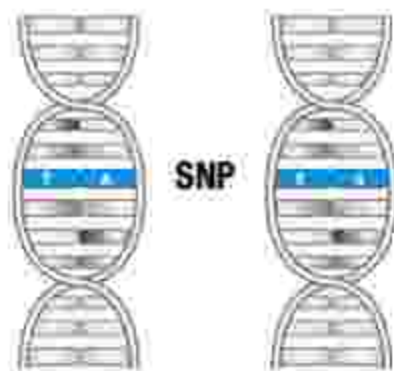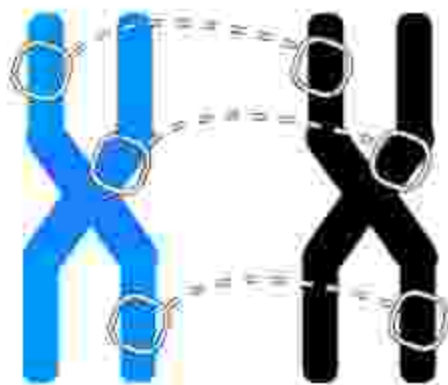# Accelerating epistasis detection on Intel CPUs and discrete GPUs with Intel® Advisor

**Aleksandar Ilic, Diogo Marques, Rafael Campos and Zakhar Matveev**

inesc id
lisboa
20 YEARS
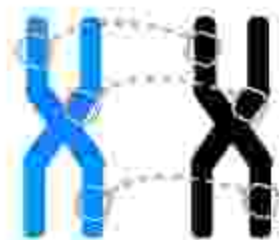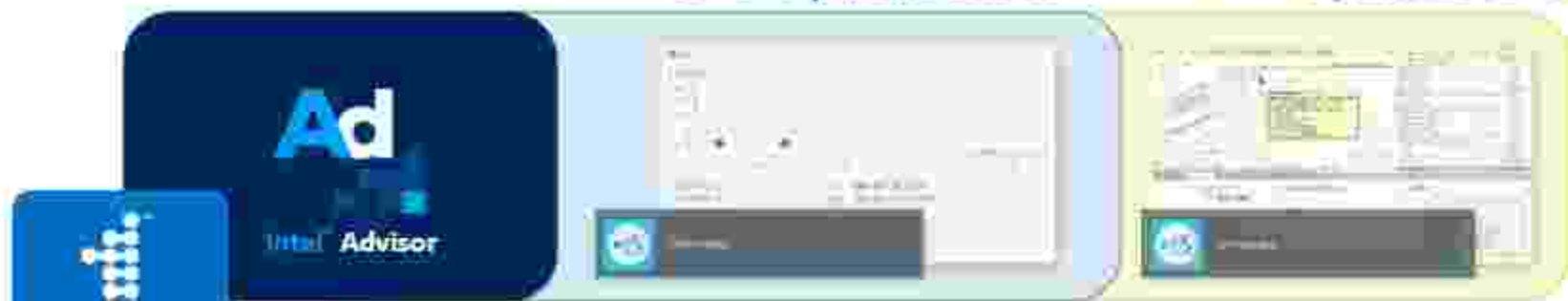DEFINING TECHNOLOGY

