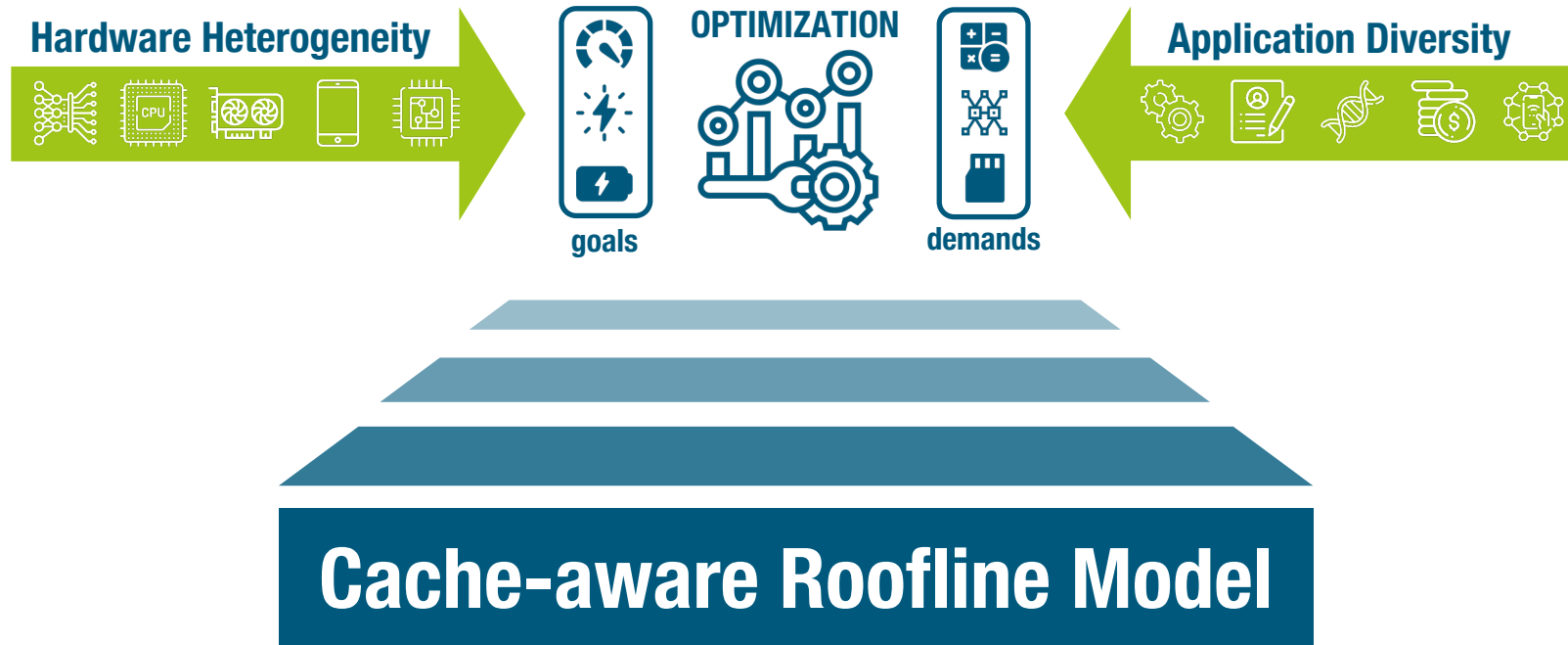# Cache-aware Roofline Model:
## Performance, Power and Energy-efficiency

**Aleksandar Ilic**

Diogo Marques, Frederico Pratas, Leonel Sousa
Rafael Campos, Ricardo Nobre, Sergio Santander-Jiménez

