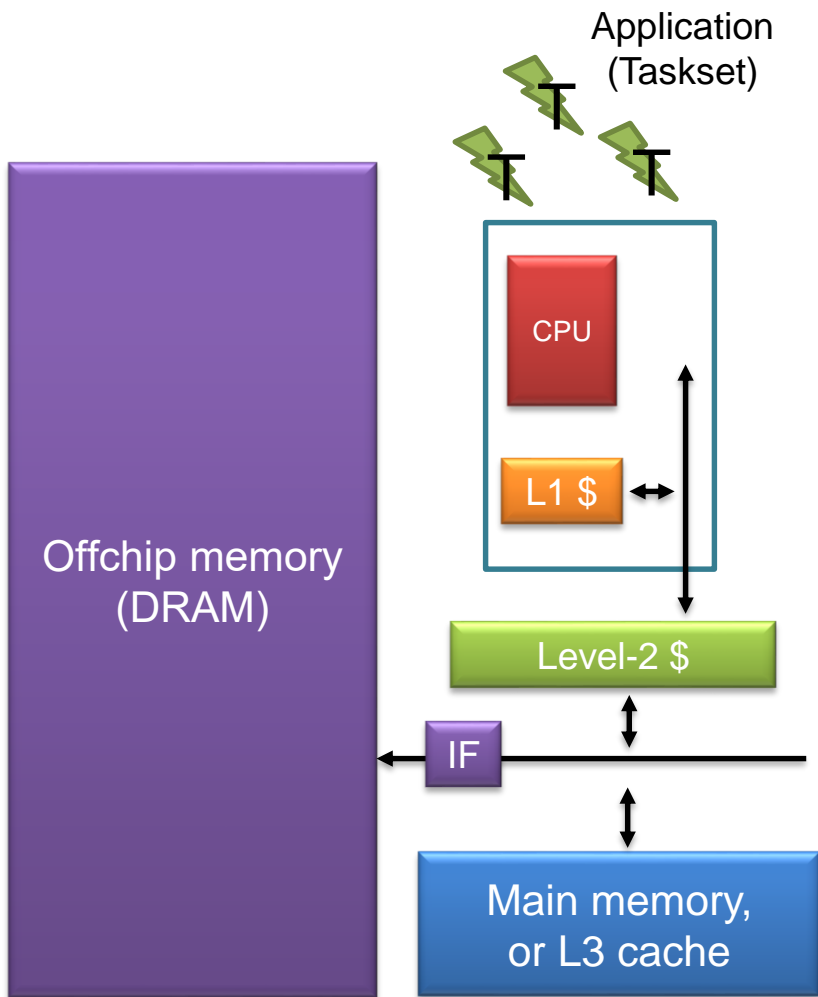
The keyword: **predictability**

- ✓ Provide the correct result....**when** expected
- ✓ The system must be simple to **analyze**





Real-Time systems – traditional approach

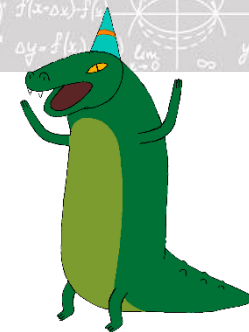


Single-core, multiple tasks/applications

1. Analyze the system (HW/SW)
2. Derive a (mathematical?) model
3. Do some magic mathematics...



...guaranteed timing bounds!



Optimal sharing of the core between task

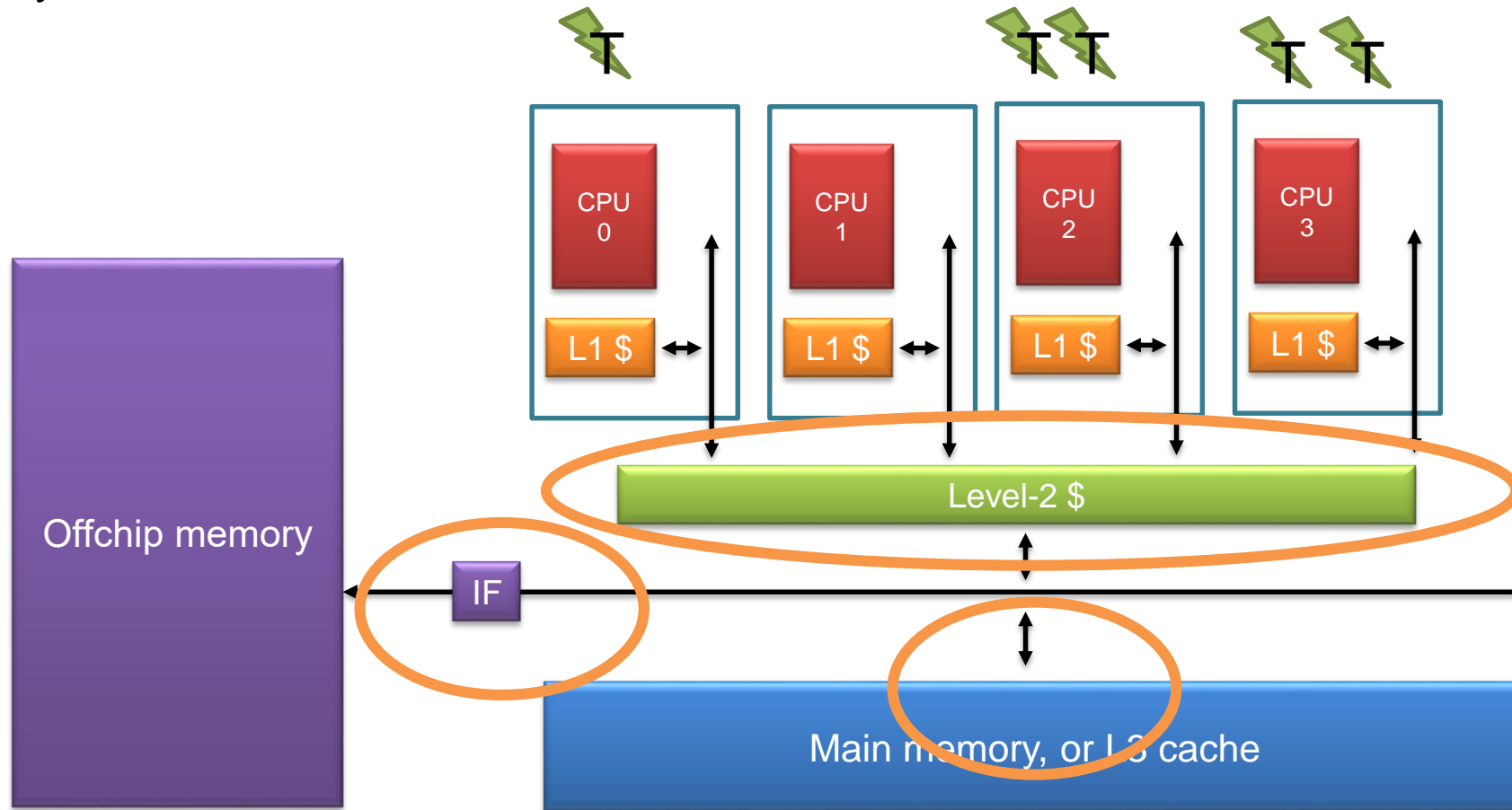
- ✓ ..and guaranteed **by construction**
- ✓ *Scheduling* (also, mapping)



Multi-core systems

Architectural bottlenecks

- ✓ Shared memory banks
- ✓ Caches (\$)
- ✓ I/Os





It's (mainly) a memory issue!

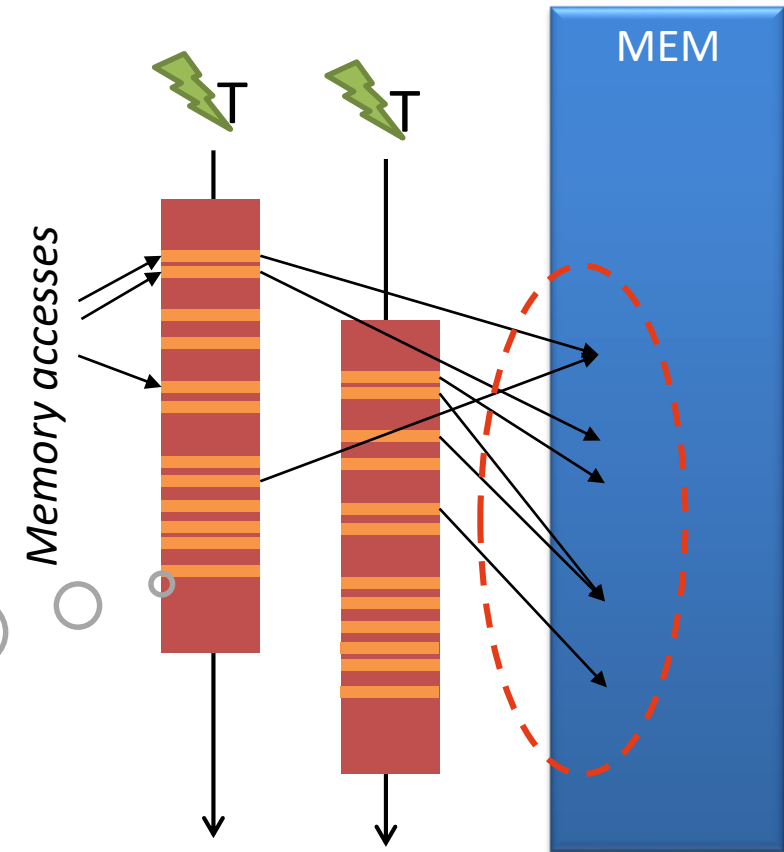
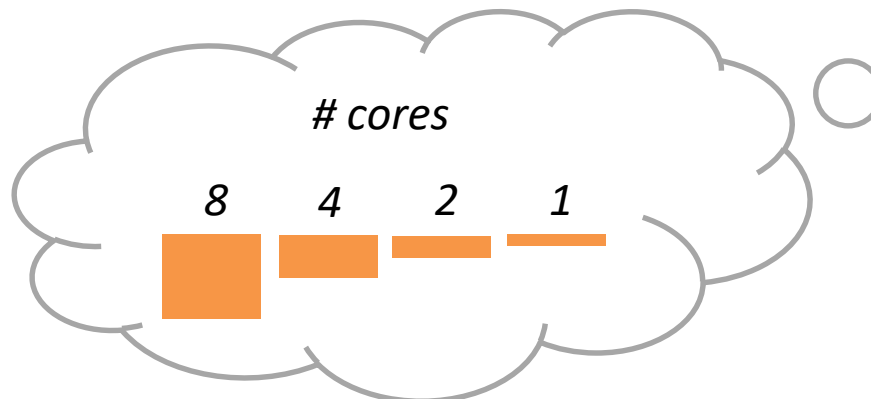
Beyond traditional techniques

1. More parameters

- Shared resources (e.g., memory, SSDs, IOs, caches..)
- The complexity of analysis grows exponentially w/number of cores

2. Mem accesses: instead of thin lines, big bars

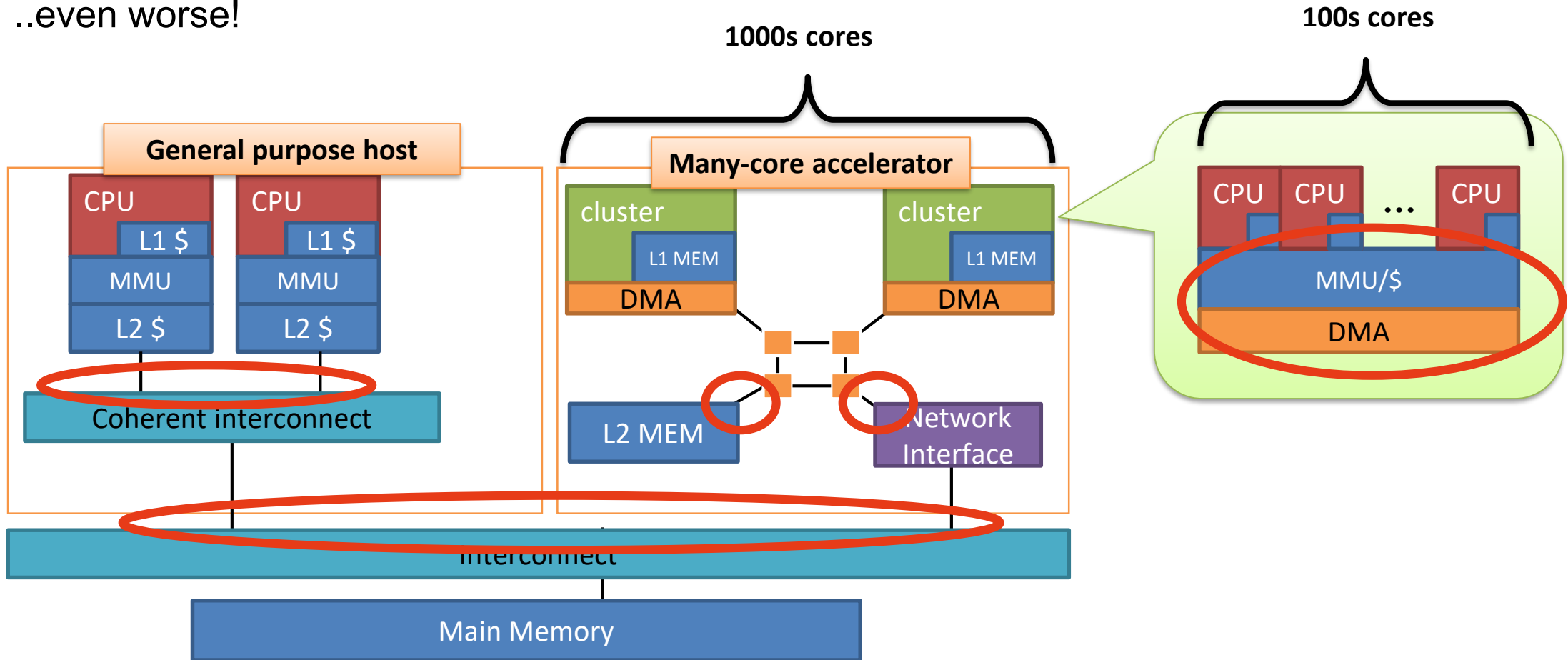
- The mostly accessed resource in the system
- Traditional techniques are too conservative (bounds too high)





Many-core systems

- ✓ Thousands cores arranged in CLUSTERS
- ✓ Host-accelerator architecture (e.g., GP-GPUs)
- ✓ ..even worse!





