# Memory system performance

Petar Radojkovic
BSC memory systems team leader
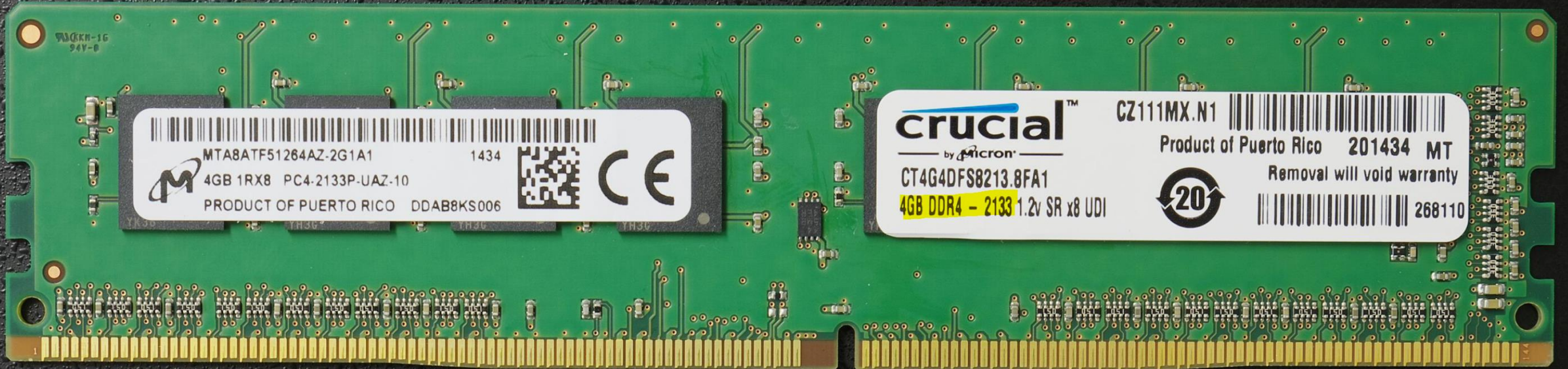
DEEP-SEA weekly seminars

5. November 2021.

# What is memory system performance?

# Memory performance = Bandwidth



- DDR4-2133
  - Bandwidth = 2133 * 64Byte * …

# Memory performance = Latency

- Memory latency is what stalls the pipeline (execution)
- Memory wall is about *latency*



In 1995, Wulf and McKee published a four-page note entitled "Hitting the Memory Wall: Implications of the Obvious" in the (un-refereed) ACM SIGARCH *Computing Architecture News* [27]. The motivation was simple: at the time, researchers were so focused on improving cache designs and developing other latency-tolerance techniques that the computer architecture community largely ignored main memory systems. The article projected the performance impact of the increasing speed gap between processors and memory. The study predicted that if the trends held, even with cache hit rates above 99%, relative memory latencies would soon be so large that the processor would essentially always be waiting for memory — which amounts to "hitting the wall".

Wm. A. Wulf and S. A. McKee.
*Hitting the memory wall: Implications of the obvious.*
Computer architecture news, 1995.

# Memory latency is not a single number

- Memory latenc<u>ies</u>
    - 1: Lead-off (unloaded) latency
    - 2: Loaded latency

- Memory latency curve
    - Memory latency = f(memory system stress)



B. L. Jacob.
*The memory system: You can't avoid it, you can't ignore it, you can't fake it.*
Synthesis Lectures on Computer Architecture, 2009.