inesc id
lisboa **20** YEARS
DEFINING TECHNOLOGY

Hardware Heterogeneity

OPTIMIZATION

Application Diversity

goals

demands

Cache-aware Roofline Model

PERFORMANCE   POWER   ENERGY-EFFICIENCY   CASE-STUDY

# Cache-aware Roofline Model: Outline

**PERFORMANCE**   **POWER**   **ENERGY-EFFICIENCY**   **CASE-STUDY**
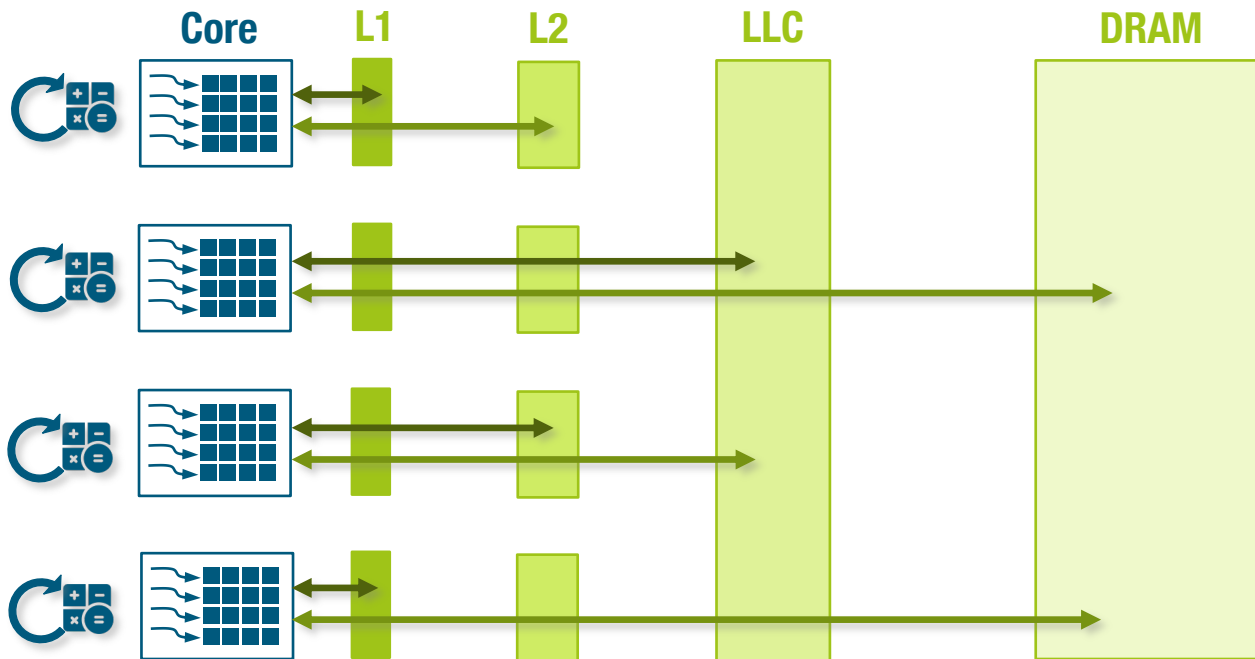
# Cache-aware Roofline Model

A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
D. Marques, A. Ilic, Z. Matveev and L. Sousa, "Application-driven Cache-Aware Roofline Model", Elsevier FGCS (2020)

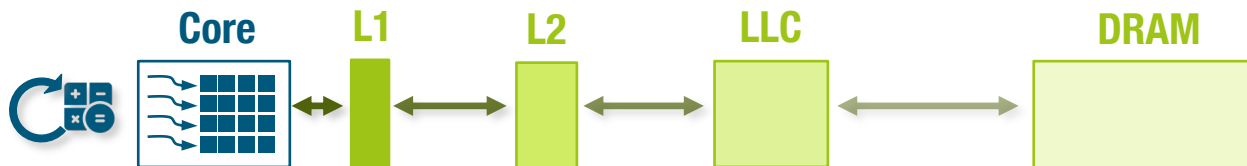# Roofline in a nutshell



Communication overlapped with computation
Max performance capped by peak compute throughput or available bandwidth (processor's view)

# What is bandwidth?



**Cache-aware Roofline Model (CARM)[1]: Bandwidth as seen by the core**
- **Obtained via micro-benchmarking**



**Original Roofline Model (ORM)[2]: Bandwidth between memory levels**
 - **Can be obtained from data-sheets**

[1] A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
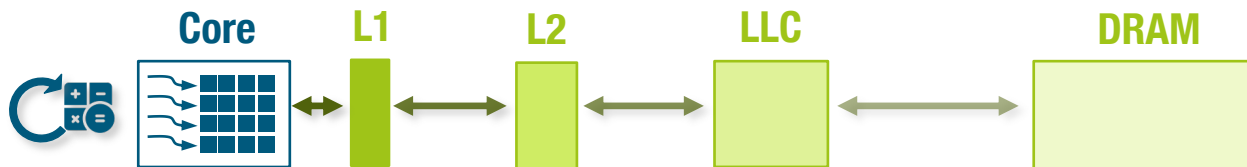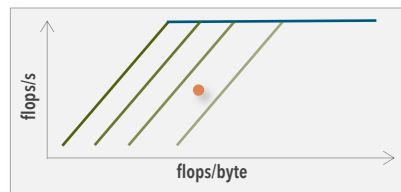[2] S. Williams, A. Waterman, D. Patterson, "Roofline: An Insightful Visual Performance Model for Multicore Architectures", Commun. ACM (2009)

# Implications …

**Core**  **L1**  **L2**  **LLC**  **DRAM**

## Cache-aware Roofline Model[1]
- One model, one arithmetic intensity
- One application "point"

**Core**  **L1**  **L2**  **LLC**  **DRAM**

## Original Roofline Model[2]
- Several models, several intensities
- Several application "points"

[1] A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)