Knowledge of the platform is power

Two motivating examples

- ✓ Both from real systems

- 1. Many-core accelerator-based platforms
 - Quad-/Octa-core as host
 - Integrated GPU – iGPU of FPGA
 - Powerful enough to run neural networks

- 2. Reference industrial system
 - Multi-core ARM
 - Multi-OS (embedded Linux + Win for UI)
 - Hypervisor-based





Testbed #1: "automotive" platforms

Qualitatively analyze and characterize the conflicts due to parallel accesses to main memory by both CPU cores and iGPU

1. NVIDIA Tegra K1 w/Kepler GPU
2. NVIDIA Tegra X1 w/Maxwell GPU
3. NVIDIA Tegra X2 w/Parker GPU – automotive-grade
4. Intel i7-6700 w/intel GPU
5. Xilinx Zynq Ultrascale multi-core + FPGA (+GPU)

HERCULES

Roberto Cavicchioli, Nicola Capodieci and Marko Bertogna, *"Memory Interference Characterization between CPU cores and integrated GPUs in Mixed-Criticality Platforms"*, 22nd IEEE International Conference on Emerging Technologies And Factory Automation

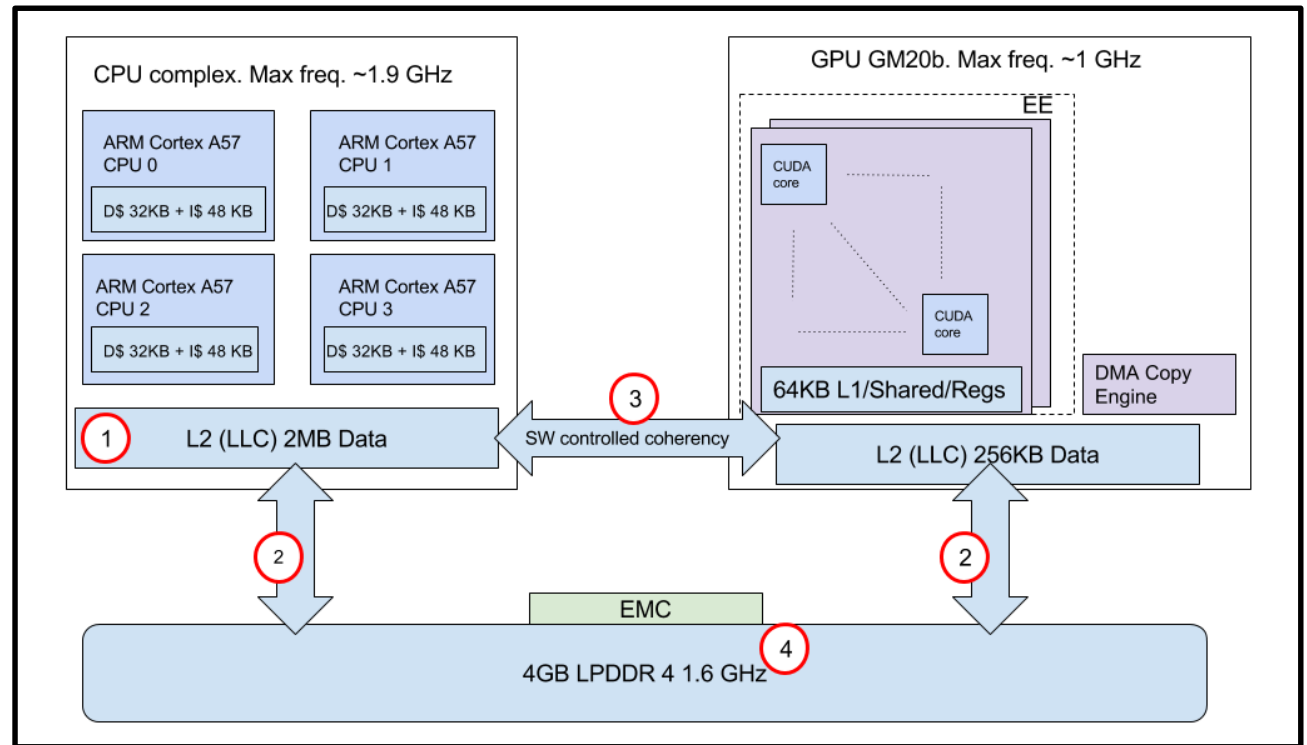


NVIDIA Tegra K2

- ✓ Shared memory between CPU/GPU complex
 - "Unified Virtual Memory"
 - Unlike traditional "discrete" GPU systems

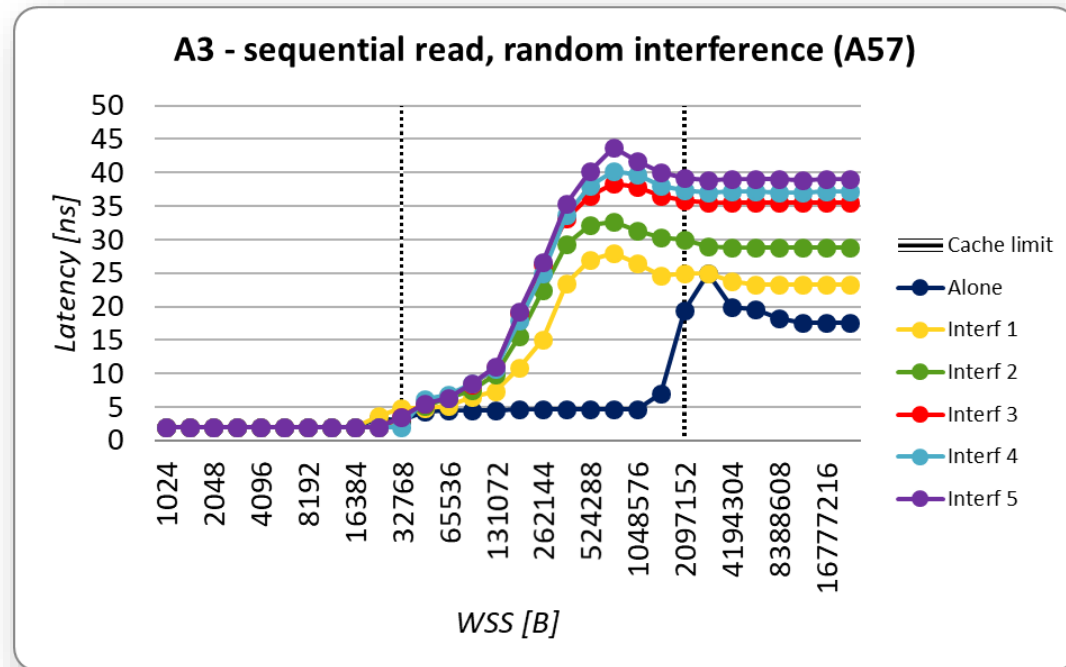
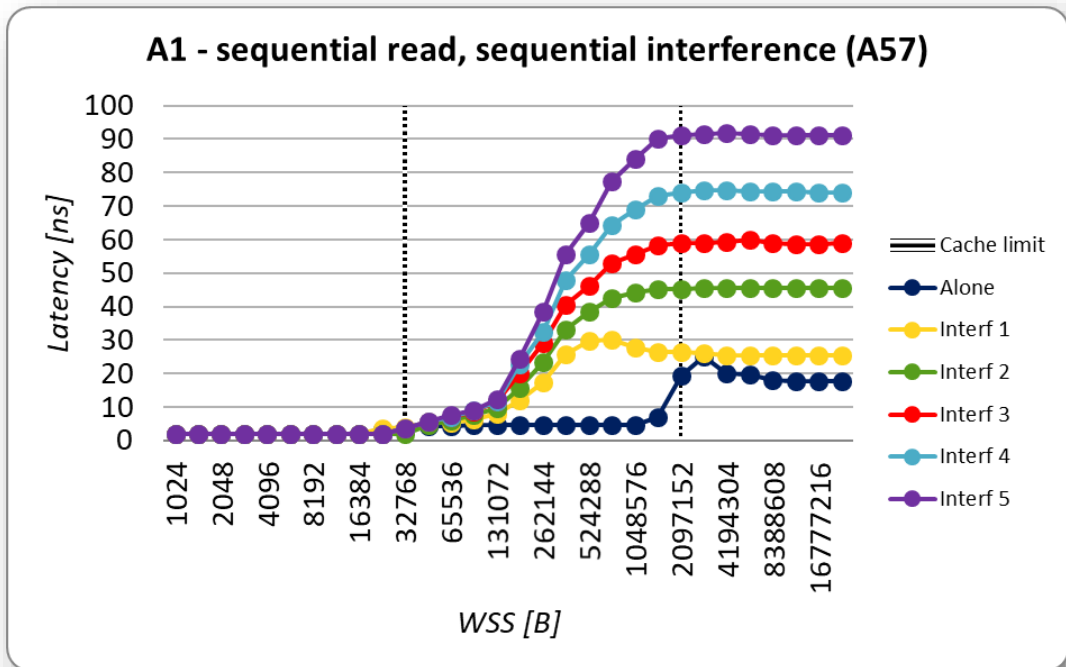
Notable contention points

1





Test 'A' - Tegra X2 – A57



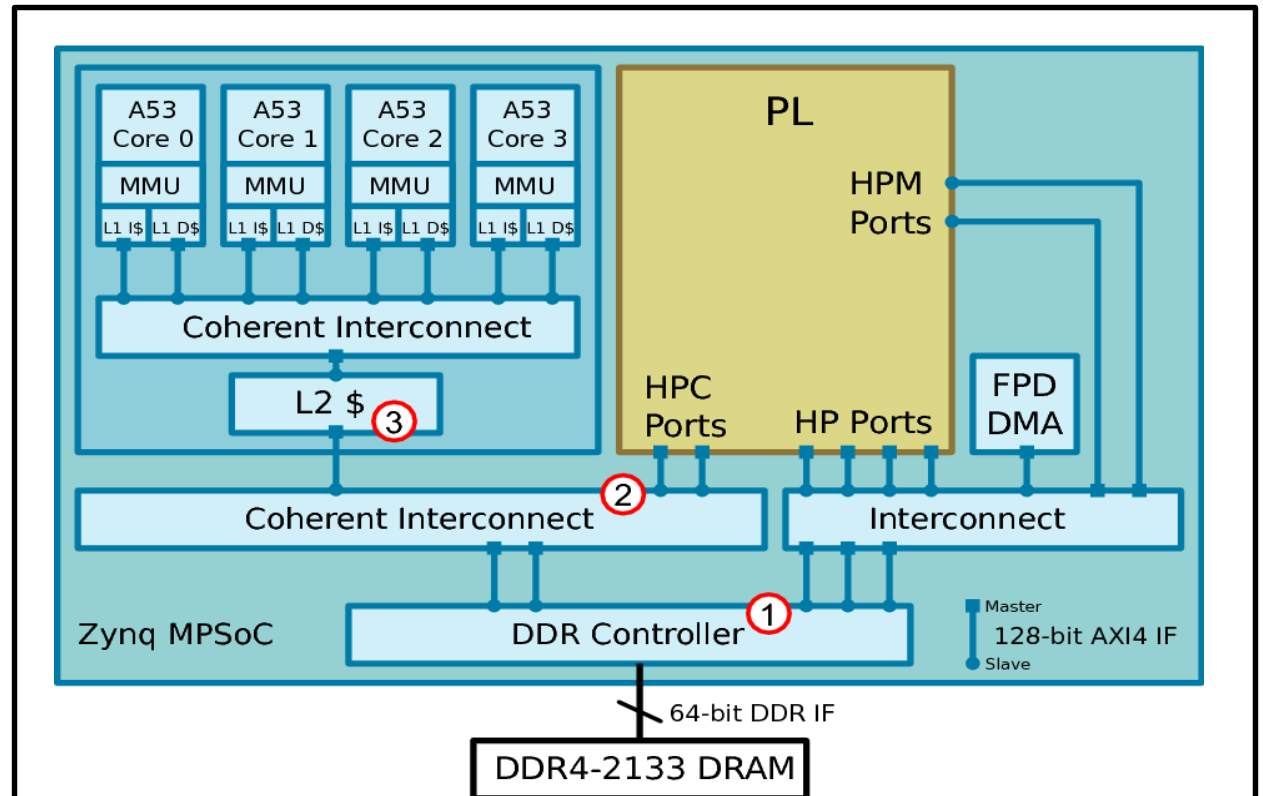


Xilinx Zynq Ultrascale

- ✓ Last-generation FPGA-based heterogeneous SoC
 - FPGA = (re-)programmability
- ✓ ARM A53 Quad-core as host "PS"
- ✓ FPGA as accelerator "PL"

