## Characterizing Applications based on Cache-Aware Roofline Model

#### Leonel Sousa,

Aleksandar Ilic and Diogo Marques





INESC-ID, Instituto Superior Técnico (IST), Universidade de Lisboa, Portugal

## Keynote PPAM, Bialystok, September 2019







- Energy-efficient architectures: Sprinting towards Exascale Systems
  - Heterogeneity: compute and memory capability, specialization (CPU, GPU, big.LITTLE, FPGA)
  - Performance, power, energy and efficiency trade-offs
- **Application design**: Exploiting architecture diversity to reach efficiency
  - Different computational and memory requirements
  - **Painful** optimization and characterization (for each architecture)
- How far can we go: Performance, Power, Energy and Efficiency Maximums?
  - In a simple, insightful and fast way (allowing the first-order analysis)









#### **RECAP: ORIGINAL ROOFLINE MODEL\***

# CACHE-AWARE ROOFLINE MODEL

- Performance
- EXTENSIONS: POWER AND ENERGY-EFFICIENCY

#### APPLICATION CHARACTERIZATION WITH CACHE-AWARE ROOFLINE MODEL

- APPLICATION-DRIVEN CARM
- ISO3DFD CASE STUDY
- PROXY-APPLICATIONS FROM EXASCALE COMPUTING PROJECT
- BAYESIAN K2 SCORE

## **ON-GOING (FUNDED) PROJECTS**

\* Williams S., et. al., "Roofline: An insightful visual performance model for multicore architectures", CACM (2009)



#### Multi-cores and the Roofline Model







- General-purpose processors with complex memory hierarchy
  - Multiple cores with powerful out-of-order engines
  - Several levels of memory hierarchy: private/shared caches + deeper (and diverse) memory levels
- **OBSERVATION:** Computations and communication (data transfers) simultaneously performed
  - The overall execution time can be limited either by the time to compute or by the time taken to transfer data
  - Different Roofline Models observe memory traffic differently!

#### **DRAM Roofline Model**







- **DRAM ORM:** DRAM-based Original Roofline Model
  - Memory traffic: bytes between DRAM and LLC
  - Memory Bandwidth: DRAM to LLC, i.e., DRAM bandwidth
  - Compute performance: Flops delivered by the core(s)
  - Intensity (x-axis): Flops/DRAMBytes





## **Original Roofline Model**



2-6

2-4

2-2

2<sup>0</sup>





## **Original Roofline Model**







- Lx ORM: Lx-based Original Roofline Model
  - Memory traffic: bytes between two subsequent memory levels
  - Memory Bandwidth: Lx bandwidth
  - Compute performance: Flops delivered by the core(s)
  - Intensity (x-axis): Flops/LxBytes





#### **Hierarchical ORM**







- Hierarchical ORM: Several Lx-based Original Roofline Models in a single plot
  - Memory Bandwidth: Lx bandwidth (for each ORM)
  - Compute performance: Flops delivered by the core(s)
  - Intensity (x-axis): Flops/LxBytes



#### **Hierarchical ORM**







- Hierarchical ORM: Several Lx-based Original Roofline Models in a single plot
  - Memory Bandwidth: Lx bandwidth (for each ORM)
  - Compute performance: Flops delivered by the core(s)
  - Intensity (x-axis): Flops/LxBytes
  - Application characterization: As many points as memory levels (one for each ORM)









#### **RECAP: ORIGINAL ROOFLINE MODEL\***

## CACHE-AWARE ROOFLINE MODEL

- PERFORMANCE\*
- EXTENSIONS: POWER AND ENERGY-EFFICIENCY
- APPLICATION CHARACTERIZATION WITH CACHE-AWARE ROOFLINE MODEL
- APPLICATION-DRIVEN CARM
- ISO3DFD CASE STUDY
- PROXY-APPLICATIONS FROM EXASCALE COMPUTING PROJECT
- BAYESIAN K2 SCORE