# Memory latency curves

- Actual measurements
  - Dual socket Intel Xeon 8260 CPU (Cascade Lake)
    - 24 cores @ 2.4 GHz
  - 6 DDR4 memory channels per socket



- Measured with memory stressing benchmark (developed by the BSC)
  - Enhanced Stream benchmark (memory stress, X-axis) & Pointer chasing benchmark (latency , Y-axis)

- Open source:
  https://github.com/bsc-mem/PROFET
  **New code release coming soon!**

# Memory performance: Looking inside the box



DDR4-2666

Datasheets: Theoretical BW

Stream: Sustained BW

Pointer chasing: Lead-off latency

latency [ns]

bandwidth [MB/s]

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Use case 1: Application profiling

# Use case 1: Application profiling

# Use case 2:
# What (not) to expect from a novel memory system

- DDR4 and MCDRAM at KNL platform
  - *"MCDRAM provides an N-fold higher performance"*



Application using DDR4
$(BW_{used}^{DDR4}, Lat_{mem}^{DDR4})$

Memory access latency

Application using MCDRAM
$(BW_{used}^{MCDRAM}, Lat_{mem}^{MCDRAM})$

Used memory bandwidth

- Similar-ish analysis can be done for DDR4 and Optane (Storage Class Memory)

# Use case 3:
# Data distribution in heterogeneous memory systems

- How to distribute red, blue and green data structures in heterogeneous memory system
    - DDR4 and MCDRAM at KNL platform



Application using DDR4
$(BW_{used}^{DDR4}, Lat_{mem}^{DDR4})$

Application using MCDRAM
$(BW_{used}^{MCDRAM}, Lat_{mem}^{MCDRAM})$

Memory access latency

Used memory bandwidth

- Similar-ish analysis can be done for DDR4 and Optane (Storage Class Memory)

# Use case 4:
# Memory system design space exploration

- Overall system performance = f (mem system)?

- *"Can we estimate application performance with <u>different (new) main memory</u>?"*
    - Example: High-end CPU with DDR4-XXXX
        - DDR4 (different frequencies)
        - DDR5
        - HBM2/3
        - LPDDR4/5
        - …
        - Optane
        - Memory over CXL (CCIX, etc.)
        - …
        - Novel memory designs

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Our idea: PROFET – PROfile & EsTimate

- **PROfile**
  - Memory latency curves
    - Conventional memory system
    - Novel memory system

  - Application running on an actual system

- **EsTimate**
  - Application performance, power and energy with the novel memory system



Memory access latency vs. Used memory bandwidth, showing Application using DDR4 $(BW_{used}^{DDR4}, Lat_{mem}^{DDR4})$ and Application using MCDRAM $(BW_{used}^{MCDRAM}, Lat_{mem}^{MCDRAM})$

# PROFET modeling

- Based on strong profiling capabilities of high-end processors (hardware counters)

- Parameters that we cannot measure
  - Set the boundaries based on the microarchitectural specification (e.g., row-buffer size)

- Analytically model performance – Equations

- Solve the equations

- Published:
  Radulovic et al., ***PROFET: Modeling System Performance and Energy Without Simulating the CPU***. SIGMETRICS, 2019.

- Open source:
  https://github.com/bsc-mem/PROFET
  **New code release coming soon(-ish)! ;-)**

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

---

PROFET: Modeling System Performance
and Energy Without Simulating the CPU
34:13

Note that $\widehat{MLP}(Ins_{ooo})$ is a point estimate for $MLP$ based on the available information. If the point estimate is outside the valid range, between the lower and upper bounds described above, then it is corrected to lie in the range.

### 4.4 Performance as a function of latency

This section completes the analysis of out-of-order processor performance as a function of latency. We start by repeating Eq. 4, which gives the predicted CPI in terms of $Stalls_{LLC}$:

$$CPI_{tot}^{(2)} = CPI_{tot}^{(1)} + \frac{Miss_{LLC}}{Ins_{tot}} \times \left( Stalls_{LLC}^{(2)} - Stalls_{LLC}^{(1)} \right) \tag{4 again}$$

As remarked at the beginning of Section 4.3, in comparison with an in-order processor, an out-of-order processor has a more complex expression for $Stalls_{LLC}$, and this was given in Eq. 12:

$$Stalls_{LLC} = \frac{1}{MLP} \times (Pen_{mem} - CPI_0 \times Ins_{ooo}) \tag{12 again}$$

Finally we replace the $MLP$ parameter in this equation with the point estimate in Eq. 15:

$$\widehat{MLP}(Ins_{ooo}) = \frac{Miss_{LLC}}{Ins_{tot}} \times Ins_{ooo} + 1 \tag{15 again}$$

In fact, as explained in Section 4.3.3, this value is restricted to lie between the lower and upper bounds given in that section. For the sake of clarity, we consider the more common case for which it is not necessary.