Notable contention points

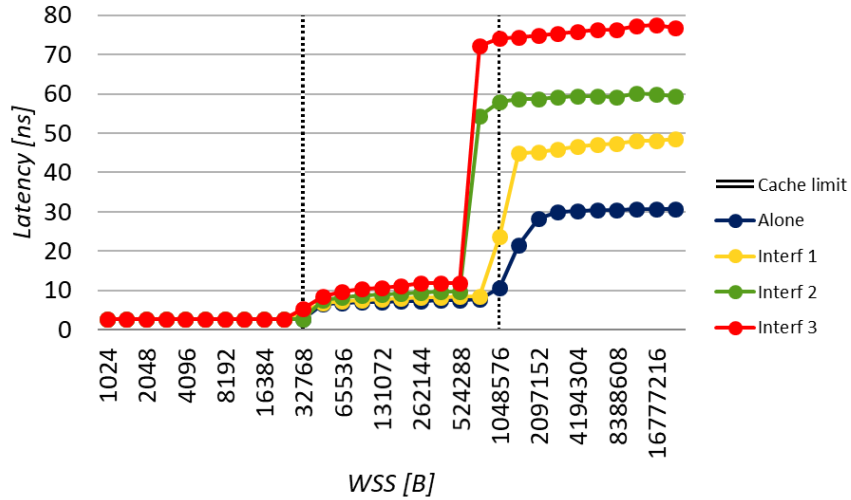
1



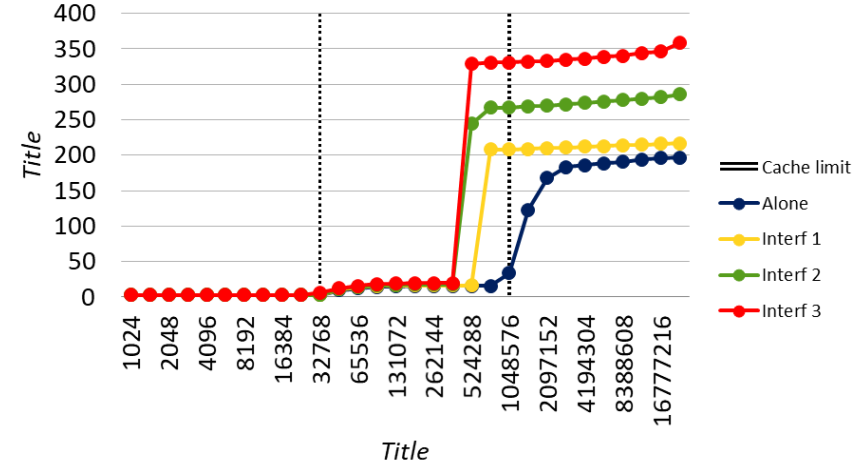


Test 'A' - Xilinx Zynq

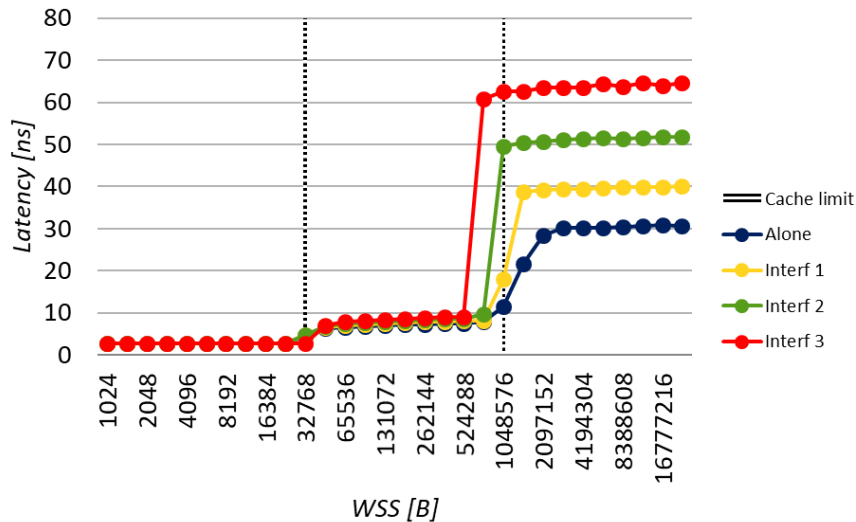
A1 - Sequential read, sequential interference



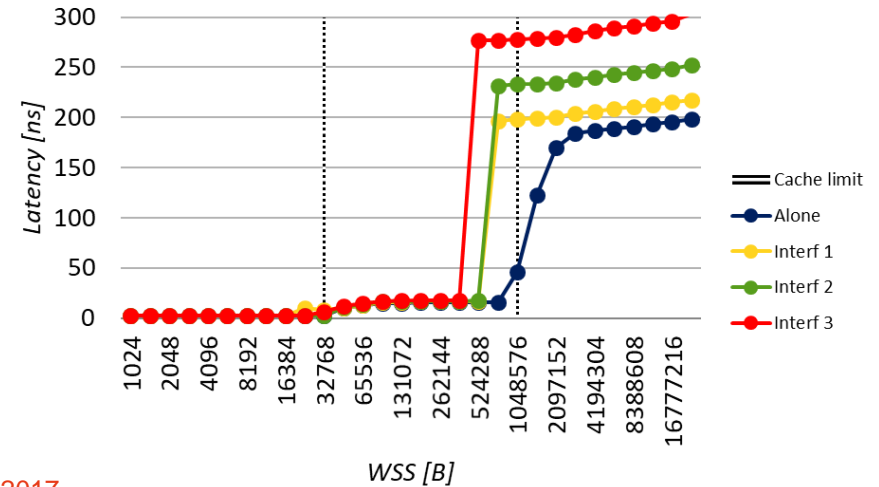
A2 - Random read, sequential interference



A3 - Sequential read, Random interference

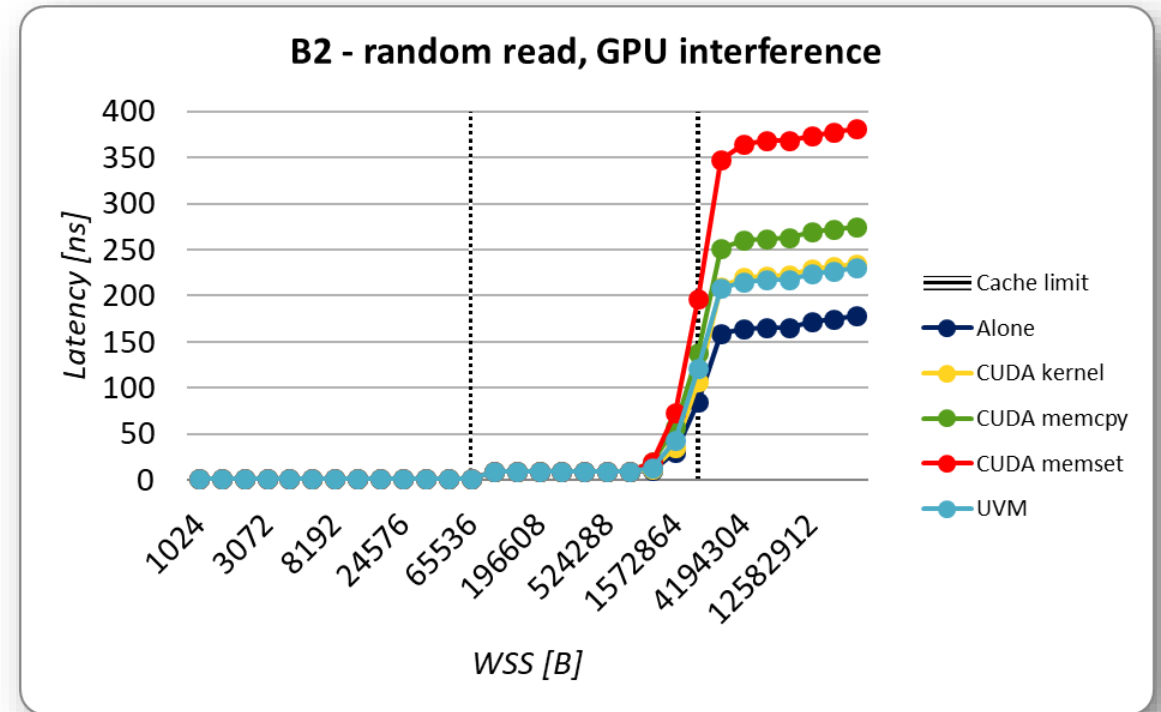
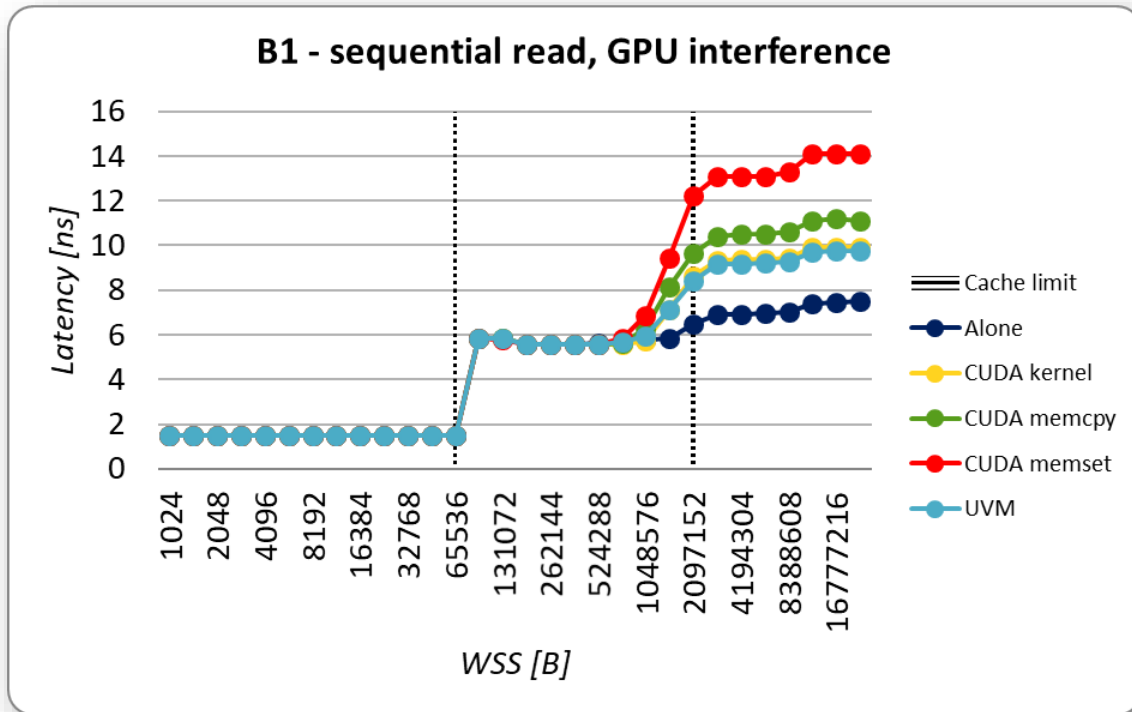


A4 - Random read, Random interference



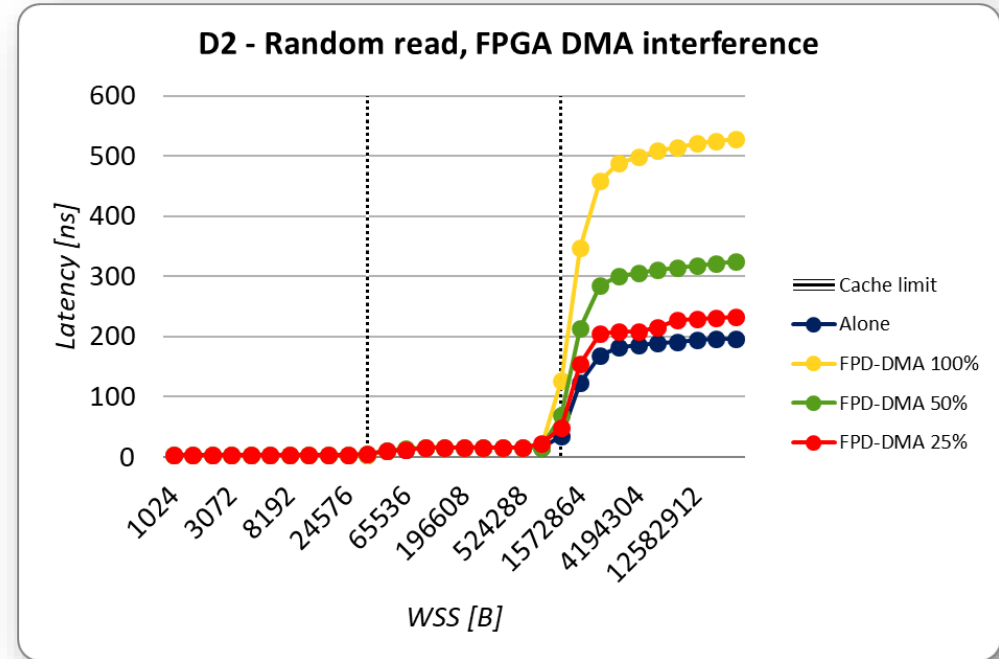
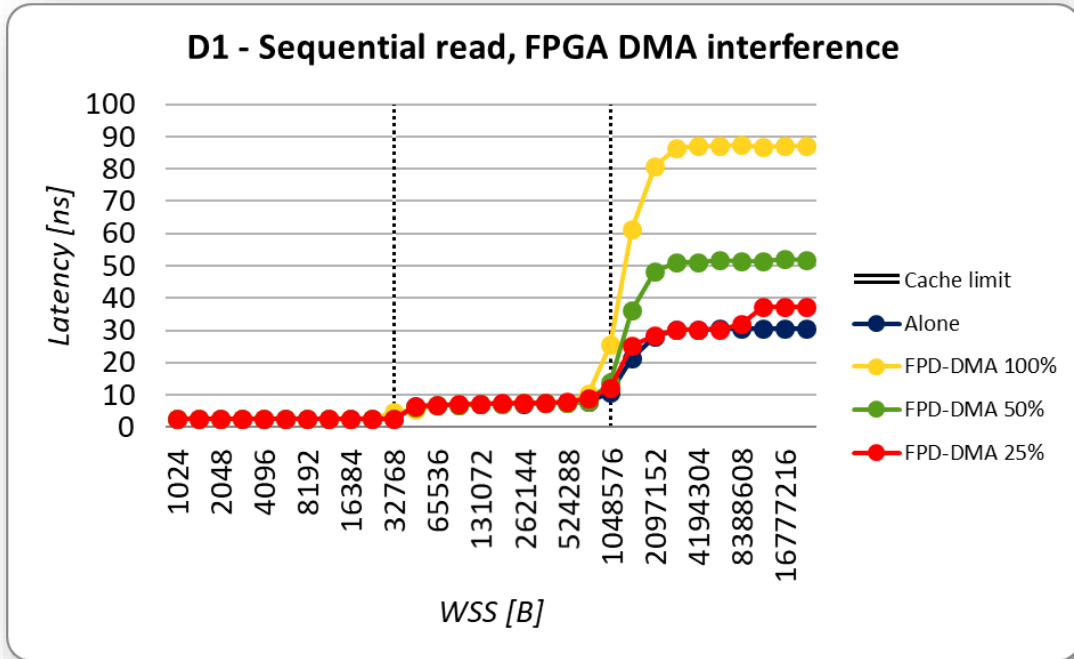