# **ON-GOING (FUNDED) PROJECTS**

\*A. Ilic, F. Pratas and L. Sousa "Cache-ware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)



## Cache-aware Roofline Model (CARM)







- CARM: visual representation of the **limits of parallel processing** on contemporary multi-cores
  - Considers memory traffic and computation from the consistent architecture point of view (cores)
  - Relates the peak compute performance and realistically achievable bandwidth with Arithmetic Intensity (AI)
  - Unifies the complete memory hierarchy in a single plot model
- **CONSTRUCTION**: How to obtain these bandwidth values? (only B<sub>L1->C</sub> directly derivable from data sheets)

\* Ilic A., Pratas F., Sousa L., "Cache-aware Roofline Model: Upgrading the loft", IEEE CAL (2014)

Iİİ

## Developed Tools: Methodology and micro-benchmarks

**Inesc id** LISBOA

- CARM-oriented tools: construction, validation and application characterization

- SchedMon<sup>1</sup>): Software tool for near-OS counter-based monitoring (multiplexing and shed events)
- LARM<sup>2</sup>): CARM bandwidth and FP performance micro-benchmarks (NUMA/KNL CARM)
- CARM validation tools and micro-benchmarks<sup>3)</sup>



<sup>1)</sup> Taniça L., et. al. "SchedMon: A Performance and Energy Monitoring Tool for Modern Multi-cores", MuCoCoS/Euro-Par (2014)
 <sup>2)</sup> Denoyelle N., et. al., "Modeling Large Compute Nodes with Heterogeneous Memories in CARM", PMBS, SC (2017) and IEEE TPDS (2018)
 <sup>3)</sup> Marques D., "Analyzing Performance of Multi-cores and Applications with Cache-aware Roofline Model", HPBench, HPCS (2017)

Iİİ

#### **Cache-aware Roofline Model**







#### Total Cache-aware Roofline Model

- Includes **all transitional states** (traversing the memory hierarchy and filling the pipeline)
- Single-plot modeling for **different** types of compute and memory **operations**

- Insightful single plot model
  - Shows performance limits of multicores
  - Redefined AI: flops and bytes as seen by core
  - Constructed once per architecture
- Considers complete memory hierarchy
  Influence of caches and DRAM to performance
- Applicable to other types of operations
  not only floating-point
- Useful for:
  - Application characterization and optimization
  - Architecture development and understanding

\* Ilic A., Pratas F., Sousa L., "Cache-aware Roofline Model: Upgrading the loft", IEEE CAL (2014)

### Cache-aware Roofline Model: Interpretation







Optimization hints

- Memory: improve access pattern, use of caches
- Compute: vectorization, FMAs, parallelization
- Shady: all kinds of everything (mem+comp)

- Performance analysis with CARM
  - Application (kernel) is a single "dot"
  - In respect to their AI and FP Performance
- Draw an imaginary vertical line at app Al
  - Arithmetic intensity: Property of the application
  - Should not change (unless the algorithm changes)
- Intersected roofs: potential bottlenecks
  - Priority to the roofs above
  - Roofs below are also important!
- Optimization: "Break the above roofs"
  - Optimizations should improve the performance
  - Points move up in the CARM plot

\* Marques D., et. al., "Performance Analysis with Cache-Aware Roofline Model in Intel Advisor", HPCS (2017)



### Cache-aware Roofline Model in Intel Advisor





Performance (GFLOPS)





- Automatic construction (from NHL to KBL-X)
- Break-down by application phases, loops and functions (hierarchical feature)
- In-depth application profiling and optimization hints

- Performance analysis with CARM
  - Applications (kernels) as single "dots"
  - In respect to their AI and FP Performance
- Draw an imaginary vertical line at app Al
  - Arithmetic intensity: Property of the application
  - Should not change (unless the algorithm changes)
- Intersected roofs: potential bottlenecks
  - Priority to the roofs above
  - Roofs below are also important!
- Optimization: "Break the above roofs"
  - Optimizations should improve the performance
  - Points move up in the CARM plot

"incredibly useful diagnosis tool (that can guide the developers in the application optimization process), ensuring that they can squeeze the maximum performance out of their code with minimal time and effort."