inesc-id.pt

# 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!

PERFORMANCE   POWER   ENERGY-EFFICIENCY   CASE-STUDY

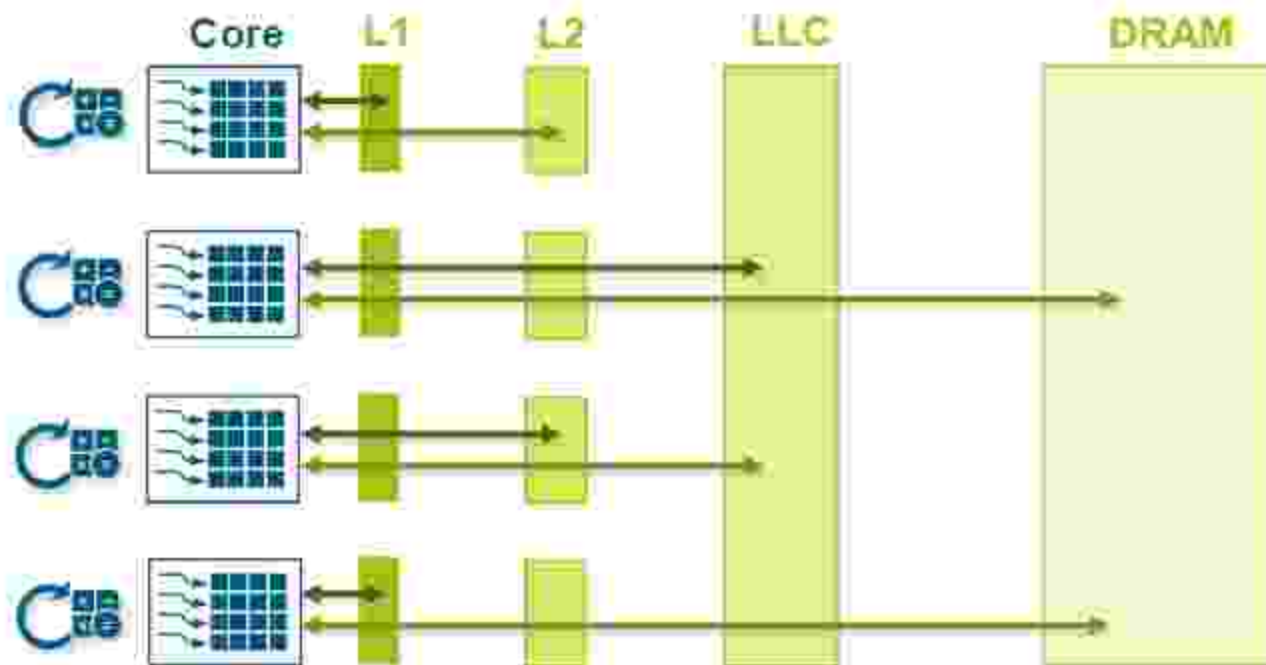# Outline

PERFORMANCE   POWER   ENERGY-EFFICIENCY

# 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



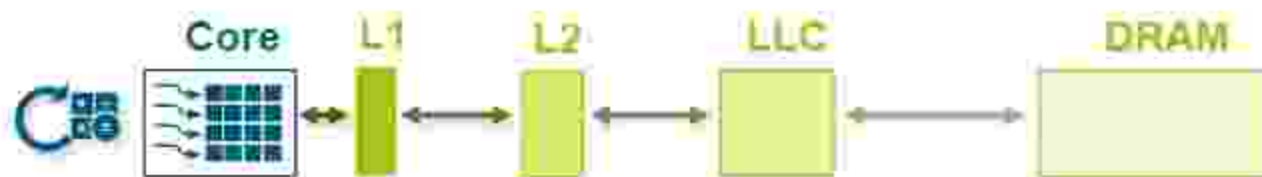Communication overlapped with computation
Max performance capped by peak compute throughput or available bandwidth (processor's view)

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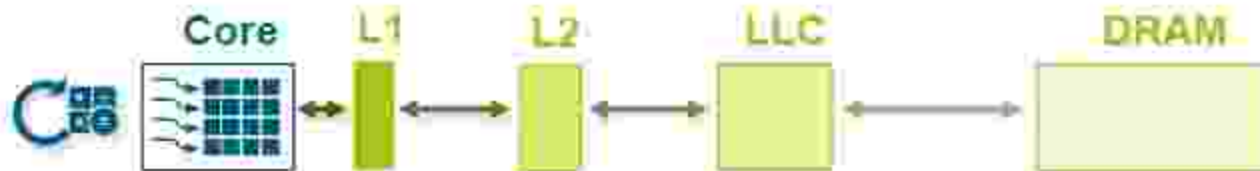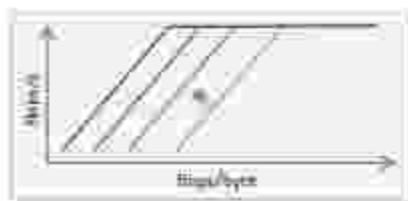[2] S. Williams, A. Waterman, D. Patterson, "Roofline: An Insightful Visual Performance Model for Multicore Architectures", Commun. ACM (2009)