Test 'B' - Tegra X2





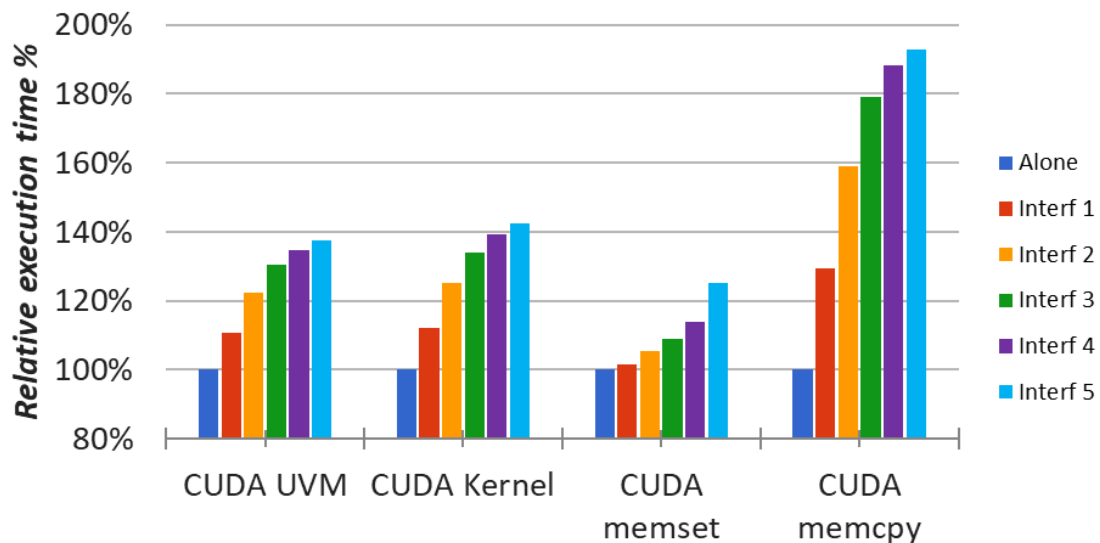
Test 'B' - Xilinx Ultrascale



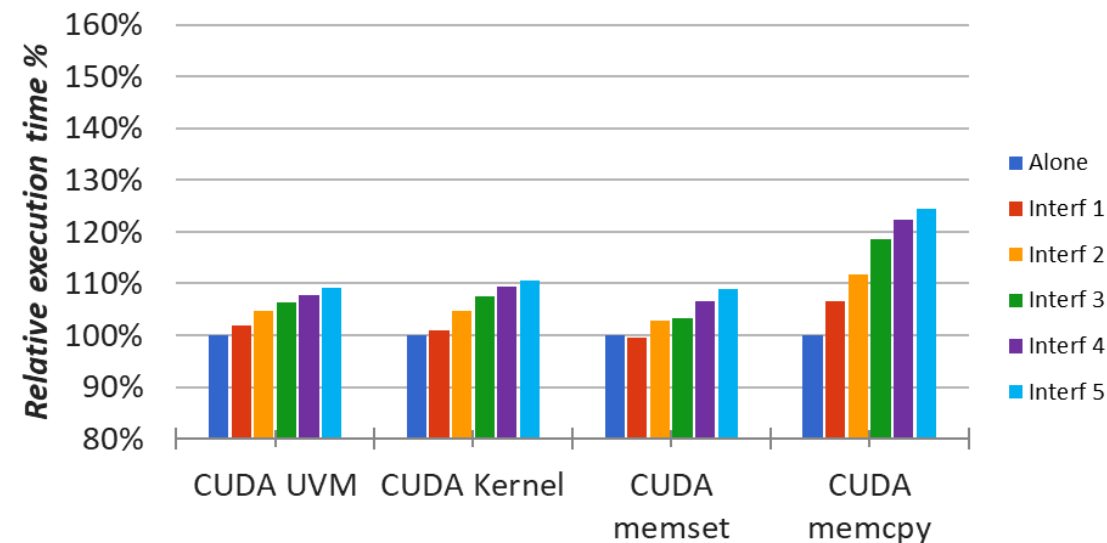


Test 'C' - Tegra X2 – A57

C1 - GPU (GP10b) with CPU interf. (A57, sequential)



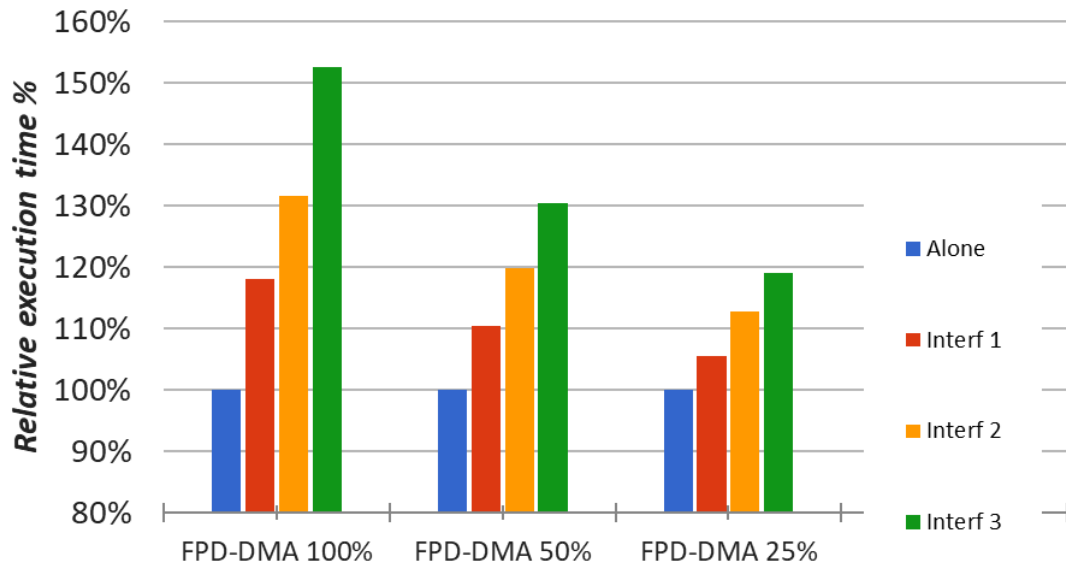
C2 - GPU (GP10b) with CPU interf. (A57, random)



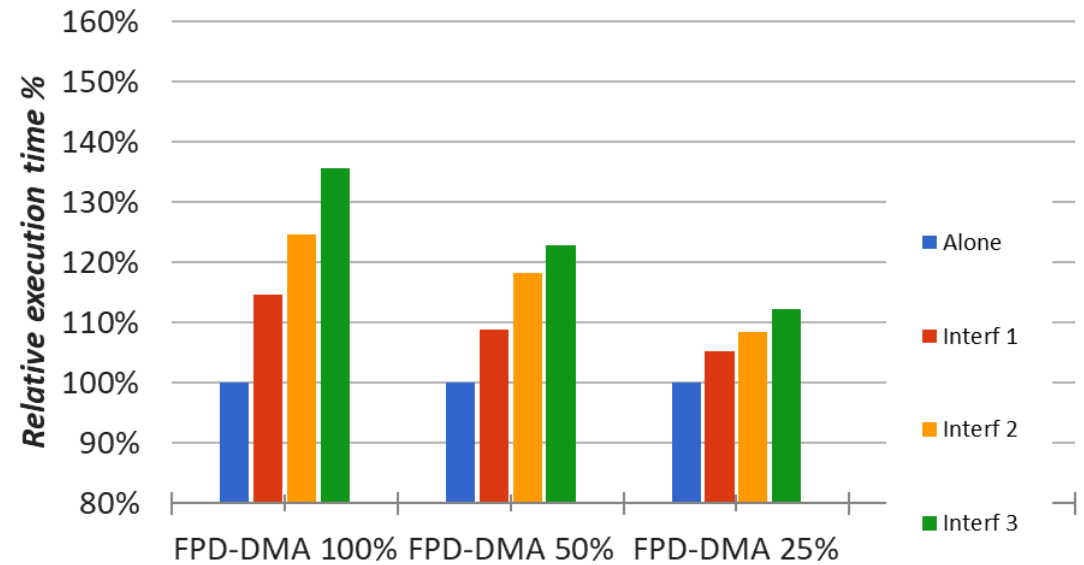


Test 'C' - Xilinx Ultrascale

E1: FPGA DMA Activity with sequential CPU interference



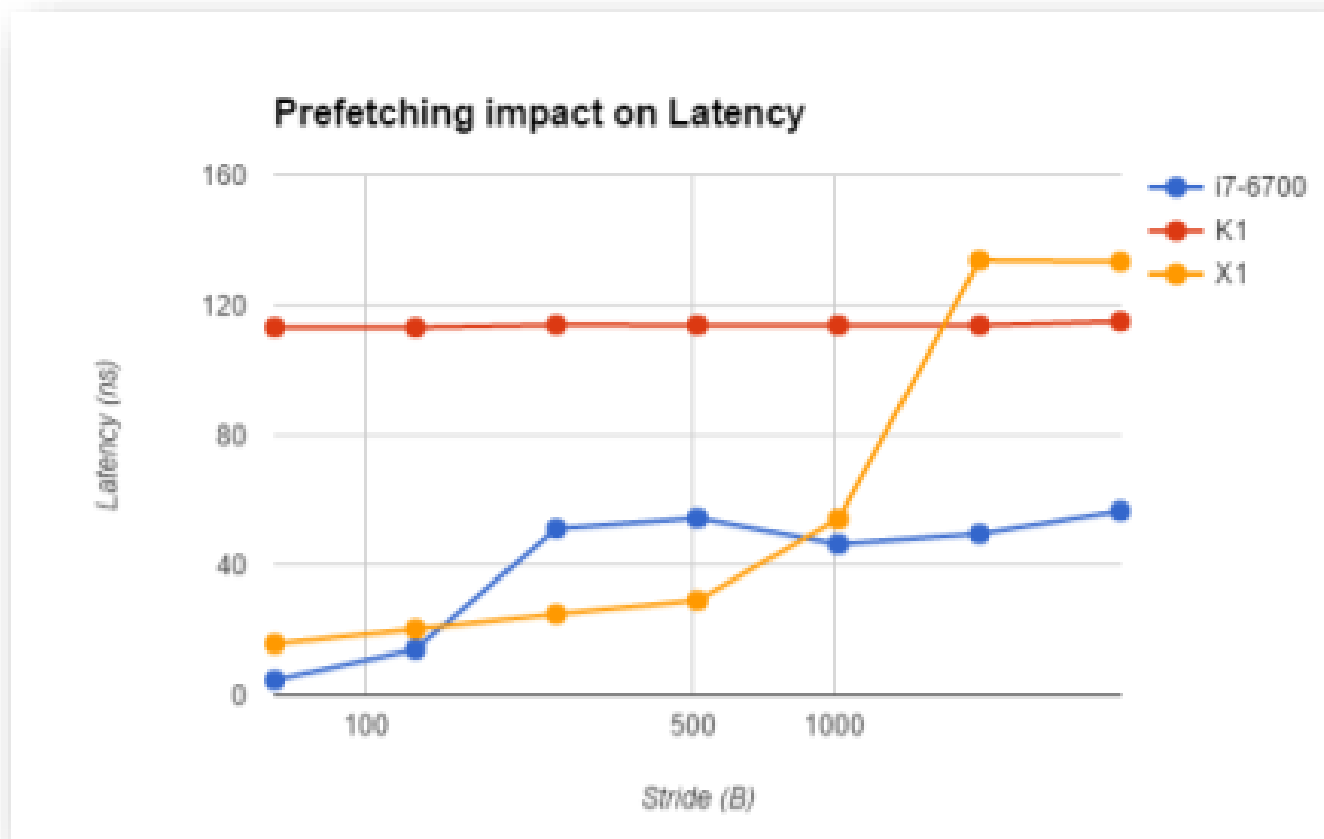
E2: FPGA DMA Activity with random CPU interference





Prefetching

- ✓ Interfere with prefetching mechanism
- ✓ Interfering cores read at increasing strided addresses