#### **Cache-aware Roofline Model: Power Consumption**

- **TÉCNICO** LISBOA IJ inescid
- **Performance:** Computations (*flops*) and communication (*bytes*) overlap in time ۲
- **Power consumption**: Superposed contributions of flops and bytes ۲



\* Ilic, A., Pratas, F. and Sousa, L., "Beyond the Roofline: Cache-aware Power & Energy-Efficiency Modeling...", IEEE Transactions on Computers (2016)

#### **Efficiency CARMs**







![](_page_16_Picture_4.jpeg)

#### **Cache-aware Roofline Model: Use Cases**

![](_page_17_Figure_1.jpeg)

TÉCNICO LISBOA

IJ

Application Characterization

![](_page_17_Figure_3.jpeg)

![](_page_17_Picture_4.jpeg)

\* Ilic, A., Pratas, F. and Sousa, L., "Beyond the Roofline: Cache-aware Power & Energy-Efficiency Modeling…", IEEE Transactions on Computers (2017) \* Antão, D., et.al., "Monitoring Performance and Power for Application Characterization with CARM", PPAM'13

#### Cache-Aware Roofline Model: Extensions

![](_page_18_Picture_1.jpeg)

![](_page_18_Picture_2.jpeg)

![](_page_18_Figure_3.jpeg)

\* Ilic A., et.al., "Beyond the Roofline: Cache-aware Power & Energy-Efficiency...", IEEE Transactions on Computers (2017)

\* Lopes A., et.al, "Exploring GPU performance, power and energy-efficiency bounds with CARM"", ISPASS (2017)

\* Denoyelle N., et.al., "Modeling Non-Uniform Memory Access on Large Compute Nodes with the Cache-Aware Roofline Model", IEEE

Transactions on Parallel and Distributed Systema (2018) 9/10/2019

ĮĬ

![](_page_19_Picture_0.jpeg)

![](_page_19_Picture_1.jpeg)

![](_page_19_Picture_2.jpeg)

- (Total) Performance CARM
- (Total) Power CARM: for several domains, i.e., power of cores, uncore power and complete package power
- Energy-Efficiency CARM: Performance + Power Domains
- Energy, Power-efficiency and EDP-based CARMs
- DVFS, GPU and NUMA CARMs

**ON-GOING WORK** 

- CARM FOR ARM, FPGAS, COMPLETE SYSTEM ...

![](_page_19_Picture_10.jpeg)

#### Outline

![](_page_20_Picture_1.jpeg)

![](_page_20_Picture_2.jpeg)

**RECAP: ORIGINAL ROOFLINE MODEL\*** 

CACHE-AWARE ROOFLINE MODEL

- PERFORMANCE\*
- EXTENSIONS: POWER AND ENERGY-EFFICIENCY

APPLICATION CHARACTERIZATION WITH CACHE-AWARE ROOFLINE MODEL

- APPLICATION-DRIVEN CARM
- ISO3DFD CASE STUDY
- PROXY-APPLICATIONS FROM EXASCALE COMPUTING PROJECT
- BAYESIAN K2 SCORE

ıſi

# **ON-GOING (FUNDED) PROJECTS**

\*A. Ilic, F. Pratas and L. Sousa "Cache-ware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014) Recently submitted publication

#### **Application-driven CARM**

![](_page_21_Picture_1.jpeg)

![](_page_21_Picture_2.jpeg)

#### **Applications**

Can they exploit the modeled absolute maximums?

- Large set of phases with diverse characteristics
- Different instruction mix, vectorization, SP/DP and LD/ST balance, FP share, memory access pattern...

#### Hardware

Do the Rooflines reflect the application demands?