[2] S. Williams, A. Waterman, D. Patterson, "Roofline: An Insightful Visual Performance Model for Multicore Architectures", Commun. ACM (2009)
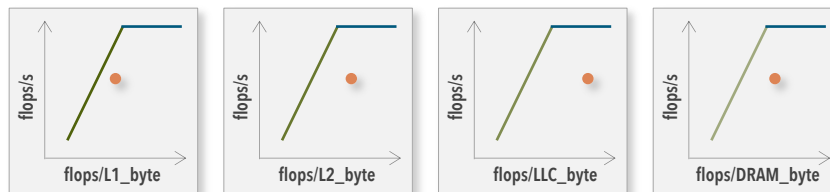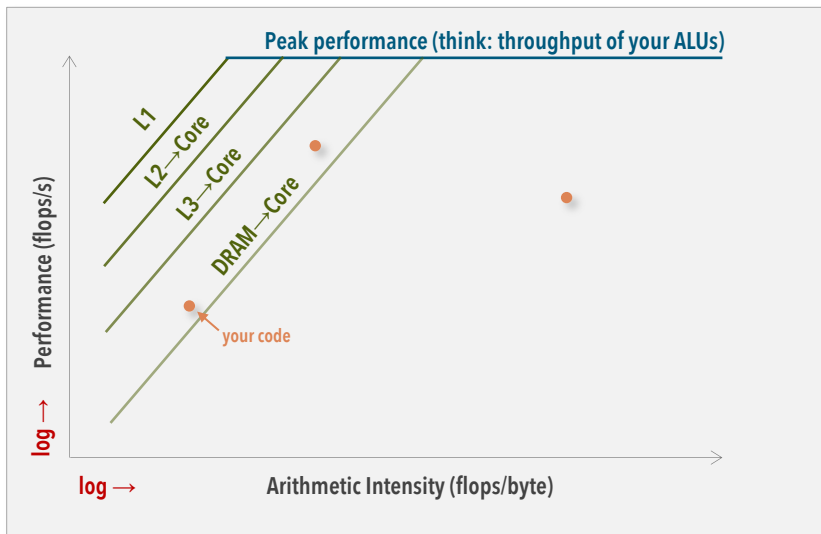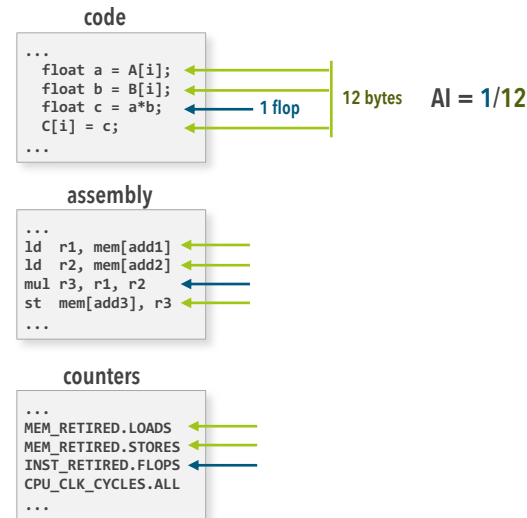
# Implications … bring cool features



Peak performance (think: throughput of your ALUs)

L1
L2→Core
L3→Core
DRAM→Core

Performance (flops/s)

log ↑

your code

log →

Arithmetic Intensity (flops/byte)
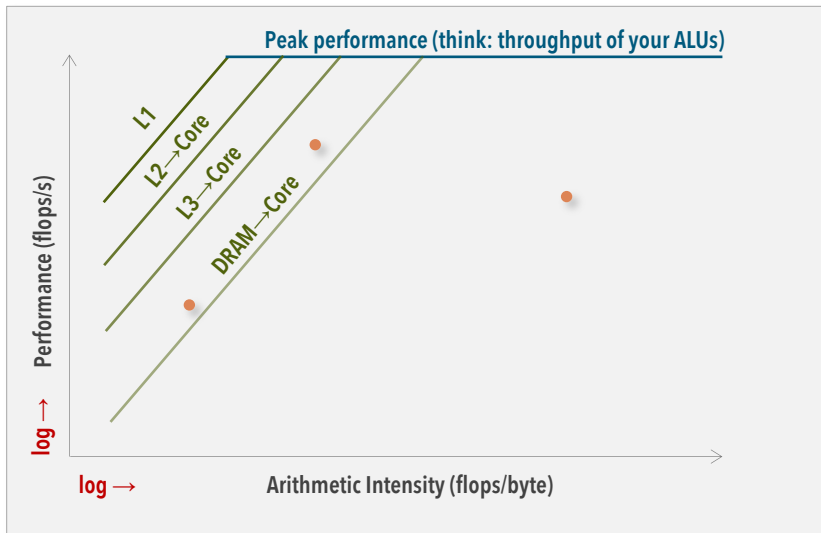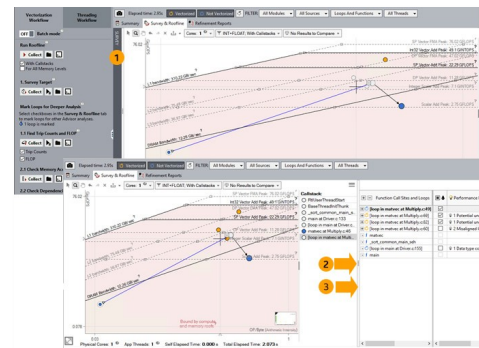
## Cache-aware Roofline Model
- **Shows absolute <u>architecture</u> maximums\***
(You can't break them! Can your application exploit them?)