# Implications ...



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

**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).
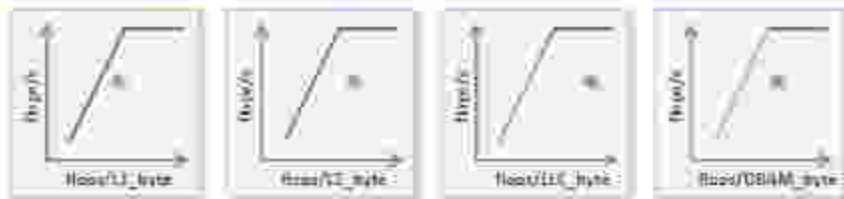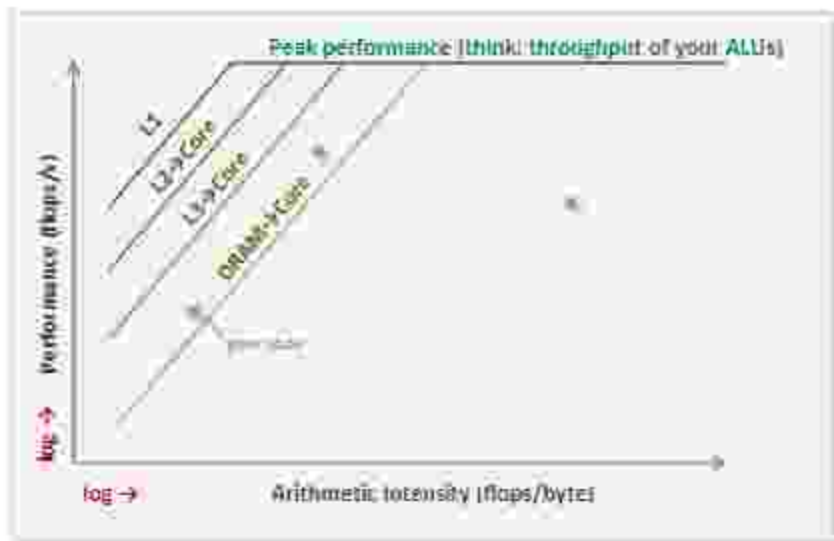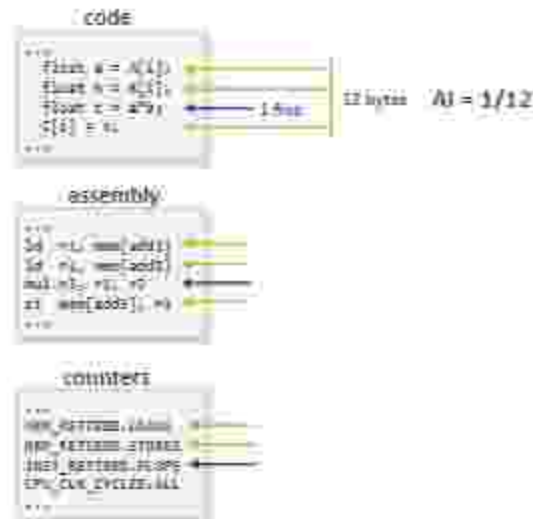
# Implications ... bring cool features
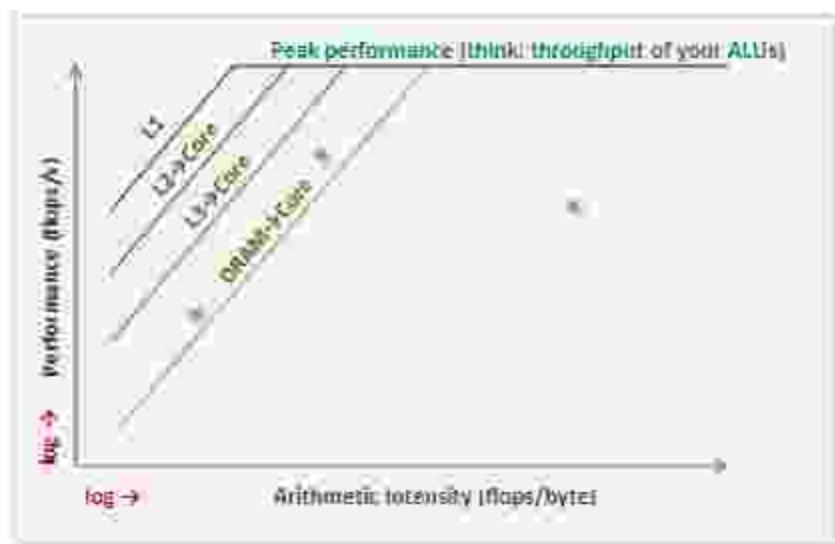