- Maximums vary with utilization/execution scenario
- **Components/subsystems** differently exercised: ports, compute units, front-end, back-end, sockets ...
- **Memory** subsystem: deep and diverse hierarchy, caches (private/shared), DRAM, NUMA, HBM...

#### Insightful Micro-architecture Modeling Existing approaches model the absolute maximums

- Disjoint Roofline Methodologies (Cache-aware, Classic ORM, Hierarchical, Integrated,...)
- May provide misleading optimization guidelines
- Inconclusive bottleneck detection

![](_page_21_Picture_16.jpeg)

### **Application-driven CARM**

![](_page_22_Picture_1.jpeg)

![](_page_22_Picture_2.jpeg)

#### **Applications**

Can they exploit the modeled absolute maximums?

- Large set of phases with diverse characteristics
- Different instruction mix, vectorization, SP/DP and LD/ST balance, FP share, memory access pattern...

#### Hardware

#### Do the Rooflines reflect the application demands?

- Maximums vary with utilization/execution scenario
- **Components/subsystems** differently exercised: ports, compute units, front-end, back-end, sockets ...
- **Memory** subsystem: deep and diverse hierarchy, caches (private/shared), DRAM, NUMA, HBM...

#### Insightful Micro-architecture Modeling Existing approaches model the absolute maximums

- Disjoint Roofline Methodologies (Cache-aware, Classic ORM, Hierarchical, Integrated,...)
- May provide misleading optimization guidelines
- Inconclusive bottleneck detection

![](_page_22_Picture_16.jpeg)

#### 9/10/2019

#### Application-driven Roofline

![](_page_22_Figure_19.jpeg)

#### **Application-driven CARM**

![](_page_23_Picture_1.jpeg)

![](_page_23_Picture_2.jpeg)

State-of-the-art CARM Scalar Add Peak: 82.38 GELOPS? 82.38 bound by both memory and compute? 0.4 -FLOP/Byte (Arit etic Intensity) 0.004 2.09  $\begin{array}{c} \mbox{FP Performance} \left[ \mbox{GFLOPS/s} \right] \\ \mbox{FP Performance} \left[ \mbox{GFLOPS/s} \right] \\ \mbox{Sec} \left[ \mbox{G$ ZZZ DP SP AVX-512 Scalar FMA ADD ADD 2LD/ST Memory Bandwidth [GB/s] Memory Bandwidth [GB/s]  $^{70}_{12}$   $^{70}_{11}$   $^{70}_{12}$   $^{70}_{$ LD AVX-512

![](_page_23_Picture_4.jpeg)

L1

L2 L3 DRAM

![](_page_24_Figure_0.jpeg)

IJÎ

#### ISO-3DFD Case study

![](_page_25_Picture_1.jpeg)

![](_page_25_Picture_2.jpeg)

![](_page_25_Picture_3.jpeg)

ISO-3DFD: 3D Finite Difference Code with an Isotropic
 helps solving differential equations (seismic apps, wave propagation)
 <a href="https://software.intel.com/en-us/articles/eight-optimizations-for-3-dimensional-finite-difference-3dfd-code-with-an-isotropic-iso">https://software.intel.com/en-us/articles/eight-optimizations-for-3-dimensional-finite-difference-3dfd-code-with-an-isotropic-iso</a>

![](_page_25_Figure_5.jpeg)

DISCLAIMER: Optimization courtesy of Cédric Andreolli (Intel Corporation)

Experiments ran on Intel Xeon Gold 6140 (18 cores @ 2.3GHz) 4x16GB DDR4

![](_page_25_Picture_8.jpeg)

#### **ISO-3DFD Case study**

![](_page_26_Picture_1.jpeg)

![](_page_26_Picture_2.jpeg)

![](_page_26_Figure_3.jpeg)

![](_page_26_Picture_4.jpeg)

![](_page_27_Picture_1.jpeg)

![](_page_27_Picture_2.jpeg)

![](_page_27_Figure_3.jpeg)