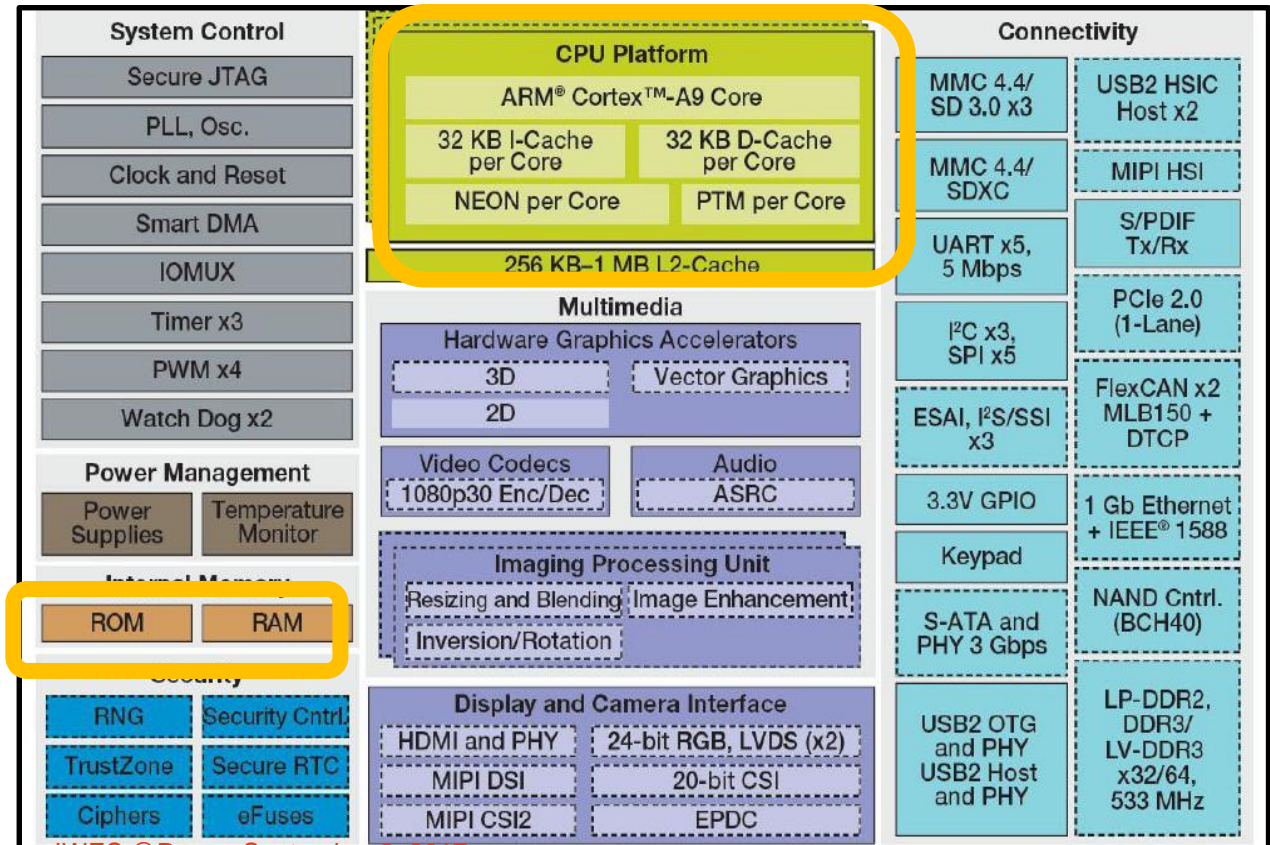




Testbed #2: industrial platform



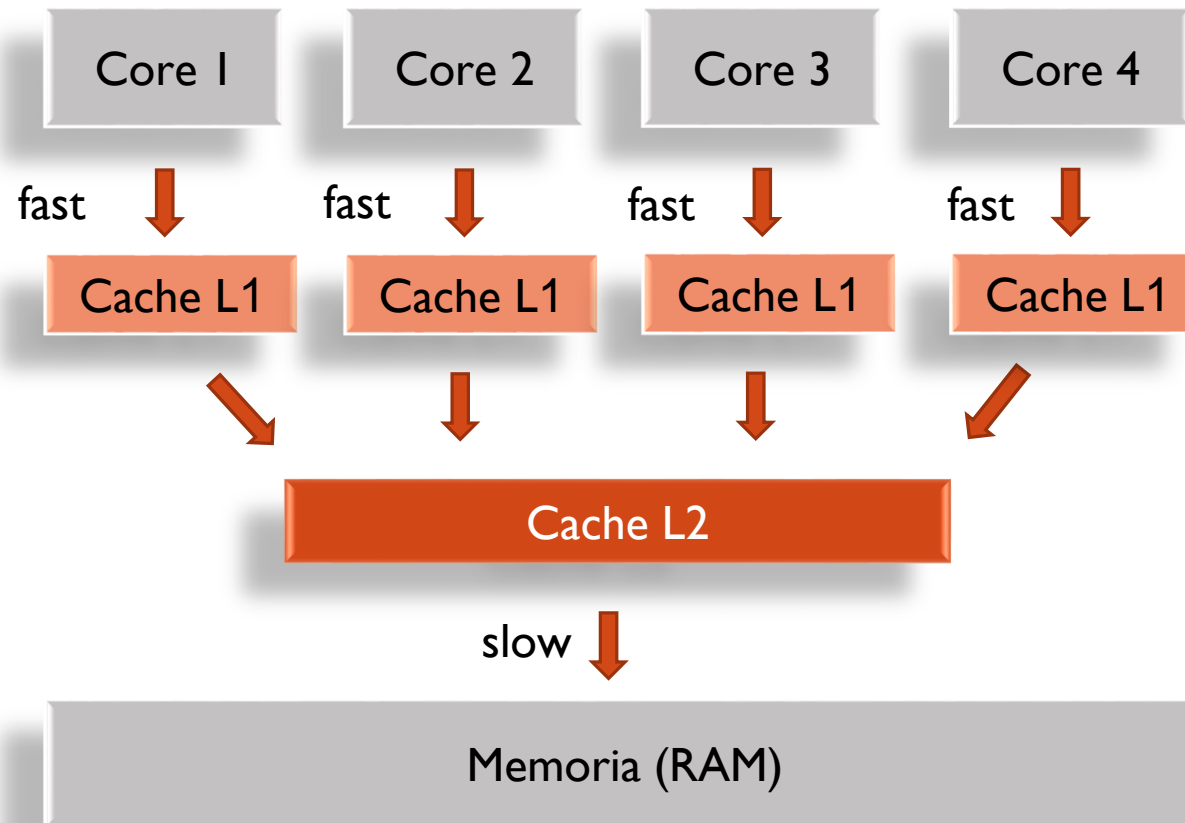
- ✓ NXP iMX6 from Egicon
 - Components for F1 teams, industrial telescopic arms
 - Credits to Francesco Bellei





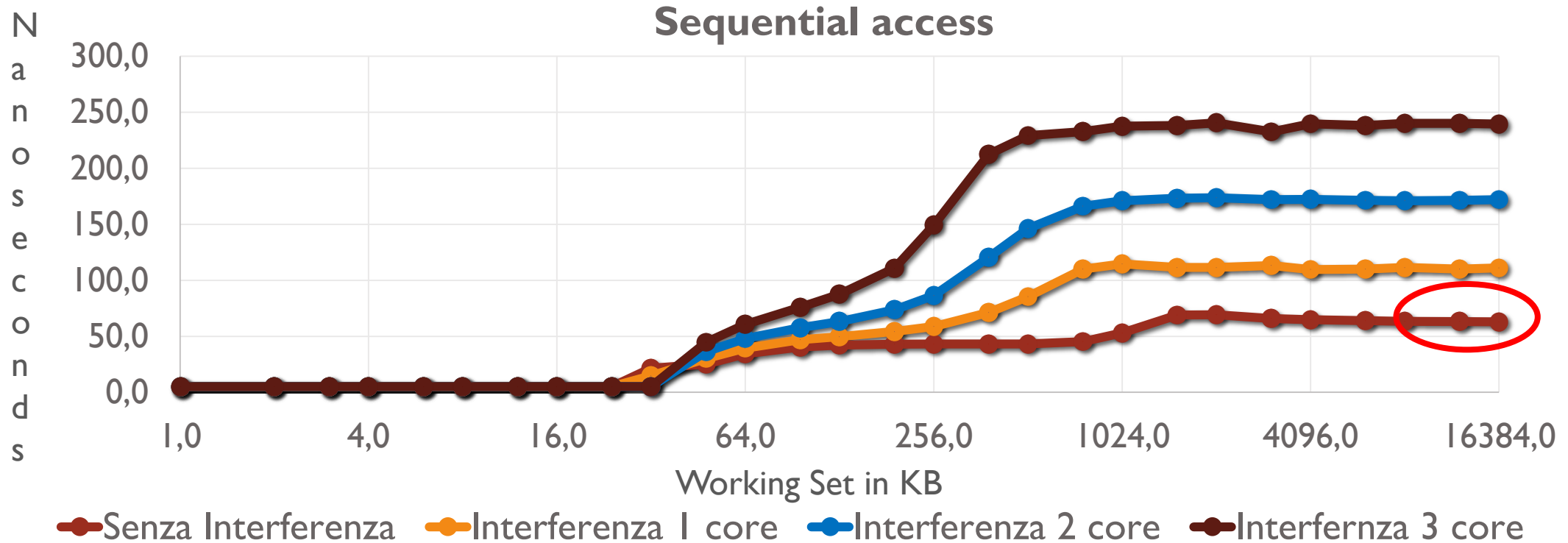
iMX6 mem hierarchy

✓ More "traditional"





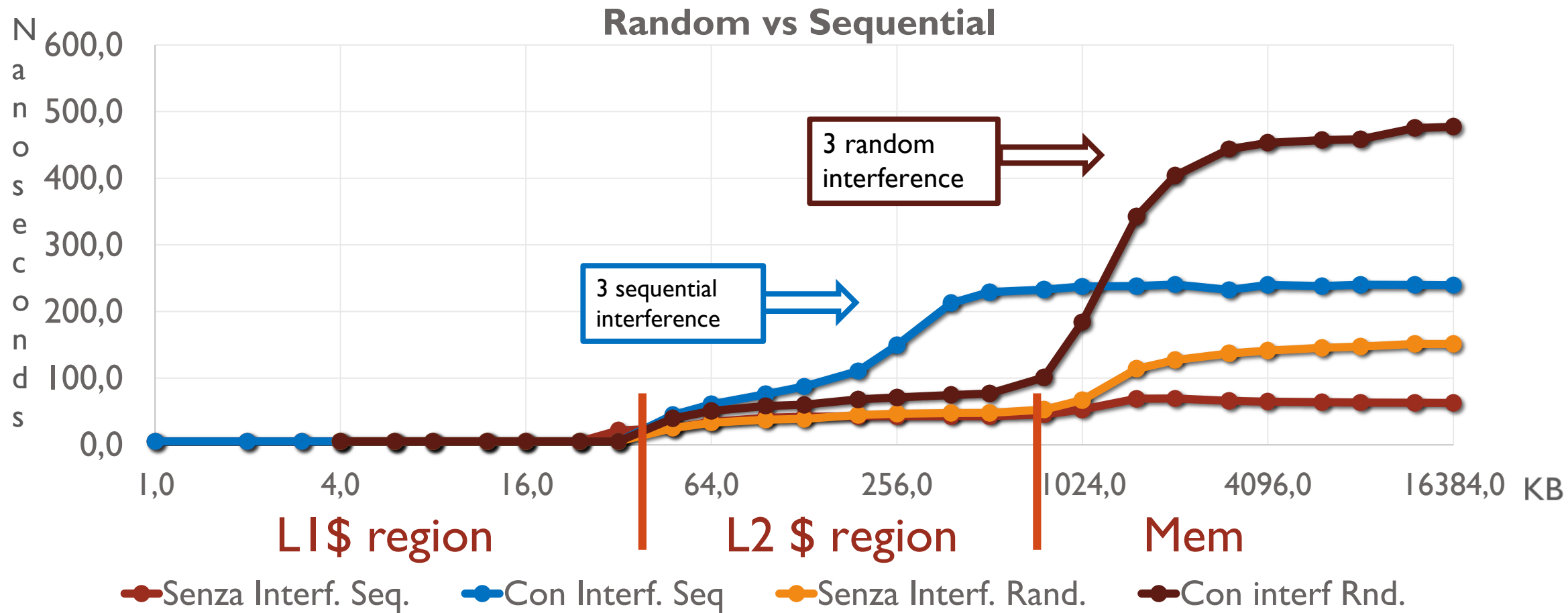
Memory latency - sequential (ns)



$$\boxed{\text{Max BW}} = \frac{\boxed{\text{Cache Line Size}}}{\boxed{\text{Lat}}} = \frac{\boxed{32 \text{ byte}}}{\boxed{62,8 \text{ ns}}} \approx \boxed{0,5 \text{ GB/s}}$$



Memory interference impact





What do we do with this knowledge?