## How to "plot" my code?
- CARM arithmetic intensity is exactly what you expect it to be!

**code**
```
...
  float a = A[i];
  float b = B[i];
  float c = a*b;
  C[i] = c;
...
```
1 flop    12 bytes    **AI = 1/12**

**assembly**
```
...
ld  r1, mem[add1]
ld  r2, mem[add2]
mul r3, r1, r2
st  mem[add3], r3
...
```

**counters**
```
...
MEM_RETIRED.LOADS
MEM_RETIRED.STORES
INST_RETIRED.FLOPS
CPU_CLK_CYCLES.ALL
...
```

**\* We will relax this requirement in the next part of the talk**

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

# Implications … bring cool features



Peak performance (think: throughput of your ALUs)

L1
L2→Core
L3→Core
DRAM→Core

Performance (flops/s)

log →

log →    Arithmetic Intensity (flops/byte)

## Cache-aware Roofline Model
 - Shows absolute underline{architecture} maximums
(You can't break them! Can your application exploit them?)
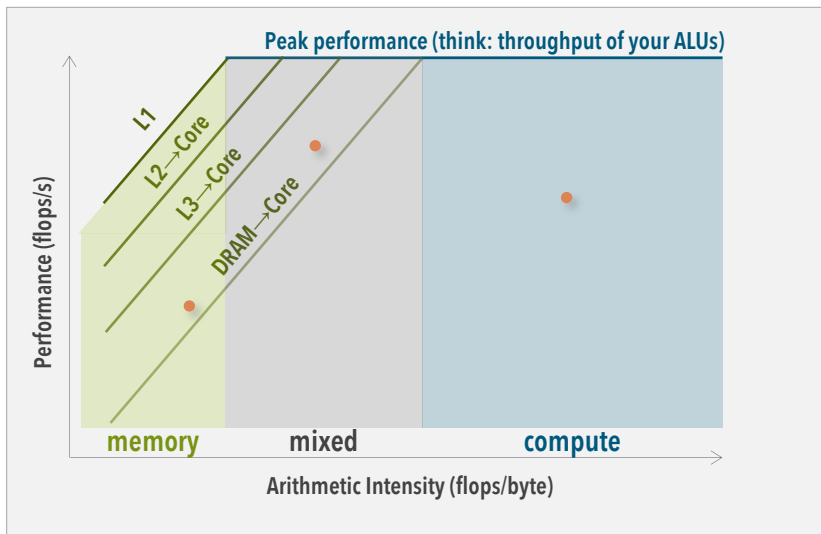
## How to "plot" my code?
 - CARM arithmetic intensity is exactly what you expect it to be!

## Intel Advisor Roofline feature
 - CARM is there since 2017



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

# Implications … bring cool features



Peak performance (think: throughput of your ALUs)

L1
L2→Core
L3→Core
DRAM→Core

Performance (flops/s)

Arithmetic Intensity (flops/byte)

memory · mixed · compute

**memory bound**
(improve access pattern, use of caches)

**mixed**
(all kinds of everything)

**compute bound**
(vectorize, parallelize…)

## Cache-aware Roofline Model
- **Shows absolute <u>architecture</u> maximums**
(You can't break them! Can your application exploit them?)

## How to "plot" my code?
- CARM arithmetic intensity is exactly what you expect it to be!

## How to use CARM?

① **Detect the boundness region**
- What are my expected maximums?
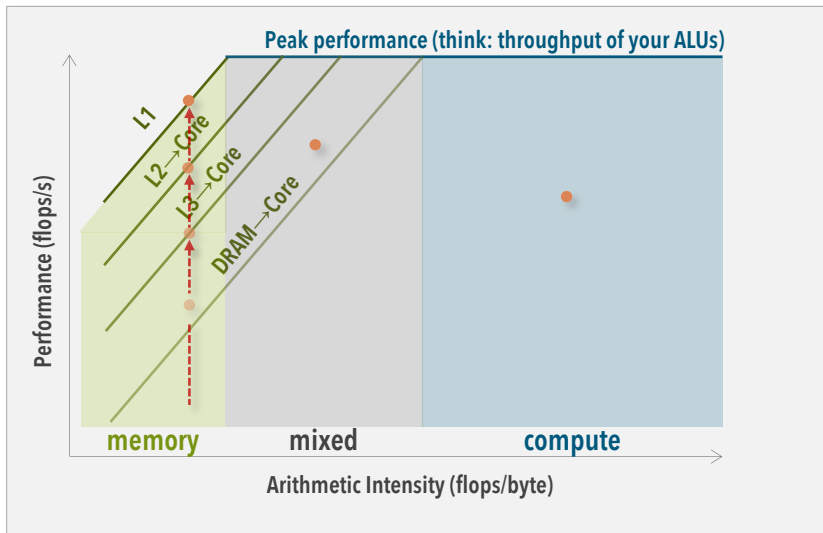- Provides first optimization hints

② **Draw an imaginary vertical line**
- What are my main bottlenecks? (observe intersected lines)
- Focus your optimization (aim at surpassing the line above)

③ **Optimize your code: Break above roofs!**
- You should move up (as your performance improves)
- Unless you restructure the code, or your compiler decides so…

A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
D. Marques, et.al., "Performance analysis with Cache-aware Roofline Model in Intel Advisor", HPCS (2017)

# Implications … bring cool features



**memory bound**
(improve access pattern, use of caches)

**mixed**
(all kinds of everything)

**compute bound**
(vectorize, parallelize…)

## Cache-aware Roofline Model
- Shows absolute <u>architecture</u> maximums
(You can't break them! Can your application exploit them?)