![](_page_27_Picture_4.jpeg)

![](_page_28_Picture_1.jpeg)

![](_page_28_Picture_2.jpeg)

![](_page_28_Figure_3.jpeg)

![](_page_28_Figure_4.jpeg)

![](_page_28_Picture_5.jpeg)

![](_page_29_Picture_1.jpeg)

![](_page_29_Picture_2.jpeg)

![](_page_29_Figure_3.jpeg)

![](_page_29_Figure_4.jpeg)

![](_page_29_Picture_5.jpeg)

![](_page_30_Picture_1.jpeg)

![](_page_30_Picture_2.jpeg)

![](_page_30_Figure_3.jpeg)

# EXAMPLES SERVICE SERVI

![](_page_30_Figure_5.jpeg)

![](_page_30_Picture_6.jpeg)

82.38

![](_page_31_Picture_1.jpeg)

![](_page_31_Picture_2.jpeg)

![](_page_31_Figure_3.jpeg)

![](_page_31_Picture_4.jpeg)

![](_page_32_Picture_1.jpeg)

![](_page_32_Picture_2.jpeg)

![](_page_32_Figure_3.jpeg)

![](_page_33_Picture_1.jpeg)

![](_page_33_Picture_2.jpeg)

![](_page_33_Figure_3.jpeg)

![](_page_34_Picture_1.jpeg)

![](_page_34_Picture_2.jpeg)

![](_page_34_Figure_3.jpeg)

![](_page_35_Picture_1.jpeg)

![](_page_35_Picture_2.jpeg)

![](_page_35_Figure_3.jpeg)
# Preliminary Outcomes: ISO-3DFD Case study







# Preliminary Outcomes: ISO-3DFD Case study







# Preliminary Outcomes: ISO-3DFD Case study















### Intel Advisor CARM

(product release version)



### sw4lite (LLNL, US) Proxy version of SW4 (3-D seismic modeling)

- Again bunch of stencils ....
- 6 main hotspots (loops)

### Intel Advisor CARM:

- All loops are bound by both (mem and comp)
- Loops 1 and 2: Mainly limited by L3
- Loop 3: Between DRAM and L3 (some locality)
- Loops 4, 5 and 6: DRAM bound







### Intel Advisor CARM

(product release version)



### sw4lite (LLNL, US) Proxy version of SW4 (3-D seismic modeling)

- Again bunch of stencils ....
- 6 main hotspots (loops)

### Improved interpretation methodology in action:

- **Memory traffic shares**: Additional performance insights (explains the dot position)
- **Performance impact**: Improved optimization hints (decouples the bottlenecks by their importance)









### Intel Advisor CARM

(product release version)



### sw4lite (LLNL, US) bare bone version of SW4 (3-D seismic modeling)

- Again bunch of stencils ....
- 6 main hotspots (loops)

### Improved interpretation methodology in action:

- **Memory traffic shares**: Additional performance insights (explains the dot position)
- **Performance impact**: Improved optimization hints (decouples the bottlenecks by their importance)









### Intel Advisor CARM

(product release version)



### ExaMiniMD (SNL) Molecular Dynamics

- 2 main hotspots (loops)
- Loop 1 (memory bound), Loop 2 (comp/mem)

### Improved interpretation methodology in action:

- **Memory traffic shares**: Additional performance insights (explains the dot position)
- **Performance impact**: Improved optimization hints (decouples the bottlenecks by their importance)











How to detect which genes influence traits or diseases?







How to detect which genes influence traits or diseases?

- Relates genetic markers that are most likely to influence diseases.
- Calculates score for a combination of K genetic markers (order).
- Highest score corresponds to SNP combination that is most likely to influence the trait or disease.







How to detect which genes influence traits or diseases?