## Cache-aware Roofline Model
- Shows absolute architecture maximums*

(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!

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)

*We relaxed this requirement in our FGCS paper (2020)

# 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 architecture maximums
(You can't break them! Can your application exploit them?)
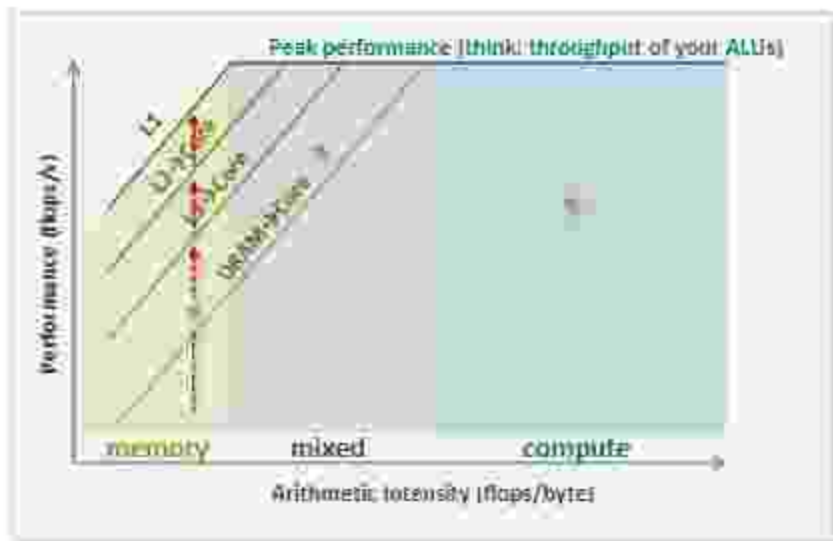
## How to "plot" my code?
- CARM arithmetic intensity is exactly what you expect it to be!

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

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

memory bound
(improve access pattern, use of caches)

mixed
(all kinds of everything)

compute bound
(vectorize, parallelize...)

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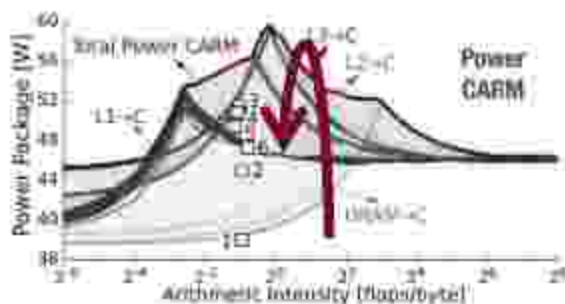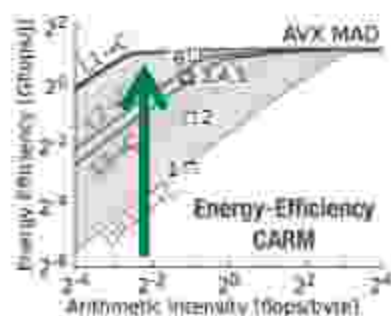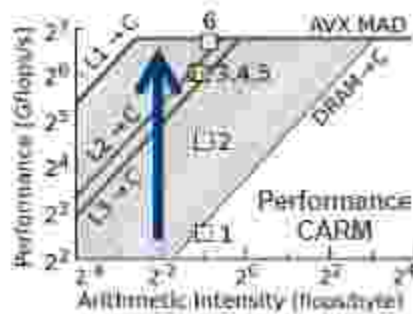(You can't break them! Can your application exploit them?)