Combining Eq. 4, Eq.12 and Eq.15, and assuming that $Ins_{tot}$, $Miss_{LLC}$, $CPI_0$, $Ins_{ooo}$ and $MLP$ do not change when moving from one memory system configuration to another, then $CPI_{tot}^{(2)}$ can be calculated as:

$$CPI_{tot}^{(2)} = CPI_{tot}^{(1)} + \frac{Pen_{mem}^{(2)} - Pen_{mem}^{(1)}}{Ins_{ooo} + Ins_{tot}/Miss_{LLC}} \tag{16}$$

This equation is written in terms of the memory access penalty, $Pen_{mem}$, but at the system level, outside a detailed analysis of a particular processor's pipeline, only $Lat_{mem}$ is relevant. Recall that $Pen_{mem}$ was defined to be the memory access latency, $Lat_{mem}$ minus the cost of an LLC hit. We note, therefore, that the expression $Pen_{mem}^{(2)} - Pen_{mem}^{(1)}$ is equal to $Lat_{mem}^{(2)} - Lat_{mem}^{(1)}$. Taking account of this and rewriting in terms of the IPC instead of the CPI gives:
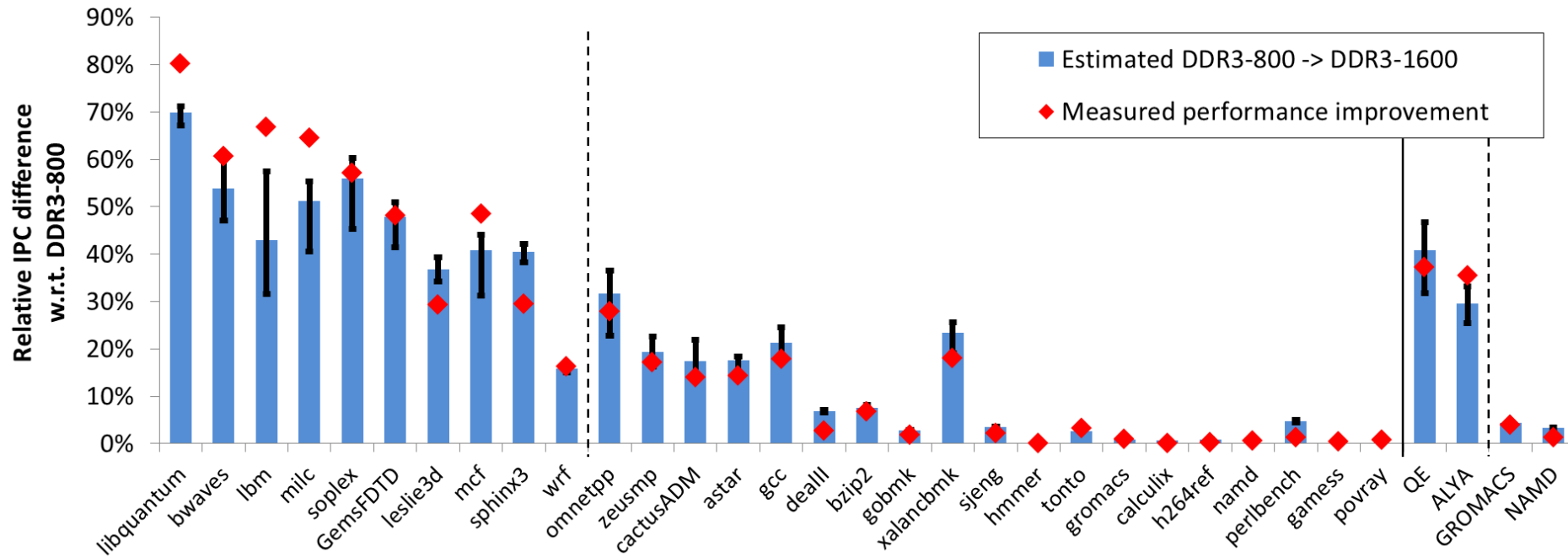
$$IPC_{tot}^{(2)} = \frac{IPC_{tot}^{(1)}}{1 + IPC_{tot}^{(1)} \times \frac{Lat_{mem}^{(2)} - Lat_{mem}^{(1)}}{Ins_{ooo} + Ins_{tot}/Miss_{LLC}}} \tag{17}$$