## How to "plot" my code?
- CARM arithmetic intensity is exactly what you expect it to be!

## How to use CARM?

① **Detect the boundness region**
- What are my expected maximums?
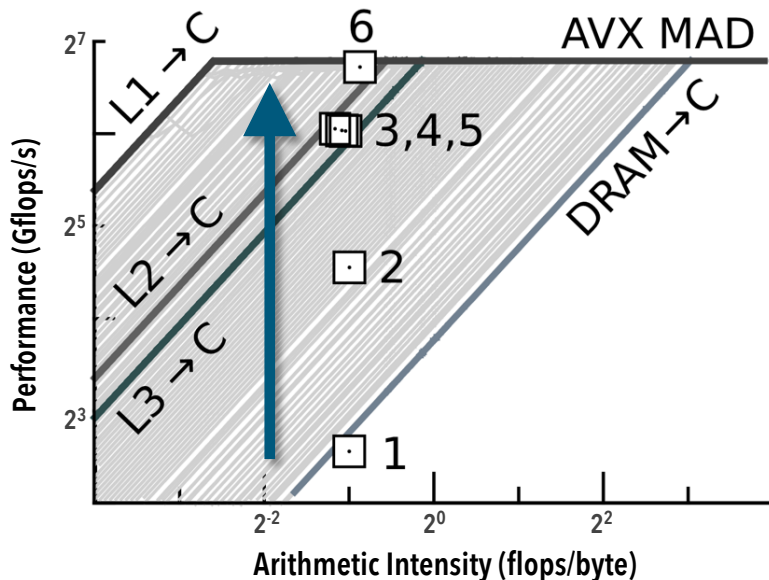- Provides first optimization hints

② **Draw an imaginary vertical line**
- What are my main bottlenecks? (observe intersected lines)
- Focus your optimization (aim at surpassing the line above)

③ **Optimize your code: Break above roofs!**
- You should move up (as your performance improves)
- Unless you restructure the code, or your compiler decides so…

A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
D. Marques, et.al., "Performance analysis with Cache-aware Roofline Model in Intel Advisor", HPCS (2017)

# Implications … bring cool features



Peak performance (think: throughput of your ALUs)

L1 → Core
L2 → Core
L3 → Core
DRAM → Core

Performance (flops/s)

memory | mixed | compute

Arithmetic Intensity (flops/byte)

**memory bound**
(improve access pattern, use of caches)

**mixed**
(all kinds of everything)

**compute bound**
(vectorize, parallelize…)

## Cache-aware Roofline Model
 - Shows absolute <u>architecture</u> maximums
(You can't break them! Can your application exploit them?)

## How to "plot" my code?
 - CARM arithmetic intensity is exactly what you expect it to be!

## How to use CARM?

① **Detect the boundness region**
   - What are my expected maximums?
   - Provides first optimization hints

② **Draw an imaginary vertical line**
   - What are my main bottlenecks? (observe intersected lines)
   - Focus your optimization (aim at surpassing the line above)

③ **Optimize your code: Break above roofs!**
   - You should move up (as your performance improves)
   - Unless you restructure the code, or your compiler decides so…

A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
D. Marques, et.al., "Performance analysis with Cache-aware Roofline Model in Intel Advisor", HPCS (2017)

# Matrix Multiplication

* A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
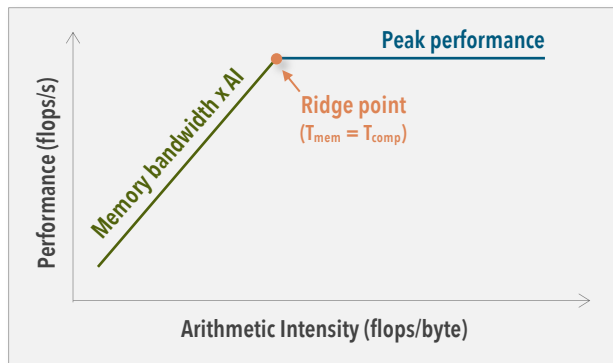* A. Ilic, F. Pratas and L. Sousa, "Beyond the Roofline: Cache-Aware Power and Energy-Efficiency Modeling for Multi-Cores", IEEE Trans. on Computers (2017)

## All codes AVX vectorized!*
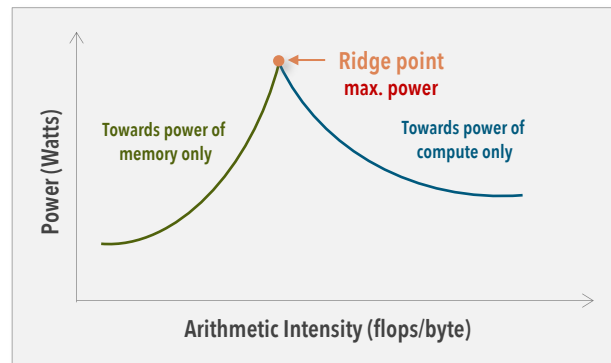
[1] **Basic implementation (row major)**



[2] **Transposed B (improved mem. access)**



[3,4,5] **Cache blocking: L3, L2, L1**



[6] **Intel MKL**

**Variable Architecture Maximums**

THROUGHPUT  BANDWIDTH  CORE COUNT  UTILIZATION

**Diverse Application Characteristics**

INST MIX  AVX/SSE  LD/ST  FP SHARE

# Application-driven CARM

**(scaling rooflines to meet application demands)**

D. Marques, A. Ilic, Z. Matveev and L. Sousa, "Application-driven Cache-Aware Roofline Model", Elsevier FGCS (2020)

# ISO-3DFD: Quite optimized 3D stencil (scalar)

**Absolute CARM**



mixed region
(bound by both memory and compute)