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

## 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

Performance CARM



Energy-Efficiency CARM



Power CARM

[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. 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

# Cache-aware Roofline Model: Extensions

## CARM-based DVFS analysis



## 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
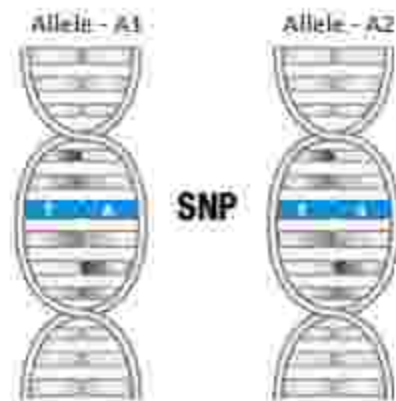
# Epistasis Detection

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)

# Binarizing your genotype
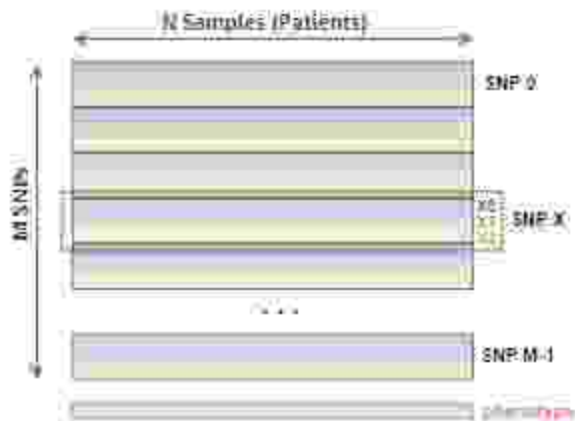


Think: Patient 1 (P1) with genotype 2 has disease (case)

# Dataset structure

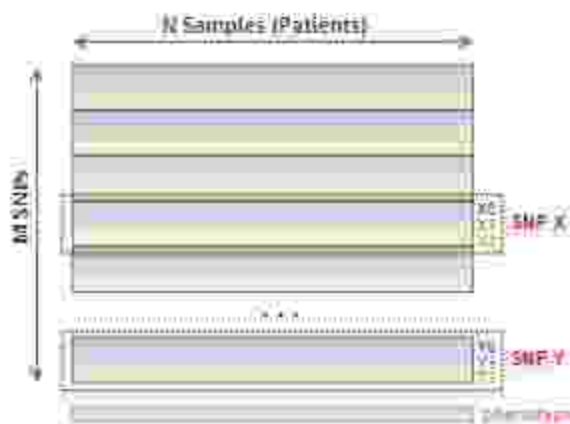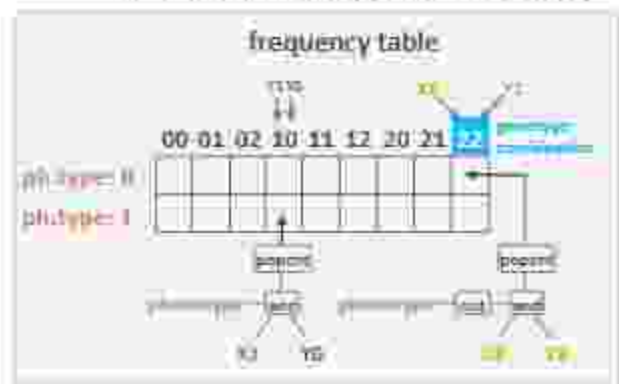| SNP X | P0 | P1 | P2 | P3 | P4 | P5 | ... | PN |
|-------|----|----|----|----|----|----|-----|----|
| X0 | 0 | 0 | 0 | 1 | 0 | 0 | ... | 1 |
| X1 | 1 | 0 | 1 | 0 | 0 | 1 | ... | 0 |
| X2 | 0 | 1 | 0 | 0 | 1 | 0 | ... | 0 |
| phenotype | 0 | 1 | 1 | 1 | 0 | 0 | ... | 1 |