The various values in Eq. 17, $IPC_{tot}^{(1)}$, $Ins_{tot}$, and $Miss_{LLC}$ are known because they were measured on the baseline memory configuration. All other inputs to PROFET, such as $Ins_{ROB}$, $MSHR$, $CPI_{min}$ (see Table 3) appear in the upper and lower bounds of $Ins_{ooo}$.[5]

Eq. 17 is plotted in Figure 7. The $x$ axis is the target system memory latency, $Lat_{mem}^{(2)}$, and the $y$ axis is the predicted IPC, $IPC_{tot}^{(2)}$. Eq. 17 is a function of the independent parameter $Ins_{ooo}$, which we cannot measure or calculate exactly. We bounded its value in the previous section, and varying it between the lower and upper bounds gives the family of curves shown in the figure. Note that the case of $Ins_{ooo} = 0$ corresponds to an in-order processor. This can be seen by comparing Eq. 17 and Eq. 6. As indicated on the figure, when the target memory latency is the same as the baseline memory latency, $Lat_{mem}^{(1)}$, PROFET correctly "predicts" the measured IPC to be that of the baseline system, $IPC_{tot}^{(1)}$.
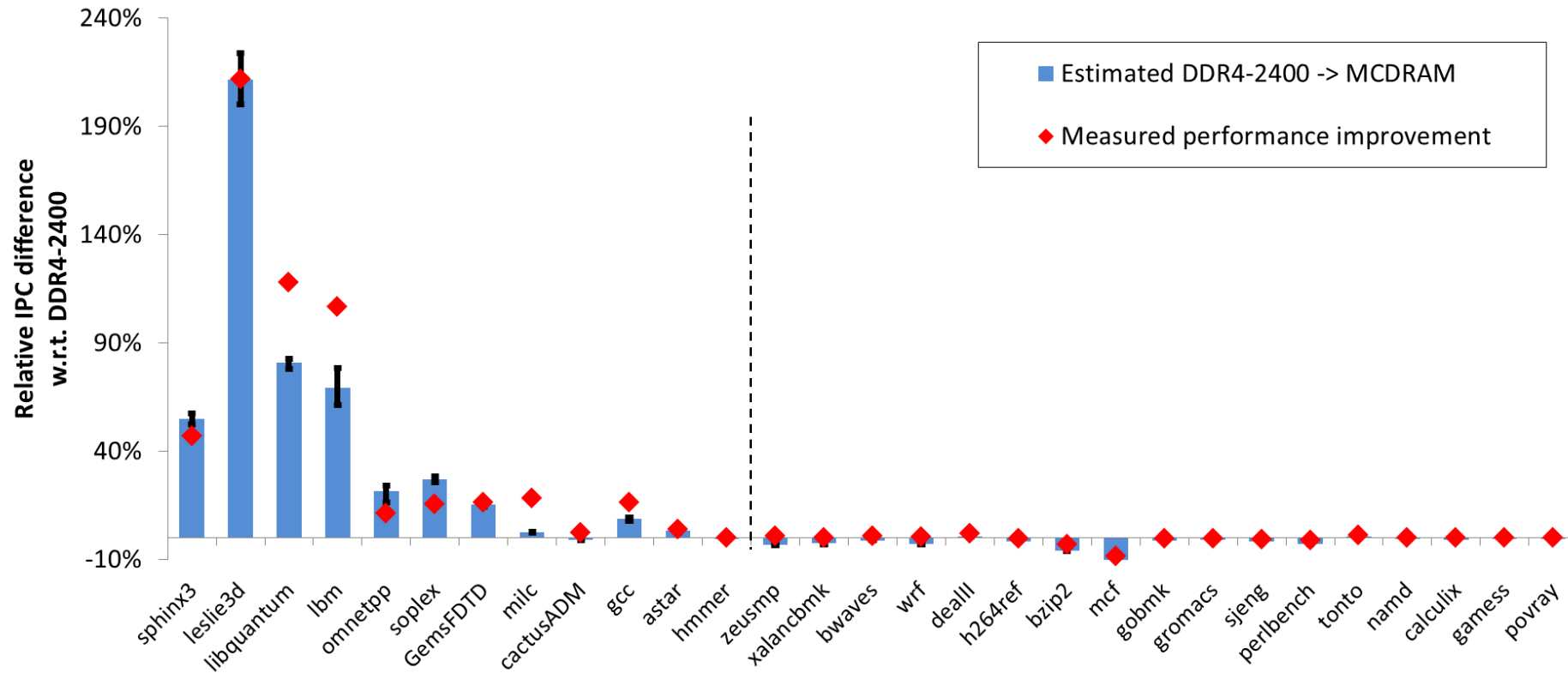
[5]Some of the input parameters appear only in the upper or lower bounds of $MLP$, which are not in Eq. 17 but considered in the full PROFET model.