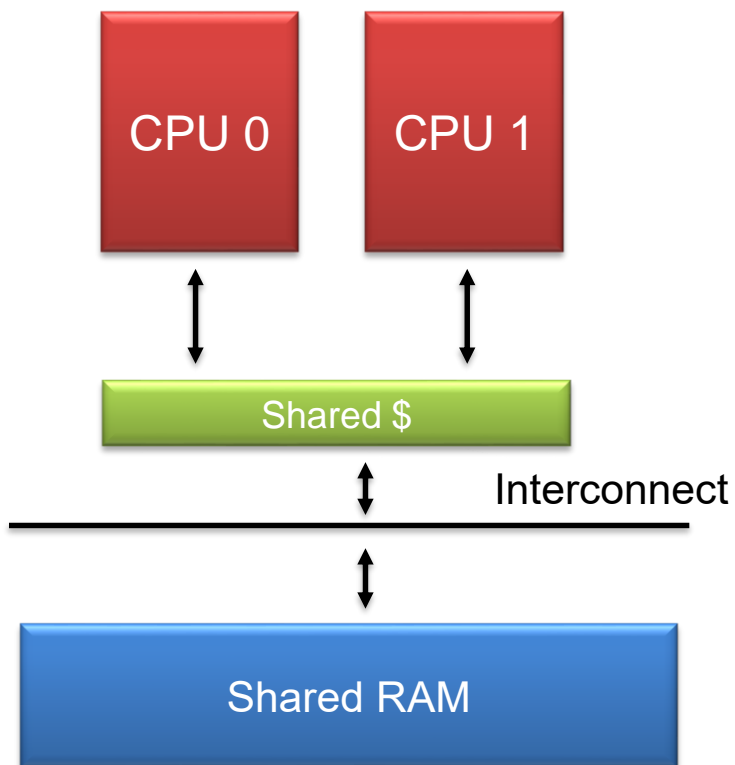


Single-core equivalence

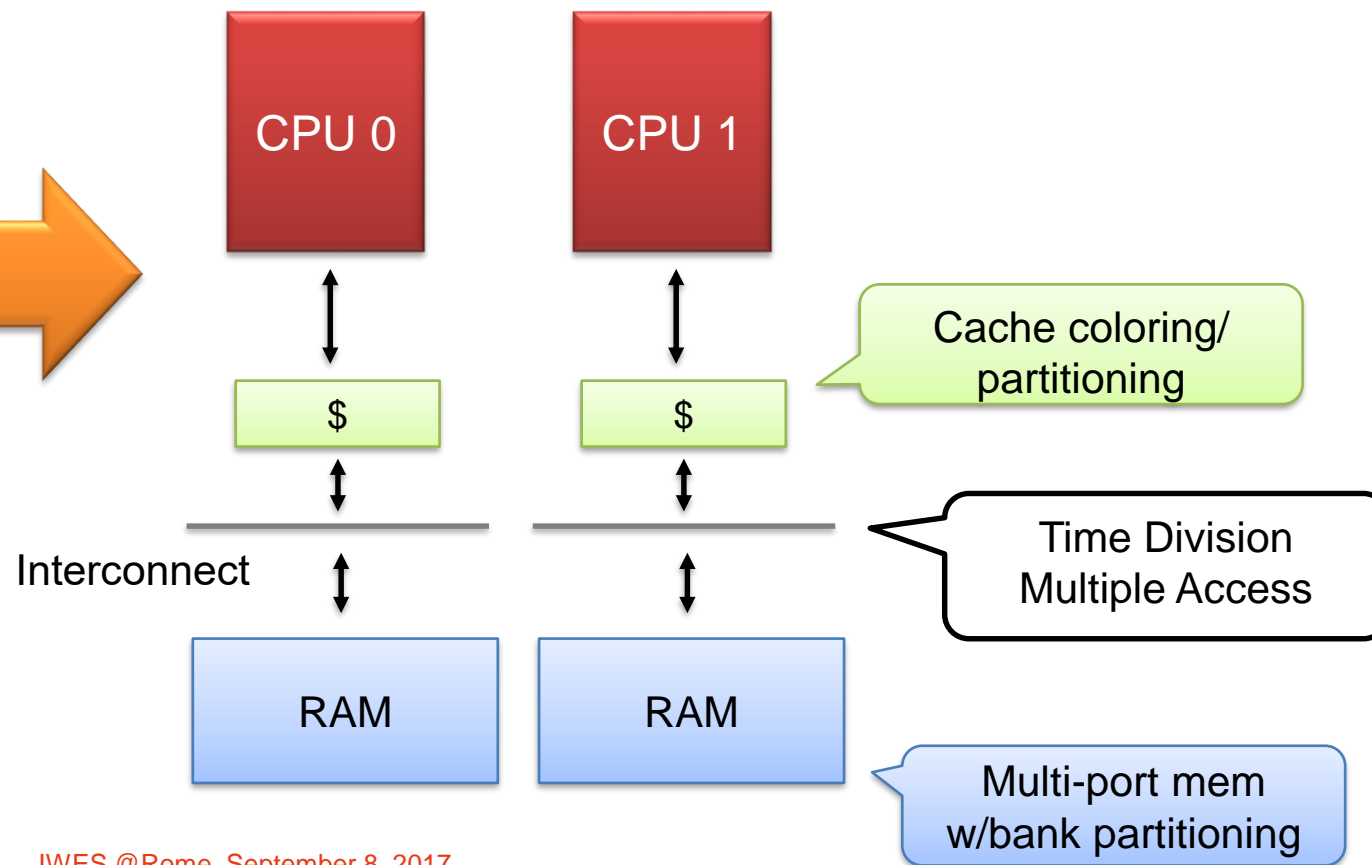


- ✓ A set of techniques to turn the **view of the system that software has..**

From this...

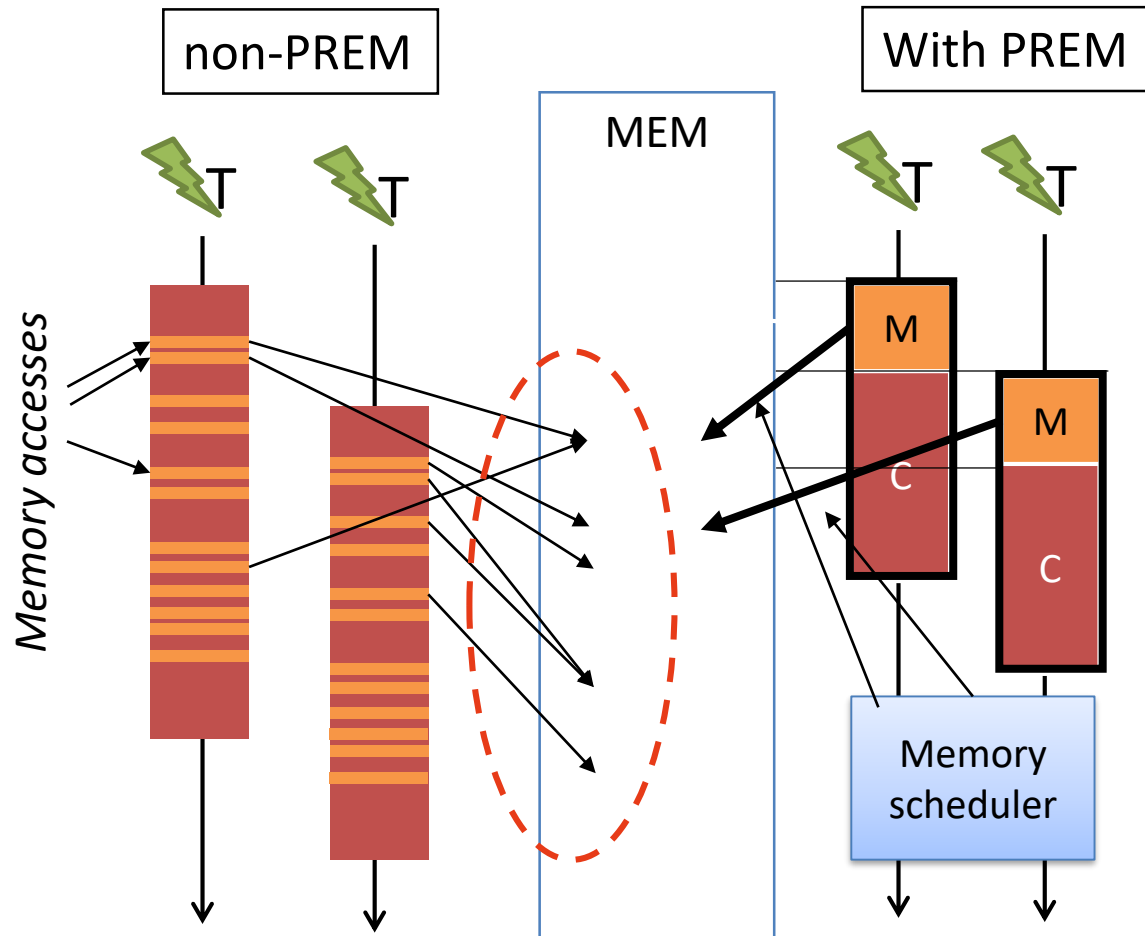


...into this





PREM - PRedictable Execution Models



- ✓ Group memory access at the beginning of every software task
- ✓ Co-schedule memory accesses and tasks-to-cores
- ✓ Greatly reduces the complexity of the scheduling problem

...and increases performance

Up to 4x predictable performance on a many-core platform

2015 paper @ RTEST

# Cores/threads	1	2	4	8
No-PREM – Worst (Analytical)	0.026	0.047	0.088	0.170
PREM – Worst (Analytical)	0.010	0.014	0.022	0.038
Speedup	2.6×	3.4×	4.0×	4.5×



Thank you!

Paolo Burgio
paolo.burgio@unimore.it





Backup



Test case A – intra-CPU interference

1. One observed core reads **sequentially** within a variable sized working set, while other cores are interfering **sequentially**
2. One observed core reads **randomly** within a variable sized working set, while other cores are interfering **sequentially**
3. One observed core reads **sequentially** within a variable sized working set, while other cores are interfering **randomly**
4. One observed core reads **randomly** within a variable sized working set, while other cores are interfering **randomly**

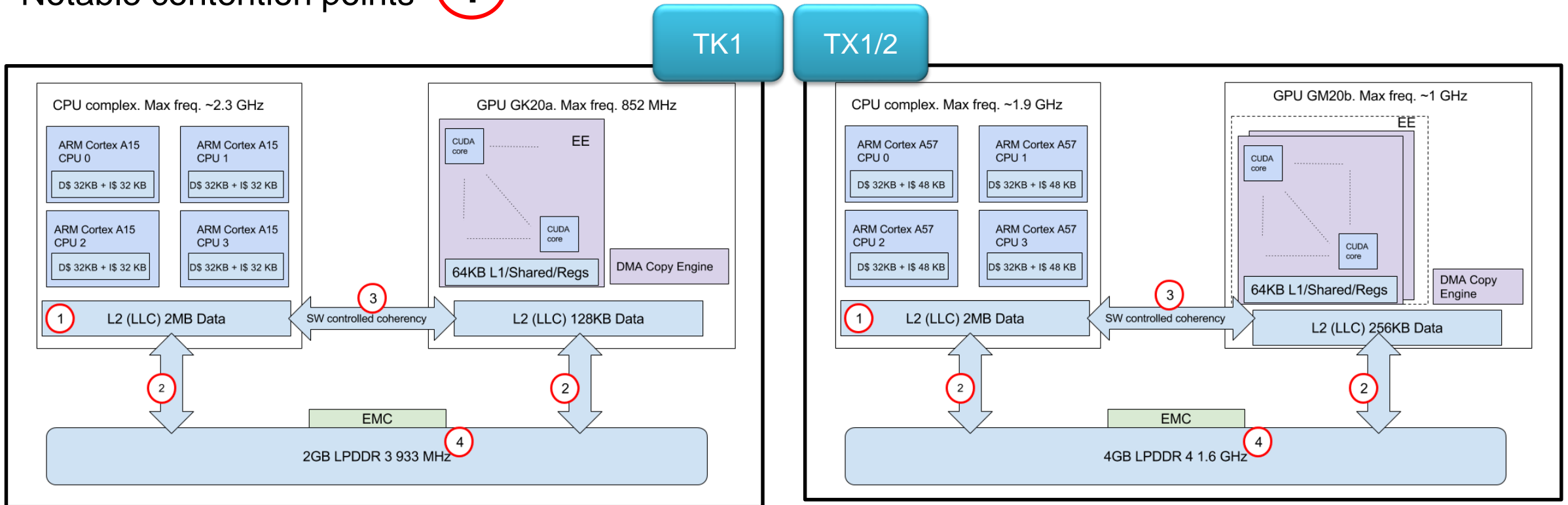


NVIDIA Tegra family

- ✓ Shared memory between CPU/GPU complex
 - "Unified Virtual Memory"

Notable contention points

1



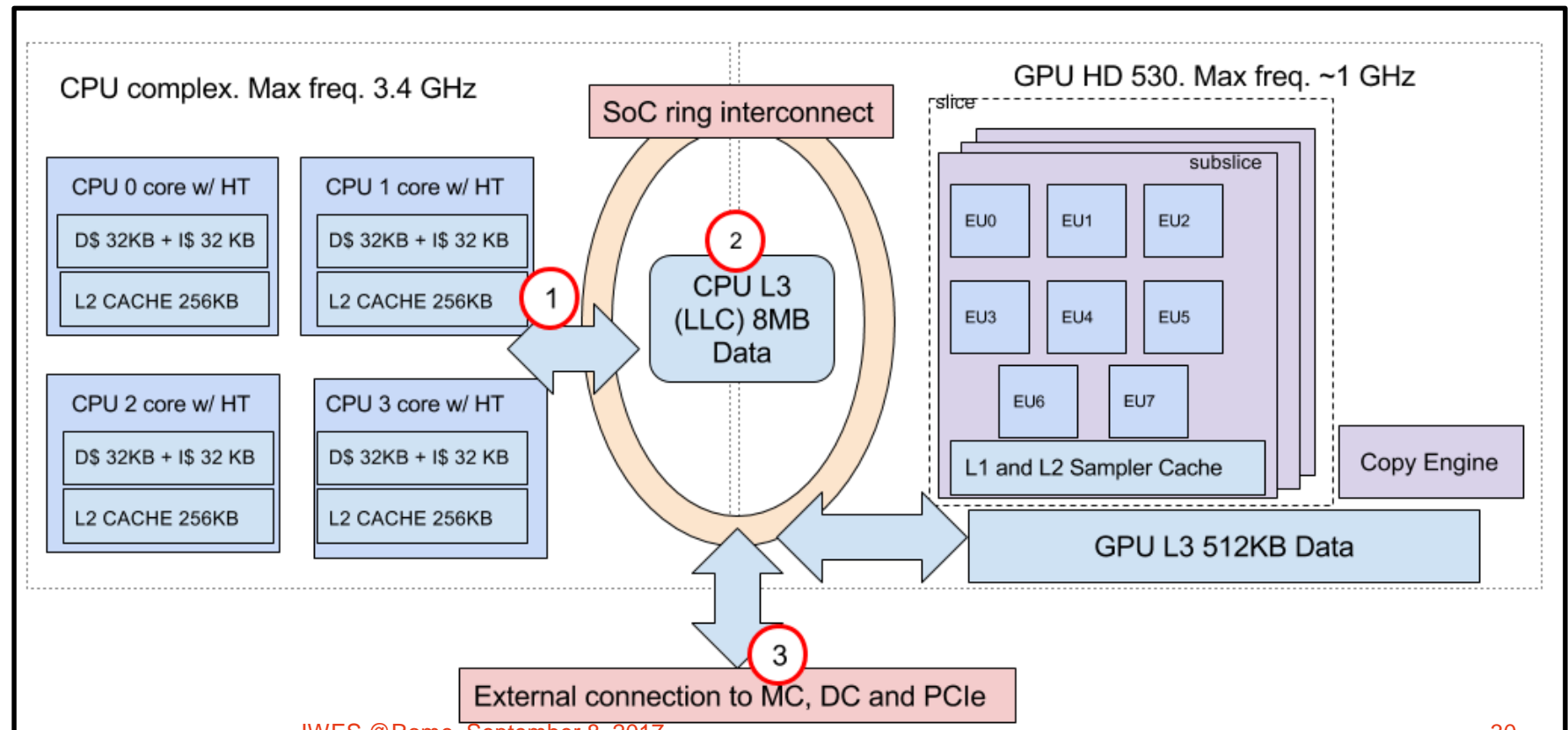


Intel i7-6700 Skylake

- ✓ x86_64 powerful host + iGPU
 - Sharing L3\$, External DRAM...