**Dataset structure**

Our dataset: 10 048 SNPs x 104 448 samples

# 2-way Epistasis Detection: Pair-wise interaction



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

frequency table

**Dataset structure**
Our dataset: 10 048 SNPs x 104 448 samples

**Search space:** All SNP combinations

$M(M-1)/2$ combinations

Our dataset: 50 476 128 combinations

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

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 Detection: GPU Optimization**

# DPC++ setup

*Host code*



Create USM buffers

Transfer data from host to device

Launch kernel to execute on device



**The kernel is launched for M*M work-items**

**Work-items without a valid combination will not do work**

# DPC++ setup

*Device code*



Kernel defined as lambda function

Obtain SMPs from work-item index



**The kernel is launched for M*M work-items**

**Work-items without a valid combination will not do work**

# DPC++ implementation

**Pair-wise interaction:** SNPs (i,j)



*Device code*

The frequency table is filled for SNPs i and j by going through all samples, disposed in columns

# GPU Roofline in Advisor

**Summary View:**

# GPU Roofline in Advisor

**Summary View:**

# GPU Roofline in Advisor

**GPU Roofline view:**

Three Genotypes + Phenotype

# Advisor in action...

Three Genotypes + Phenotype

Performance needs improvement!
The application is memory bound – let's restructure

# Restructuring our algorithm...

**Pair-wise interaction:** SNPs (i,j)



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

**Three Genotypes + Phenotype**

**Two Genotypes, No Phenotype**

## Reducing memory transfers!

R. Nobre, A. Ilic, S. Santander-Jiménez, L. Sousa: "Exploring the Binary Precision Capabilities of Tensor Cores for Epistasis Detection", IPDPS (2022)

# DPC++ implementation

**Pair-wise interaction:** SNPs (i,j)



frequency table

Three Genotypes + Phenotype

Two Genotypes, No Phenotype

*Device code*

The frequency table is filled for SNPs i and j is filled separately for each phenotype

# Advisor in action...



**Changing the algorithm moved us to the left!
Performance is worse!**

## GPU Optimization

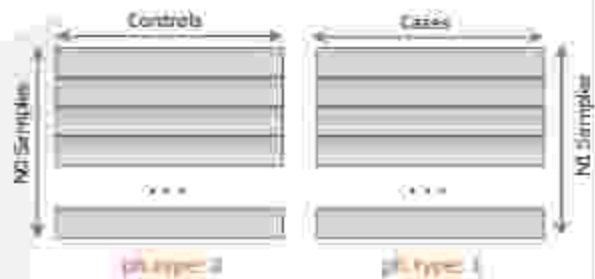### Three Genotypes + Phenotype



### Two Genotypes, No Phenotype

# Let's continue optimizing...

**GPU Optimization**

- Three Genotypes + Phenotype
- Two Genotypes, No Phenotype
- Transposed Dataset

**Transposed Dataset**
(inverted SNPs and samples)

Improving memory accesses by dataset transposition
Coalescing of memory accesses by the work-items

# Advisor in action...



Memory accesses are now more efficient!
We have now "block loads" versus a lot of gathers!

## GPU Optimization

**Three Genotypes + Phenotype**



**Two Genotypes, No Phenotype**



**Transposed Dataset**



1.7x

1.7x

# Advisor in action...

**This allows to perform disease detection 1.7x faster!**

**How can we further improve this result?**

## GPU Optimization

### Three Genotypes + Phenotype

### Two Genotypes, No Phenotype

### Transposed Dataset

# Data set tiling

**GPU Optimization**

Tiling our dataset to squeeze the maximums!
Using a tile size of $B_{SNP}$ we can maintain constant access stride

PERFORMANCE     POWER     ENERGY-EFFICIENCY     CASE-STUDY

# Conclusions

# Conclusions



The insights of Intel Advisor can allow to vastly improve performance

This allows to perform epistasis detection 3.6x faster!

These techniques can be applied to applications in CPU or GPU

# Acknowledgements

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, 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)

**Thank you!**

inesc id lisboa
20 YEARS
DEFINING TECHNOLOGY

**FCT** Fundação
para a Ciência
e a Tecnologia

**EUROPEAN UNION**
European Regional Development Fund