**Application-driven CARM**



memory bound

### CARM characterization cheat-sheet

|  | Absolute | Application-driven |
|---|---|---|
| region | mixed | memory |
| max. perf. | compute (add) | memory (L1) |
| bottleneck | memory/compute | memory |
| optimize | everything | memory (or nothing) |

## Application-driven CARM
- models architecture maximums exploitable by your application
- improves characterization and hints (bottlenecks, optimization)
- provides consistent characterization during optimization process

# ISO-3DFD: Scalar (left) vs. AVX512 (right)

**PERFORMANCE**

**POWER**

**ENERGY-EFFICIENCY**

**CASE-STUDY**

# Cache-aware Roofline Model

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

# CARM: Power Consumption

## Performance CARM
- Contributions of comps and mops <u>overlapped</u> (in time)



## Power CARM
- Contributions of comps and mops <u>superposed</u> (average power)



comps = compute operations
mops = memory operations

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

# CARM: Power Consumption

**Performance CARM**
- Contributions of comps and mops <u>overlapped</u> (in time)

**Power CARM: Cores**
- Contributions of comps and mops <u>superposed</u> (average power)



comps = compute operations
mops = memory operations

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

# Total Power CARM: Defining envelope



Intel 3770K
Ivy Bridge

4 Cores
(AVX MAD)

L2→C    L3→C

Total Power Roofline

L1→C

DRAM→C

Power Package [W]

Arithmetic Intensity [flops/byte]

## CARM for different RAPL domains

### Cores



### Uncore



### Package



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

# Energy-efficiency CARM



Maximum efficiency for infinite arithmetic intensity!

# Matrix Multiplication



Performance CARM

Energy-Efficiency CARM

Power CARM

**All codes AVX vectorized!***

**[1]  Basic implementation (row major)**

A x B = C

**[2]  Transposed B (improved mem. access)**

A x B = C

**[3,4,5]  Cache blocking: L3, L2, L1**

A x B = C

**[6]  Intel MKL**

A x B =

* A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)
* A. Ilic, F. Pratas and L. Sousa, "Beyond the Roofline: Cache-Aware Power and Energy-Efficiency Modeling for Multi-Cores", IEEE Trans. on Computers (2017)

# Cache-aware Roofline Model: Extensions

## CARM-based DVFS analysis



**A** PERFORMANCE CARM

**B** POWER CARM (PACKAGE DOMAIN)

**C** ENERGY-EFFICIENCY CARM (PACKAGE)

## GPU CARM: Performance, Power, DVFS



## NUMA CARM: Multi-socket, KNL

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

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

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

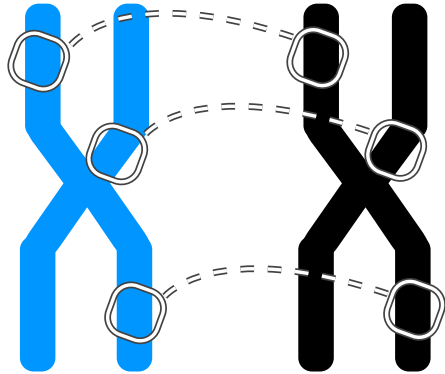**PERFORMANCE**  **POWER**  **ENERGY-EFFICIENCY**  **CASE-STUDY**

# Epistasis Detection: CARM-driven Optimization

R. Nobre, A. Ilic, S. Santander-Jiménez, L. Sousa, "Exploring the Binary Precision Capabilities of Tensor Cores for Epistasis Detection", IPDPS (2020)
R. Campos, D. Marques, S. Santander-Jiménez, L. Sousa, A. Ilic, "Heterogeneous CPU+ iGPU Processing for Efficient Epistasis Detection", EuroPar (2020)
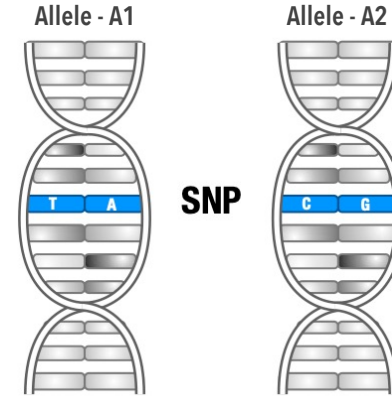
# Epistasis in a nutshell



**SNP**

Some SNP interactions may cause life-threating diseases (e.g., Alzheimer, breast cancer)
Discovering which and how many is important, but challenging task!