Notable contention points

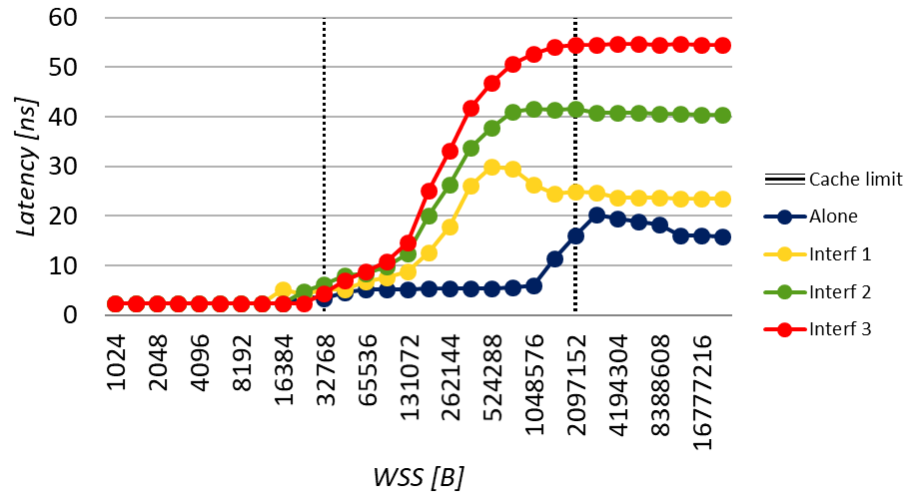
1



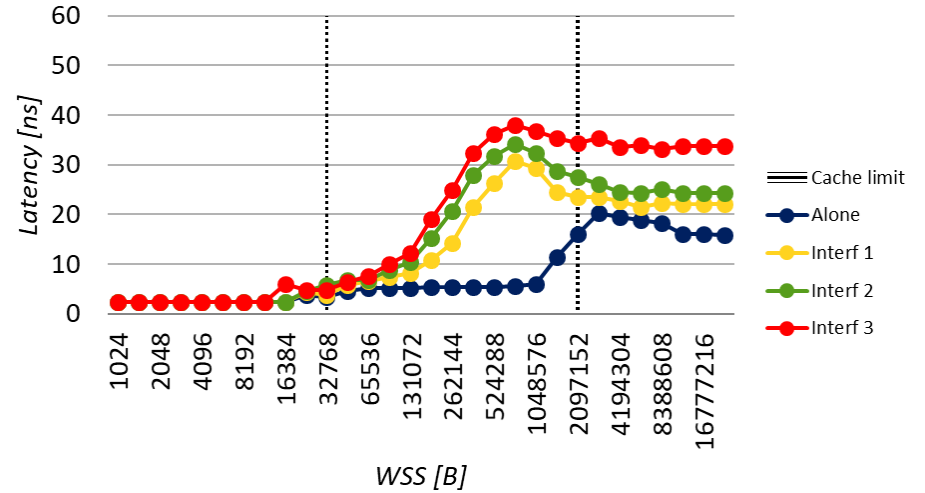


Tegra X1

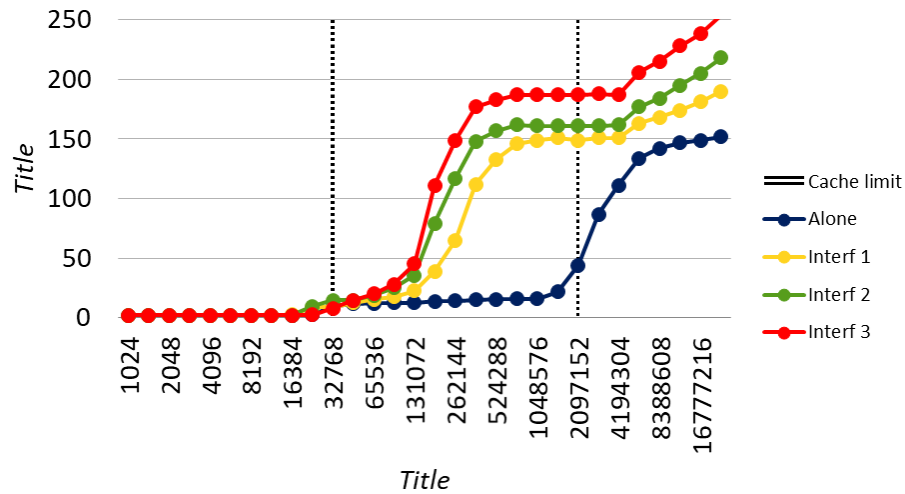
A1 - sequential read, sequential interference



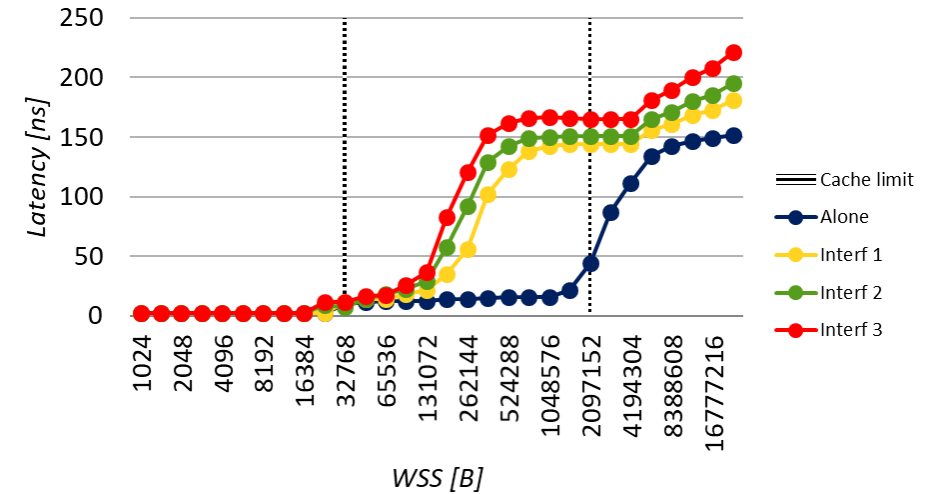
A3 - sequential read, random interference



A2 - random read, sequential interference

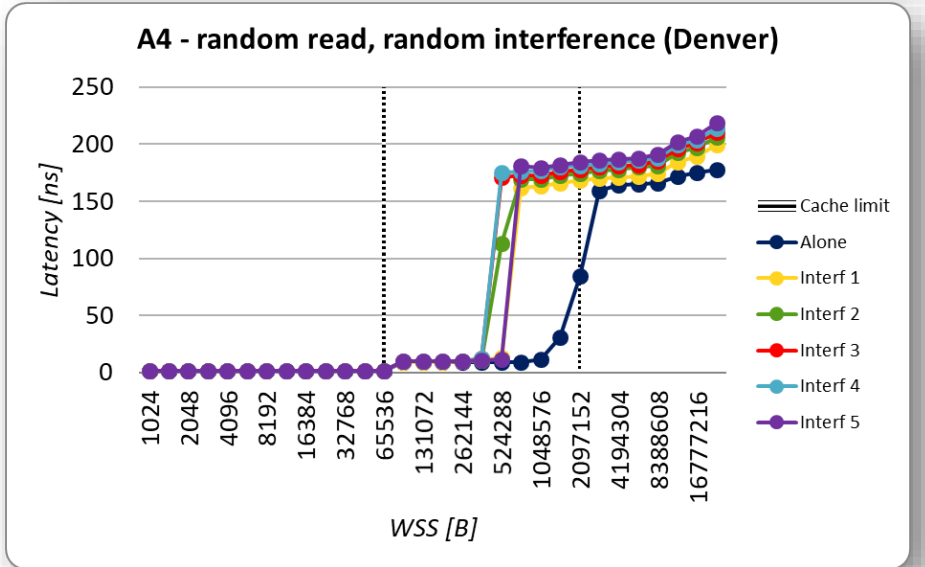
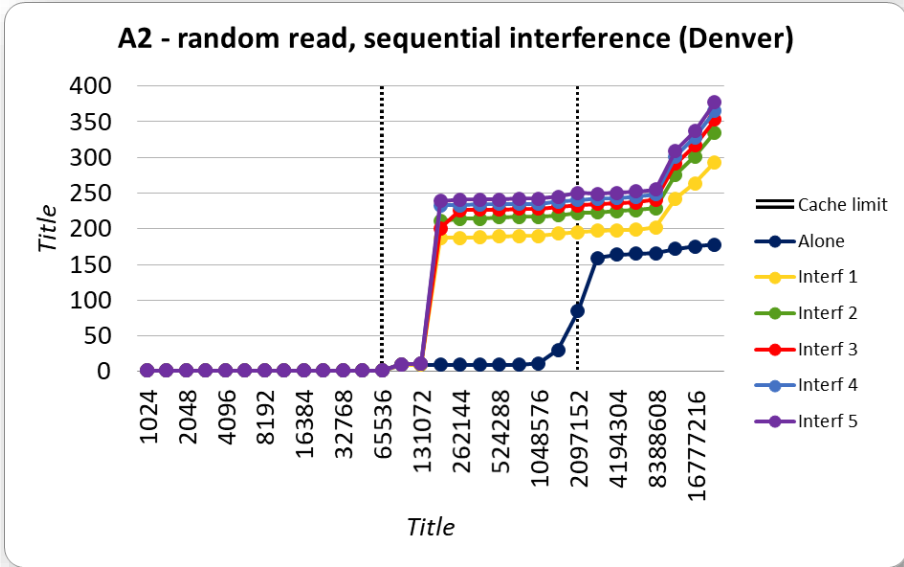
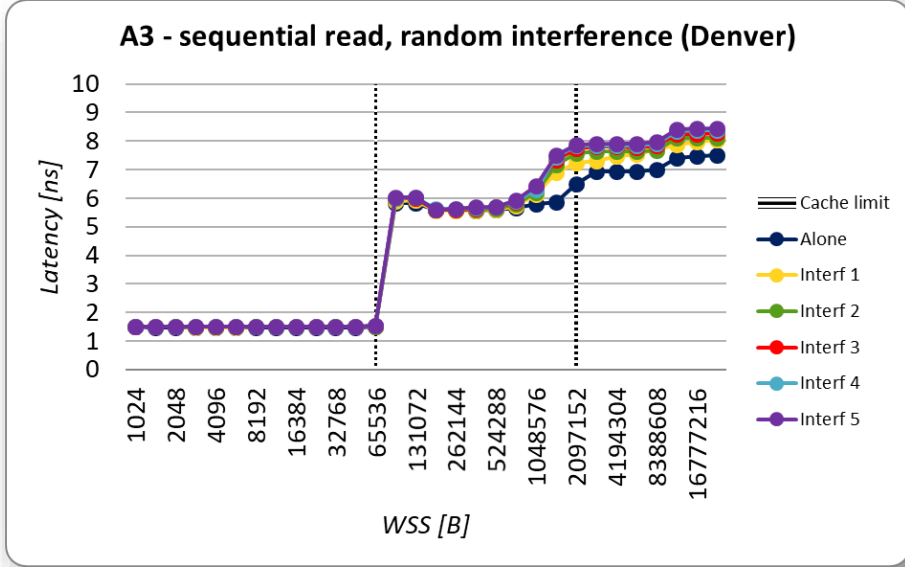
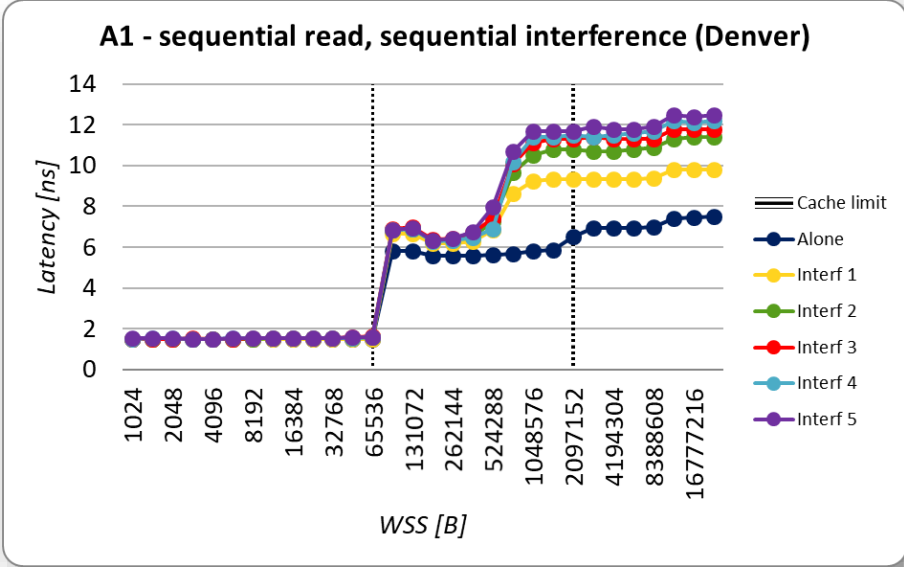