			SN₽₀	SNP	SNP <sub>2</sub>	SNP₃	SN₽₄	 SNP <sub>M-1</sub>	Class
	$\left( \right)$	P <sub>0</sub>							
Patients <		P <sub>1</sub>							
		$P_2$							
		$P_3$							
		P_4							
		P <sub>N-1</sub>							

- Relates genetic markers that are most likely to influence traits or diseases.
- Calculates score for a combination of K genetic markers (order).
- Highest score corresponds to SNP combination that is most likely to influence the trait or disease.







How to detect which genes influence traits or diseases?



- Relates genetic markers that are most likely to influence traits or diseases.
- Calculates score for a combination of K genetic markers (order).
- Highest score corresponds to SNP combination that is most likely to influence the trait or disease.







How to detect which genes influence traits or diseases?

# Patients $P_{P_1}$ Patients $P_{P_2}$ Patients $P_{P_2}$ Patients $P_{P_1}$ Patients $P_{P_2}$ Patients $P_$

1 and 2

- Relates genetic markers that are most likely to influence traits or diseases.
- Calculates score for a combination of K genetic markers (order).
- Highest score corresponds to SNP combination that is most likely to influence the trait or disease.



















































62









### Original Implementation Intel Advisor CARM



### Bayesian K2 Score High-Order Epistasis

- 1 main hotspot

### Intel Advisor CARM:

- Loops is bound by both (mem and comp)
- Loop 1: Mainly limited by L3

### Index Calculation + Frequency Table



### **Advisor Instruction Mix**



### Completely dominated by Scalar instructions







### Original Implementation Intel Advisor CARM





### Bayesian K2 Score High-Order Epistasis

- 1 main hotspot

### Intel Advisor CARM:

- Loop is bound by both (mem and comp)
- Mainly limited by L3

- Utilization of 8-bit integers to reduce memory footprint
- Vectorization









64 indexes calculated in parallel

- Utilization of 8-bit integers to reduce memory footprint
- Vectorization







### Optimized Version – Single-Thread Intel Advisor CARM



### Bayesian K2 Score High-Order Epistasis

- 1 main hotspot

### Intel Advisor CARM:

- Loop is Bound by memory
- Between L2 and L3



64 indexes calculated in parallel

- Utilization of 8-bit integers to reduce memory footprint
- Vectorization







### Optimized Version – Single-Thread Intel Advisor CARM



### Bayesian K2 Score High-Order Epistasis

- 1 main hotspot

### Intel Advisor CARM:

- Loop is Bound by memory
- Between L2 and L3



- Utilization of 8-bit integers to reduce memory footprint
- Vectorization







### Optimized Version – 18 Threads Intel Advisor CARM













### Optimized Version – 18 Threads Intel Advisor CARM





# Outline





**RECAP: ORIGINAL ROOFLINE MODEL\*** 

CACHE-AWARE ROOFLINE MODEL

- Performance
- EXTENSIONS: POWER AND ENERGY-EFFICIENCY

APPLICATION CHARACTERIZATION WITH CACHE-AWARE ROOFLINE MODEL

- APPLICATION-DRIVEN CARM
- ISO3DFD CASE STUDY
- PROXY-APPLICATIONS FROM EXASCALE COMPUTING PROJECT
- BAYESIAN K2 SCORE

ON-GOING (FUNDED) PROJECTS





# EUROPEAN PROCESSOR INITIATIVE

# FRAMEWORK PARTNERSHIP AGREEMENT IN EUROPEAN LOW-POWER MICROPROCESSOR TECHNOLOGIES



THIS PROJECT HAS RECEIVED FUNDING FROM THE EUROPEAN UNION'S HORIZON 2020 RESEARCH AND INNOVATION PROGRAMME UNDER GRANT AGREEMENT NO 826647

# How EuroHPC will help to make us stronger





- Developing a new European supercomputing ecosystem: HPC systems, network, software, applications, access through the cloud
- Making HPC resources available to public and private users, including SMEs.
- Stimulating a technology supply industry


