# MAD7: a Memory Architecture Simulator Targeted at Design Space Exploration

Hadrien Clarke[1,2], Antoine Trouvé[1] and Kazuaki Murakami[1,2]

[1] Institute of Systems, Information Technologies and Nanotechnologies

[2] Kyushu University

## Motivation

For many-core architectures:

- The number of accesses to the common main memory potentially increases
  - ▶ exacerbation of the "memory wall." [4]
- The number of cache banks increases and banks can be either private or shared [2]
  - ▶ the design space is larger than for mono-core architectures.
- Cycle-accurate simulators are either slow and/or limited [1]
  - ▶ unadapted tools for rapid design space exploration.

## Approach

- Focus on memory architecture only. The performance of a memory architecture is supposed mostly independent of the architecture of the cores.
- Trade-off between precision and simulation speed: fast simulations are more appropriate for design space exploration. (4)
- Coherency is roughly enforced and doesn't count as overhead [3].
- Evaluations of architectures are based on access costs reflecting technological choices (e.g. cache look-up time) and yielding a score instead of absolute access times.
- Multiple architectures can be compared using their respective scores for a given workload.

## References

[1] Shwe et al., "RExCache: Rapid Exploration of Unified Last-level Cache," ASP-DAC2013

[2] Zhang et al., "Does Cache Sharing on Modern CMP Matter to the Performance of Contemporary Multithreaded Programs?" PPoPP2010

[3] Martin et al., "Why On-Chip Cache Coherence is Here to Stay," Communications of the ACM, July 2012

[4] Wulf et al., "Hitting the Memory Wall: Implications of the Obvious," Computer Architecture News, March 1995

## Flow



## Design space exploration (1)



**Cache "shareness"** — same total capacity; fully shared cache; split into private caches

**Architecture depth** — same total capacity; n levels; n+1 levels

**Interconnect topology** — same hierarchy; point-to-point; bus; ring

### Legend

- C — a core
- $ — a cache unit
- $ — a cache unit of a different level
- a network node
- a network link

### Other dimensions

Along with the above considerations, "traditional" parameters such as cache sizes, associativities and line sizes constitute additional dimensions to the design space.

## Traces (2)



We use direct feed since logs of multiple threads can get large quickly.

Traces can be fed to several instances of the simulator to simulate multiple architectures at the same time.

Traces of all threads are transmitted serially

Traces are re-parallelized on a trace frame basis within the simulator

UNIX PIPE

... 2 LD 0xf007ba11ca5cade5 ... 0 ST 0x5ca1ab1ef007ba11 ... 9 ST 0xdac0ffee15cacad0 ...

thread ID | trace type | trace operand

trace frame

## Reproducibility (3)



**Figure**: Access counts from C_0 to the RAM for 100 consecutive runs; average access count after r runs; 0.75 and 1 standard deviation;

Noise mostly comes from the host's OS' scheduler

At least 4 runs are necessary for the outputs to be within 0.8 standard deviation

**Architecture**:
2 cores
private L1 caches (WB)
shared L2 cache (4 banks)

**Benchmark**:
PARSEC Blackscholes
4 threads
Simsmall
100 runs

**Figure**: Minimum required number of runs for the average of the outputs to be within 0.8 standard deviation