# Evaluation: Sandy Bridge DDR3-800 → DDR3-1600



- Performance estimation error
  - High-bandwidth benchmarks: 5.3%
  - All benchmarks: 2.9%

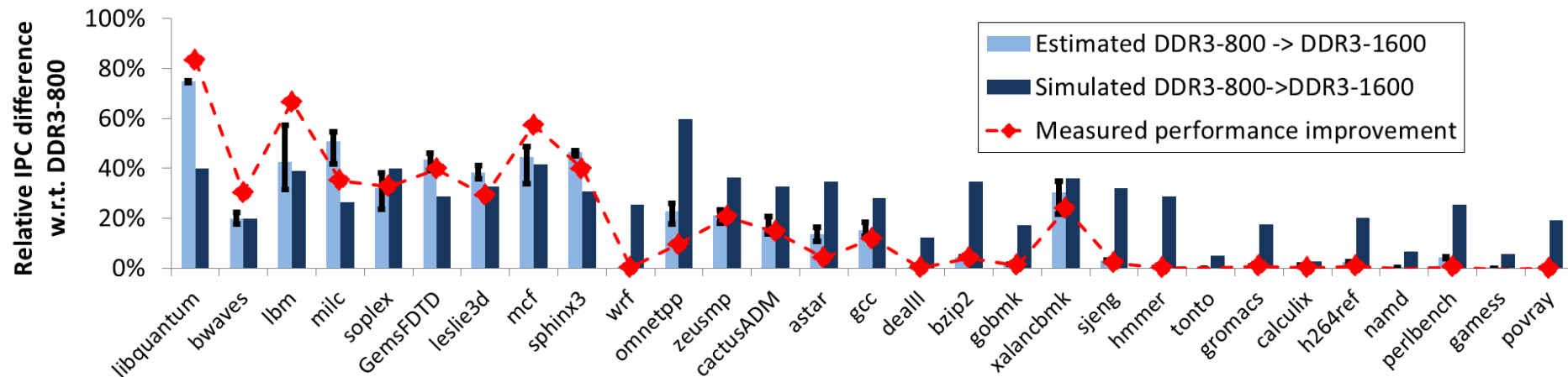- Power and energy
  - All benchmarks: 1% and 2%