### **EPI: VISION**





- High Performance Computing needs for Exascale machines beyond 2022
- Connected mobility & Autonomous Driving computing needs beyond 2023
- Low power CPU needs for Servers and Cloud
- Other markets under exploration (Server and Cloud)



# **European Processor Initiative**



- High Performance General Purpose Processor for HPC
- High-performance RISC-V based accelerator
- Computing platform for autonomous cars
- Will also target the AI, Big Data and other markets in order to be economically sustainable



### **GPP** and common architecture







**TÉCNICO** LISBOA

լյլ

# **EPAC – RISC-V Accelerator**







- EPAC EPI Accelerator
- VPU Vector Processing Unit
- STX Stencil/Tensor accelerator
- VRP VaRiable Precision co-processor



### **EPI PARTNERS**









#### **Intel Research Grant**





(intel)

IJÎ

• Started in 2018 and extended for 2019



**TÉCNICO** LISBOA

# Conclusions



SEVERAL CACHE-AWARE ROOFLINE MODELS (EXPERIMENTALLY VERIFIED)

- For several domains: performance, power and energy

# **APPLICATION-DRIVEN CARM**

- Not only considers different micro-architectures (GPU and NUMA) but also considers application requirements: ISA extensions, load/store ratio
- New visual aids for improved characterization: memory traffic and performance impact metrics
- Improved characterization of 3 applications representative of real-world scenarios: ISO-3DFD, SW4Lite and Bayesean K2 score

**ON-GOING (FUNDED) PROJECTS** 

– EPI and Intel: show the practical interest of the work

# **FUTURE WORK**

- INCLUDE ADDITIONAL INFORMATION IN APPLICATION DRIVEN CARM (INTEGERS,
  - CONVERSIONS...)



# **Questions?**

# Thank you!

#### Further readings:

A. Ilic, F. Pratas, and L. Sousa, **"Beyond the Roofline: Cache-aware Power and Energy-Efficiency** *Modeling for Multi-cores"*, *IEEE Transactions on Computers (2017)* 

A. Ilic, F. Pratas, and L. Sousa, *"Cache-aware Roofline model: Upgrading the loft",* IEEE Computer Architecture Letters, CAL (2014)

D. Marques, H. Duarte, A. Ilic, L. Sousa, R. Belenov, P. Thierry and Z. Matveev, "*Performance Analysis with Cache-Aware Roofline Model in Intel Advisor*", *tutorial paper*, HPCS (2017)

A. Lopes, F. Pratas, L. Sousa and A. Ilic, *"Exploring GPU performance, power and energy-efficiency bounds with CARM",* ISPASS (2017)

N. Denoyelle, B. Goglin, A. Ilic, E. Jeannot and L. Sousa, *"Modeling Non-Uniform Memory Access on Large Compute Nodes with the Cache-Aware Roofline Model",* IEEE TPDS (2018)

A. Ilic, F. Pratas, and L. Sousa, *"CARM: Cache-Aware Performance, Power and Energy-Efficiency Roofline Modeling",* Intel CATC (2015)

L. Taniça, A. Ilic, P. Tomás, and L. Sousa, *"SchedMon: A Performance and Energy Monitoring Tool for Modern Multi-cores",* MuCoCoS/Euro-Par (2014)



### **Cache-Aware Roofline Model**









# Application Characterization: DRAM





<sup>1)</sup> Ilic A., Pratas F., Sousa L., "Cache-aware Roofline Model: Upgrading the loft", IEEE CAL (2014)

<sup>2)</sup> Williams S., et. al., "Roofline: An insightful visual performance model for multicore architectures", Communications of the ACM (2009)

ļİ

# Application Characterization: L1 cache





<sup>1)</sup> Ilic A., Pratas F., Sousa L., "Cache-aware Roofline Model: Upgrading the loft", IEEE CAL (2014)

<sup>2)</sup> Williams S., et. al., "Roofline: An insightful visual performance model for multicore architectures", Communications of the ACM (2009)

ļİ

# Application Characterization: CARM and ORM