# Short Bio Recap: Codifying your genotype



| Genotype | A1 A2 | |
|----------|-------|---|
| 0 | ◆ ◆ | Homozygous Major |
| 1 | ◆ ◆ | Heterozygous |
| 2 | ◆ ◆ | Homozygous Minor |

◆ dominant allele
◆ recessive allele

Allele - A1    Allele - A2

SNP

T  A          C  G

# Binarizing your genotype

| SNP X | P0 | | Genotype | A1 | A2 | |
|-------|-----|-----|----------|-----|-----|-----|
| X0 | 0 | | 0 | ◆ | ◆ | Homozygous Major |
| X1 | 1 | | 1 | ◆ | ◆ | Heterozygous |
| X2 | 0 | | 2 | ◆ | ◆ | Homozygous Minor |
| phenotype | 0 | | | | | |

◆ dominant allele
◆ recessive allele

**Think: Patient 0 (P0) with genotype 1 does <u>not</u> have disease (control)**

Allele - A1    Allele - A2

SNP

T  A      C  G

# Binarizing your genotype

| SNP X | P0 | P1 | | Genotype | A1 | A2 | |
|---|---|---|---|---|---|---|---|
| X0 | 0 | 0 | | 0 | ◆ | ◆ | Homozygous Major |
| X1 | 1 | 0 | | 1 | ◆ | ◆ | Heterozygous |
| X2 | 0 | 1 | | 2 | ◆ | ◆ | Homozygous Minor |
| phenotype | 0 | 1 | | | | | |

◆ dominant allele
◆ recessive allele

**Think: Patient 1 (P1) with genotype 2 has disease (case)**

Allele - A1     Allele - A2

SNP

T   A          C   G

# Dataset structure

| SNP X | P0 | P1 | P2 | P3 | P4 | P5 | ... | PN |
|-------|----|----|----|----|----|----|----|----|
| X0 | 0 | 0 | 0 | 1 | 0 | 0 | ... | 1 |
| X1 | 1 | 0 | 1 | 0 | 0 | 1 | ... | 0 |
| X2 | 0 | 1 | 0 | 0 | 1 | 0 | ... | 0 |
| phenotype | 0 | 1 | 1 | 1 | 0 | 0 | ... | 1 |



**Dataset structure**

**Our dataset: 10 040 SNPs x 104 448 samples**

# 2-way Epistasis Detection: Pair-wise interaction



**Pair-wise interaction:** SNPs (X,Y)

frequency table

ph.type: 0
ph.type: 1

00 01 02 10 11 12 20 21 **22**

X1 Y0 / X2 Y2 / genotype combination

popcnt    popcnt

phenotype — and    phenotype — not — and

X1 **Y0**    **X2 Y2**

**N Samples (Patients)**

M SNPs

SNP X: X0 X1 X2

SNP Y: Y0 Y1 Y2

phenotype

**Dataset structure**
Our dataset: 10 040 SNPs x 104 448 samples

**Search space:** All SNP combinations

(0,1) (0,2) (0,3) (0,4) (...) (0,M-1)
(1,2) (1,3) (1,4) (...) (1,M-1)
(2,3) (2,4) (...) (2,M-1)
(3,4) (...) (3,M-1)
(...)
(...) (M-2,M-1)

**M(M-1)/2 combinations**

Our dataset: 50 395 780 combinations

**Each frequency table evaluated with Bayesian K2 score**
**Epistasis: Minimum K2 score among all combinations!**

30 | R. Nobre, A. Ilic, S. Santander-Jiménez, L. Sousa, "Exploring the Binary Precision Capabilities of Tensor Cores for Epistasis Detection", IPDPS (2020)
R. Campos, D. Marques, S. Santander-Jiménez, L. Sousa, A. Ilic, "Heterogeneous CPU+ iGPU Processing for Efficient Epistasis Detection", EuroPar (2020)

# Cache-aware Roofline Model in Intel® Advisor

# Let's CARMify it!

GALUOPS

Peak: 9.13 GALUOPS

L1 Bandwidth: 59.12 GB/sec
L2 Bandwidth: 23.3 GB/sec
L3 Bandwidth: 15.08 GB/sec
DRAM Bandwidth: 4.53 GB/sec

**mixed region**
**(not easy to optimize)**
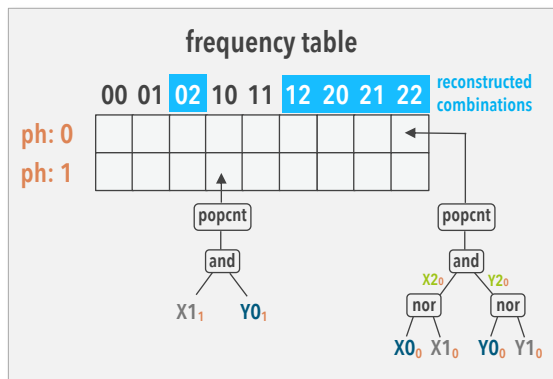
ALUOP/Byte (Arithmetic Intensity)

**Mixed region, but seems memory bound**
**Let's be smart: Restructure our algorithm!**

# Increase arithmetic intensity
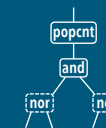
**Pair-wise interaction:** SNPs (X,Y)



frequency table

| | 00 | 01 | 02 | 10 | 11 | 12 | 20 | 21 | 22 | reconstructed combinations |

ph: 0
ph: 1

popcnt
and
$X1_1$  $Y0_1$

popcnt
and
$X2_0$  $Y2_0$
nor    nor
$X0_0$ $X1_0$  $Y0_0$ $Y1_0$



Controls     Cases

M SNPs

$X0_0$
$X1_0$          $X0_1$
$X1_1$   **SNP X**

. . .        . . .