# Evaluation: KNL DDR4-2400 → MCDRAM



- Performance estimation error
  - High-bandwidth benchmarks: 7%
  - All benchmarks: 3.8%

Barcelona
Supercomputing
Center
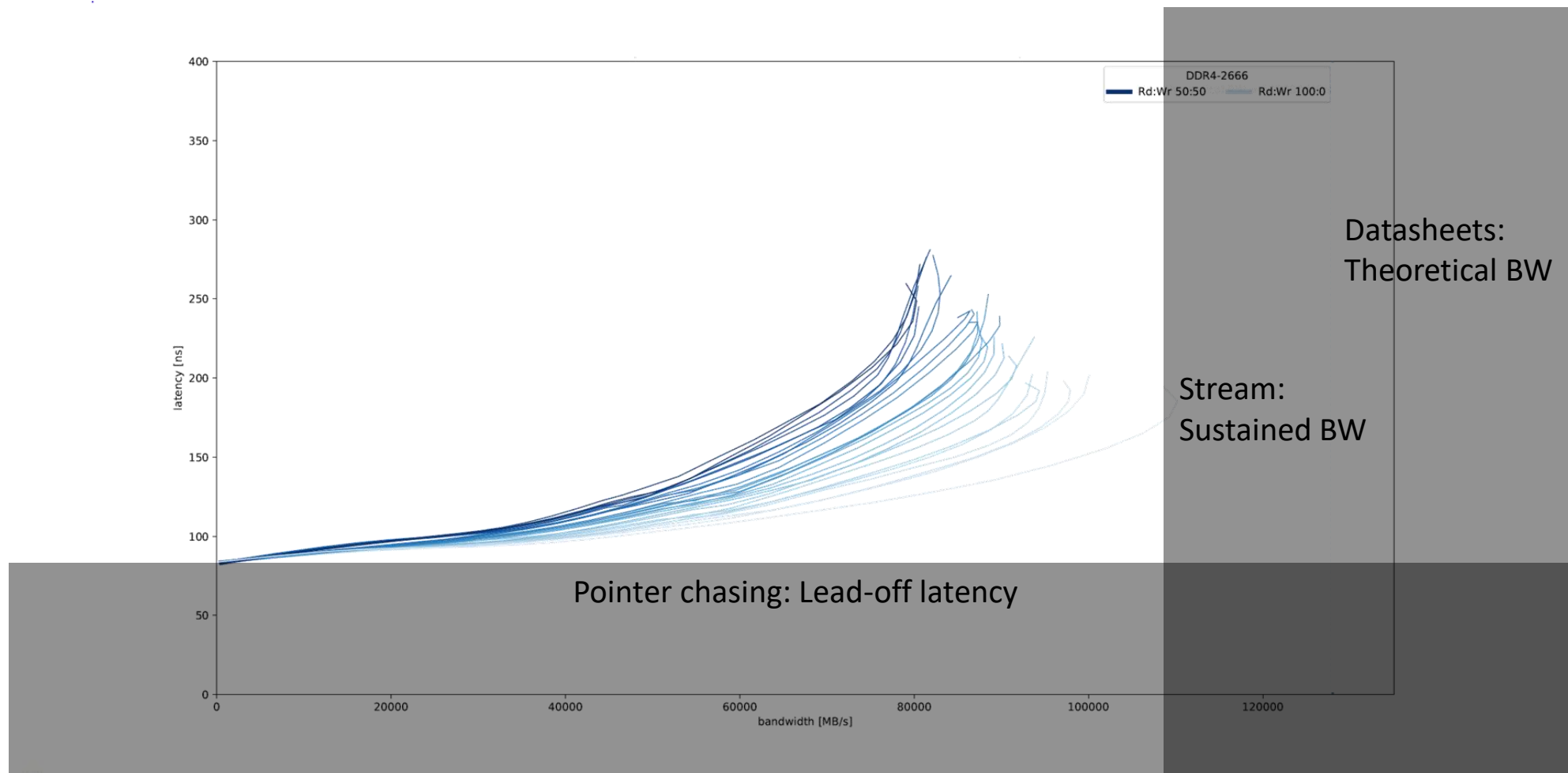Centro Nacional de Supercomputación