A4 - random read, random interference





Tegra X2 - Denver





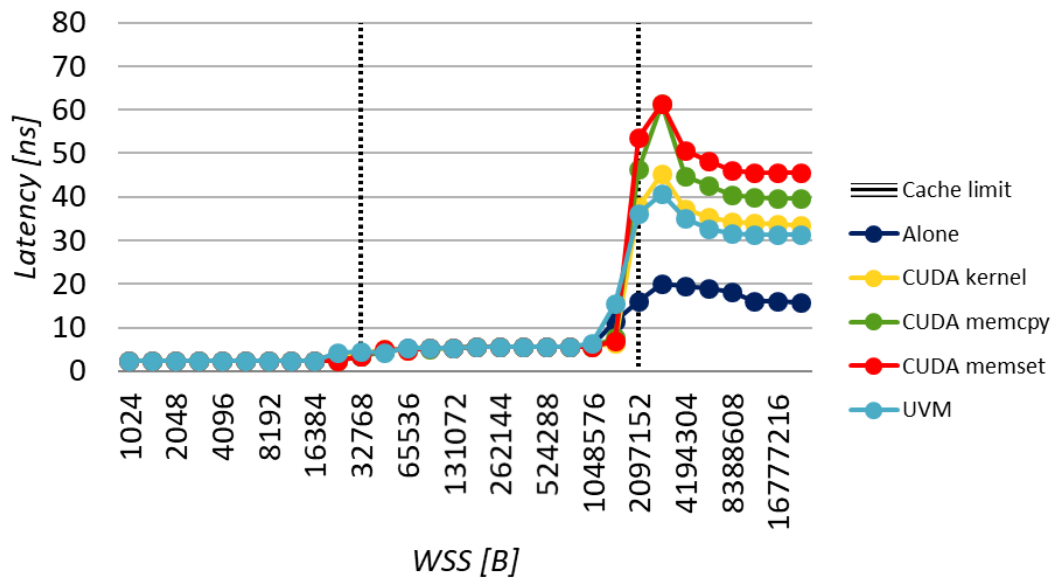
Test case B – iGPU interference on CPU

1. One CPU core reads **sequentially** within a variable working set, while the GPU accesses memory according to different paradigms:
 - CUDA memcpy
 - CUDA memcpy on UVM
 - CUDA memcpy on pinned mem
 - CUDA memset (0)
2. Same, but CPU core reads **randomly**

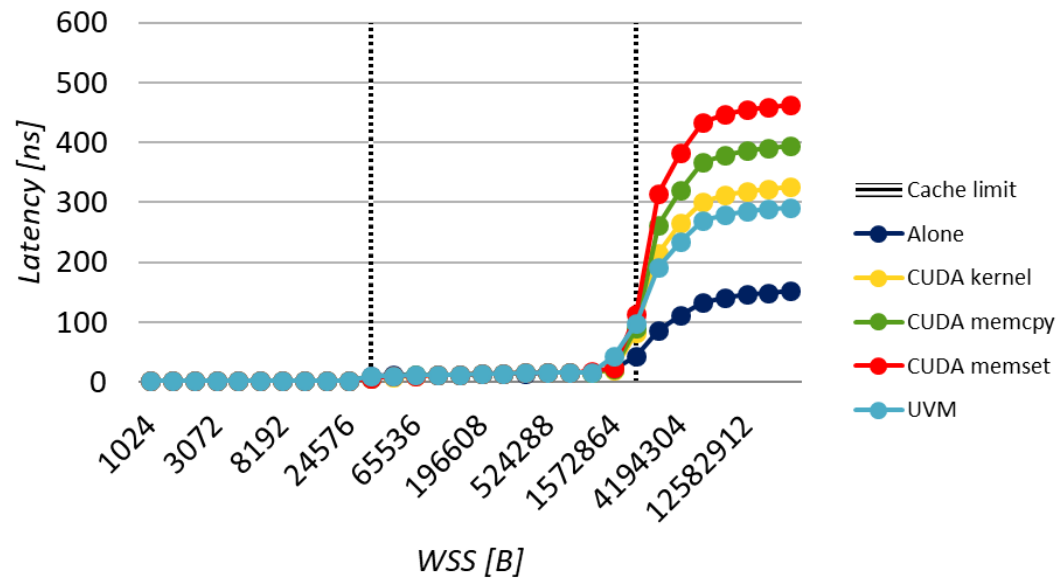


Tegra X1

B1 - sequential read, GPU interference



B2 - random read, GPU interference





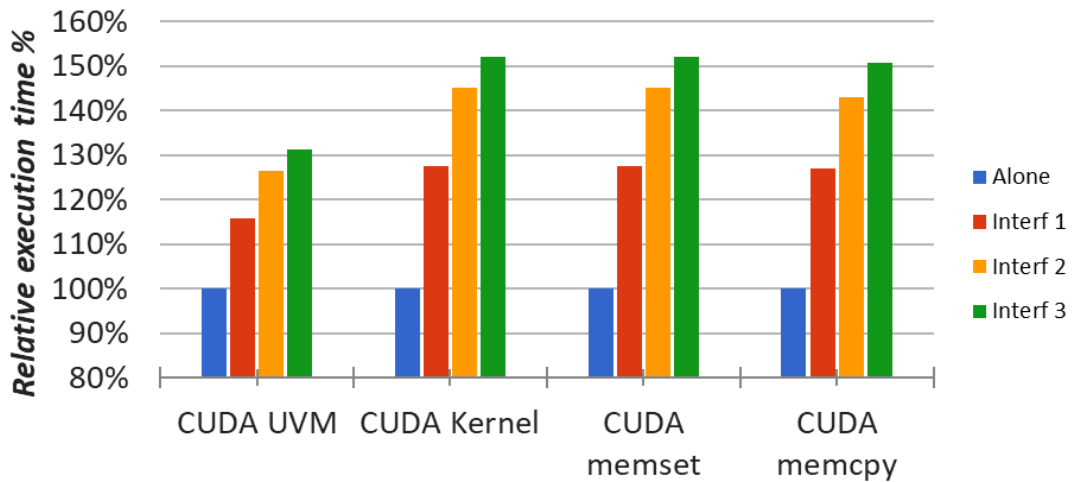
Test case C – CPU interference on iGPU

1. CPU generates **sequential** interfering mem accesses, while GPU accesses memory according to different paradigms:
 - CUDA memcpy
 - CUDA memcpy on UVM
 - CUDA memcpy on pinned mem
 - CUDA memset (0)
2. Same, but CPU core interference is **random**

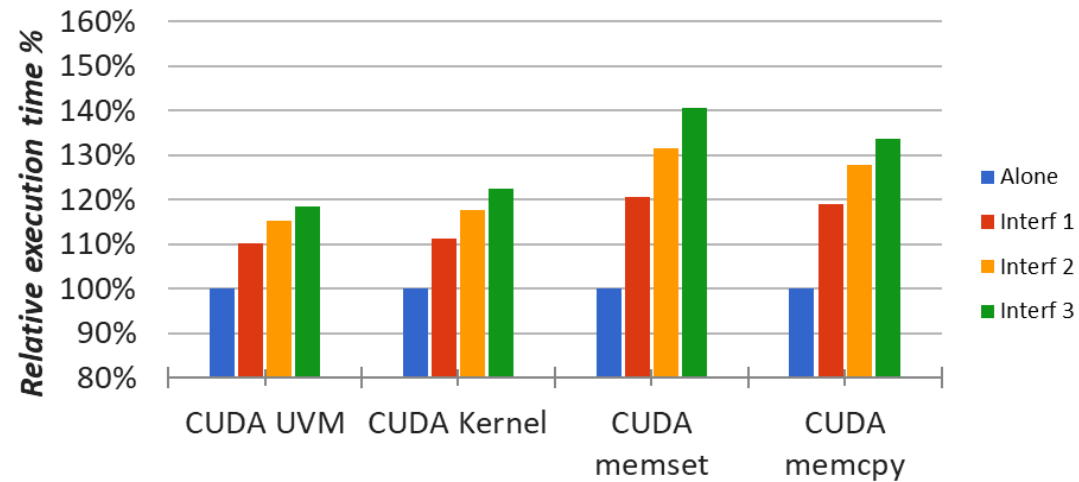


Tegra X1

C1 - GPU GM20b activity with sequential CPU interference

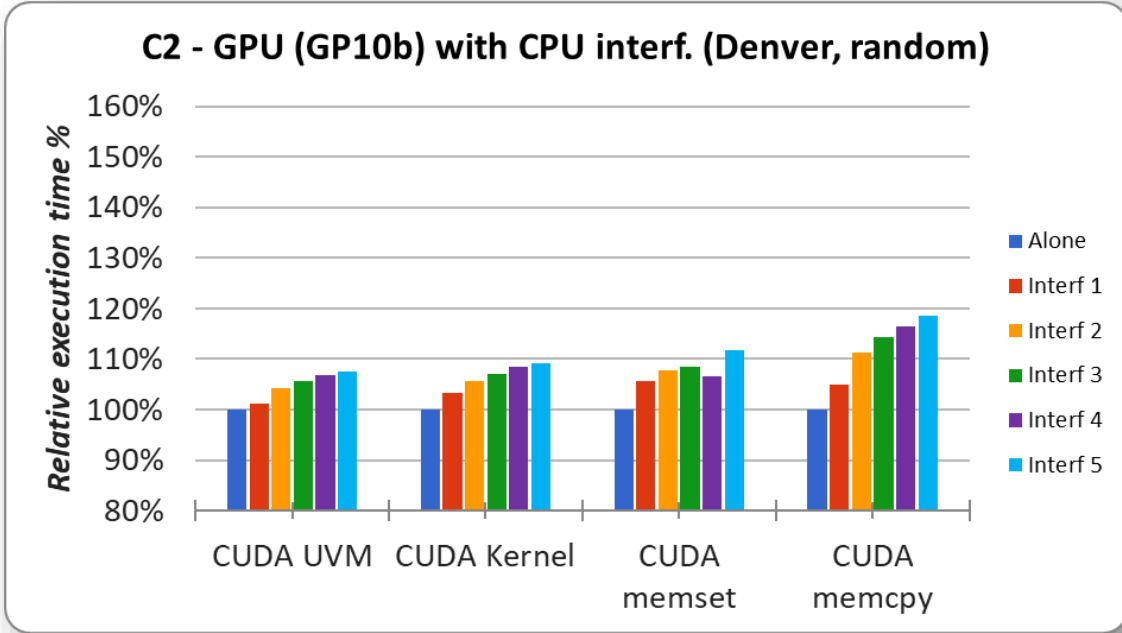
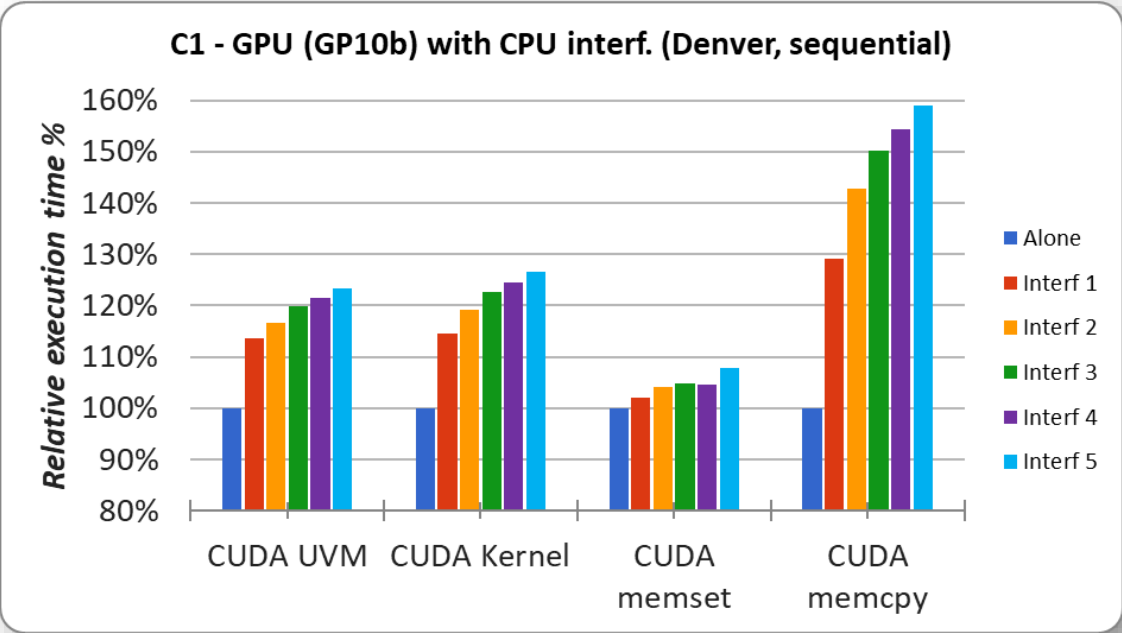


C2 - GPU GM20b activity with random CPU interference





Tegra X2 - Denver





Test 'C' - Intel i7-6700