ph.type: 0     ph.type: 1

**"New" Dataset structure**
**(removed: phenotype and genotype 2)**

**Reducing memory transfers!**
**Boosting our arithmetic intensity!**


Three Genotypes + Phenotype
popcnt
not — and


Two Genotypes, No Phenotype
popcnt
and
nor    nor

 R. Nobre, A. Ilic, S. Santander-Jiménez, L. Sousa, "Exploring the Binary Precision Capabilities of Tensor Cores for Epistasis Detection", IPDPS (2020)

# Let's CARMify it (again)!



**Wait! Being smart decreases performance!
How come?!**

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, single core

# Let's CARMify it (again)!



Wait! Being smart decreases performance!
How come?!

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, single core

## Three Genotypes + Phenotype

## Two Genotypes, No Phenotype

# Let's continue optimizing…



**CARM and perf. decrease may suggest memory issues!**
**Let's "tile" our dataset for caches!**

**Three Genotypes + Phenotype**

**Two Genotypes, No Phenotype**

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, single core

# Improvements, at last!!!

**Tiling worked!**
**We now have both: performance increase and speedup!**

# Improvements, at last!!!



**Mixed region, but close to "compute" roof!
Let's vectorize!**

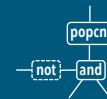**Three Genotypes + Phenotype**

**Two Genotypes, No Phenotype**

**Cache tiling**

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, single core

# CARM in action …



**Let's multi-thread it!**

**Three Genotypes + Phenotype**
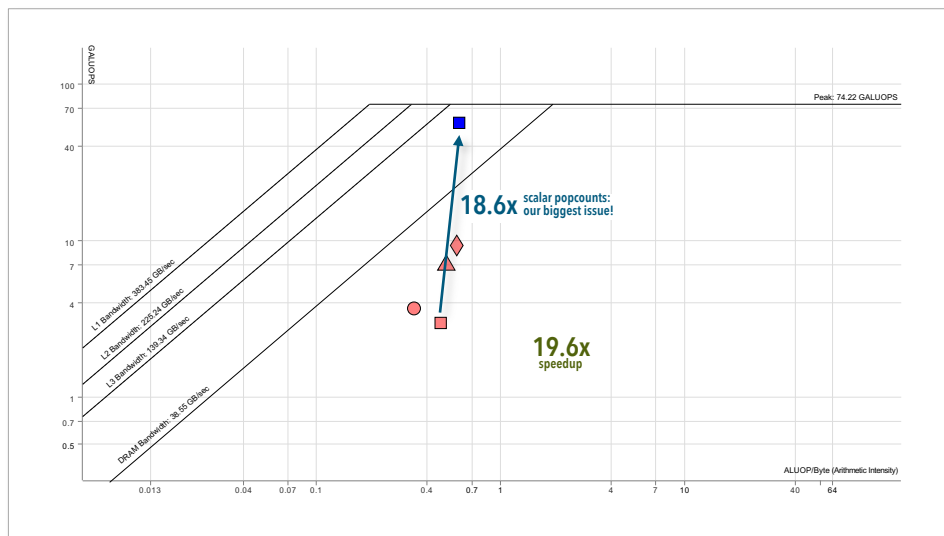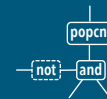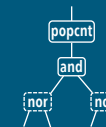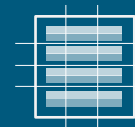
**Two Genotypes, No Phenotype**

**Cache tiling**

**AVX2 Vectorization**

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, single core

Results obtained with a "special version" of Intel® Advisor | Platform: Intel® i7-8700K (3.7GHz) with HT/Prefetching/TurboBoost disabled, six cores

# Epistasis Detection on Intel CPU+iGPU



Best Performance

Best Power and Energy-Efficiency

R. Campos, D. Marques, S. Santander-Jiménez, L. Sousa, A. Ilic, "Heterogeneous CPU+ iGPU Processing for Efficient Epistasis Detection", EuroPar (2020)

**PERFORMANCE**

**POWER**

**ENERGY-EFFICIENCY**

**CASE-STUDY**

# Cache-aware Roofline Model: Conclusions

D. Marques, A. Ilic, Z. Matveev and L. Sousa, "Application-driven Cache-Aware Roofline Model", Elsevier FGCS (2020)

R. Nobre, A. Ilic, S. Santander-Jiménez, L. Sousa, "Exploring the Binary Precision Capabilities of Tensor Cores for Epistasis Detection", IPDPS (2020)

R. Campos, D. Marques, S. Santander-Jiménez, L. Sousa, A. Ilic, "Heterogeneous CPU+ iGPU Processing for Efficient Epistasis Detection", EuroPar (2020)

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

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

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