# PROFET *vs.* Simulators

- Simulators:
  - CPU simulator: Zsim - updated and validated for Intel Sandy Bridge CPU
  - Main memory simulator: DRAMSim2
  - Workloads: SPEC CPU2006, 150 billions of instructions
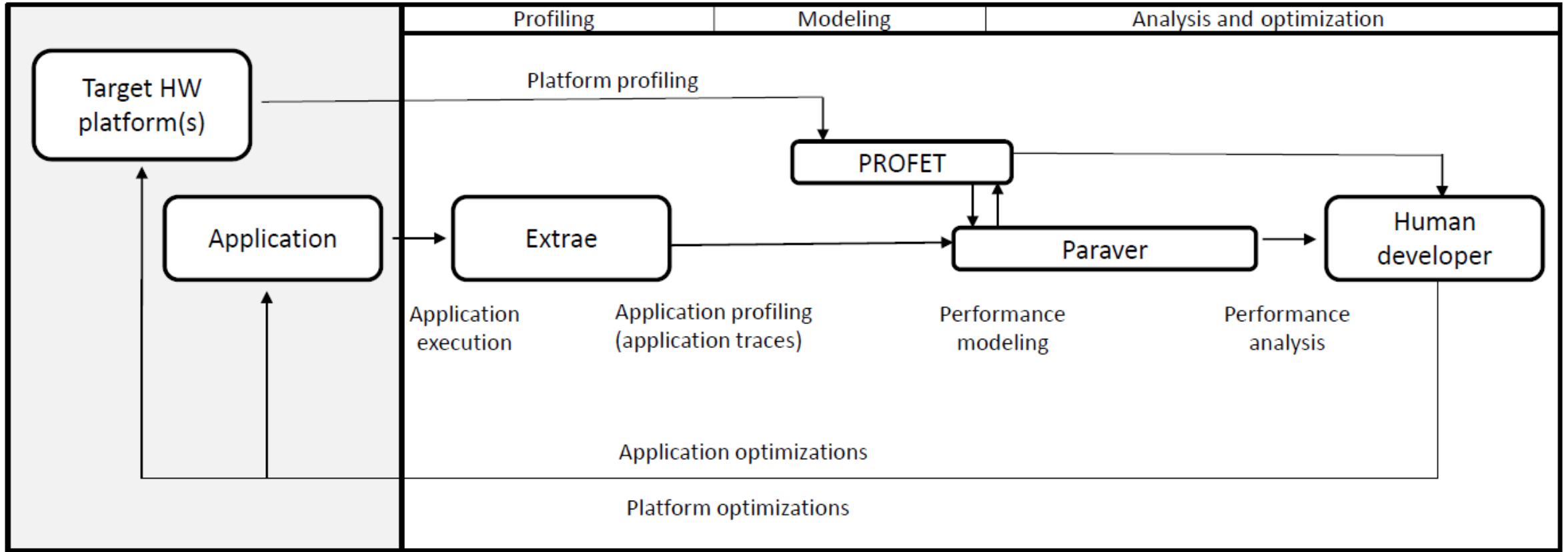
- DDR3-800 → DDR3-1600



- Average error:
  - PROFET: 3.6%
  - Simulations: 15.7%

- Model estimations follow the trend better than simulations
- PROFET faster than the simulator: **Three orders of magnitude**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Memory performance:
# Looking inside the box

# Extrae-PROFET-Paraver optimization cycle

# Thank you

petar.radojkovic@bsc.es

# The roofline model



- DDR4 and MCDRAM at KNL platform
  - *"MCDRAM provides an N-fold higher performance"*