

# Rooflines ...

**Aleksandar Ilic and Alexandre Rodrigues**

Leonel Sousa, José Morgado, Diogo Matos,  
Diogo Marques, Frederico Pratas



**EuroHPC**  
Joint Undertaking



Fundação  
para a Ciência  
e a Tecnologia



# Cache-aware Roofline Model



SOME THEORY



SOME EXAMPLE



SOME TOOLS



SOME ACTION



# Cache-aware Roofline Model: Outline



SOME THEORY



SOME EXAMPLE



SOME TOOLS



SOME ACTION

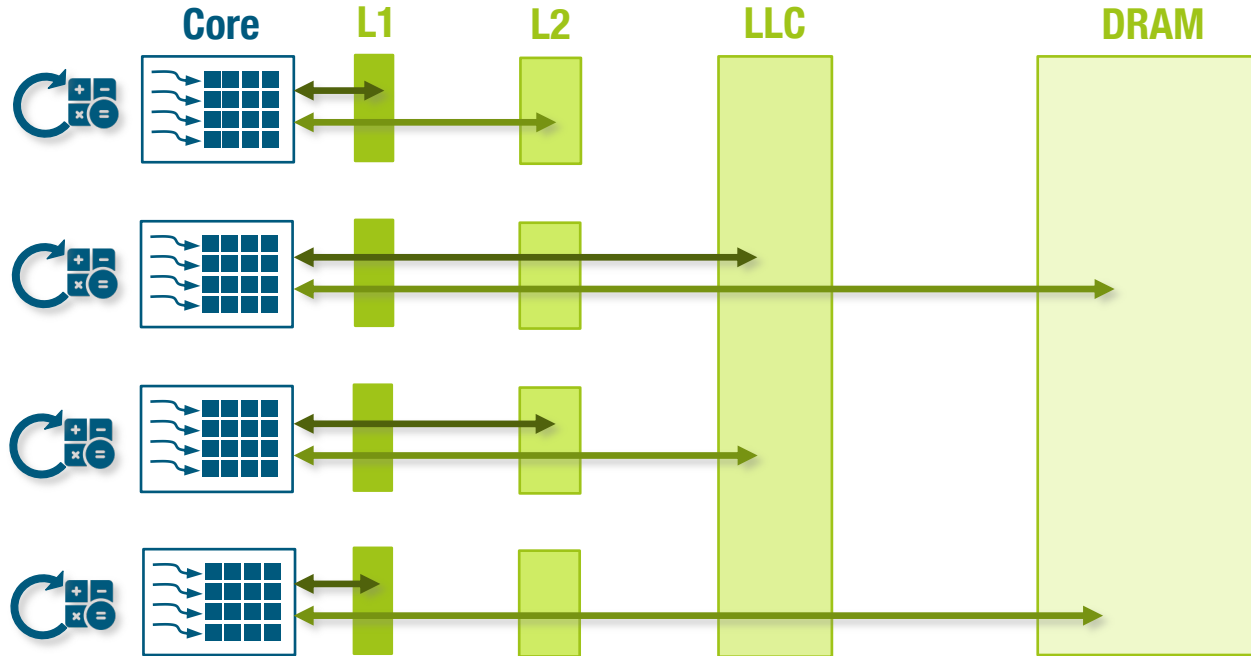


# Cache-aware Roofline Model

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

J. Morgado, L. Sousa and A. Ilic, "CARM Tool: Cache-Aware Roofline Model Automatic Benchmarking and Application Analysis", IISWC (2024)

# Roofline in a nutshell



Communication overlapped with computation

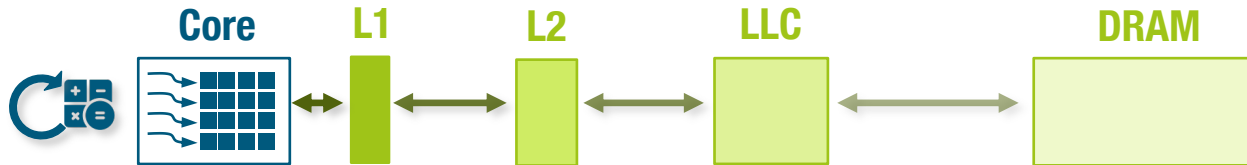
Max performance capped by **peak compute throughput** or **available bandwidth** (processor's view)

# What is bandwidth?



## Cache-aware Roofline Model (CARM)<sup>1</sup>: Bandwidth as seen by the core

- Obtained via micro-benchmarking



## Original Roofline Model (ORM)<sup>2</sup>: Bandwidth between memory levels

- Can be obtained from data-sheets

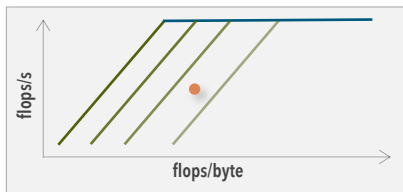
<sup>1</sup> A. Ilic, F. Pratas and L. Sousa, "Cache-aware Roofline Model: Upgrading the Loft", IEEE Computer Architecture Letters (2014)

# Implications ...



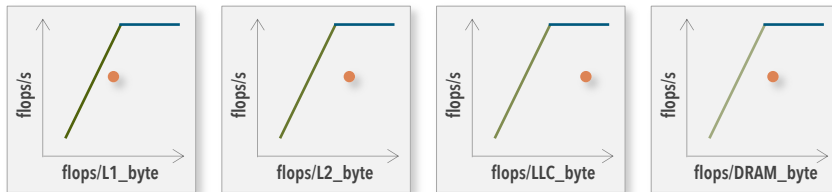
## Cache-aware Roofline Model<sup>1</sup>

- One model, one arithmetic intensity
- One application "point"



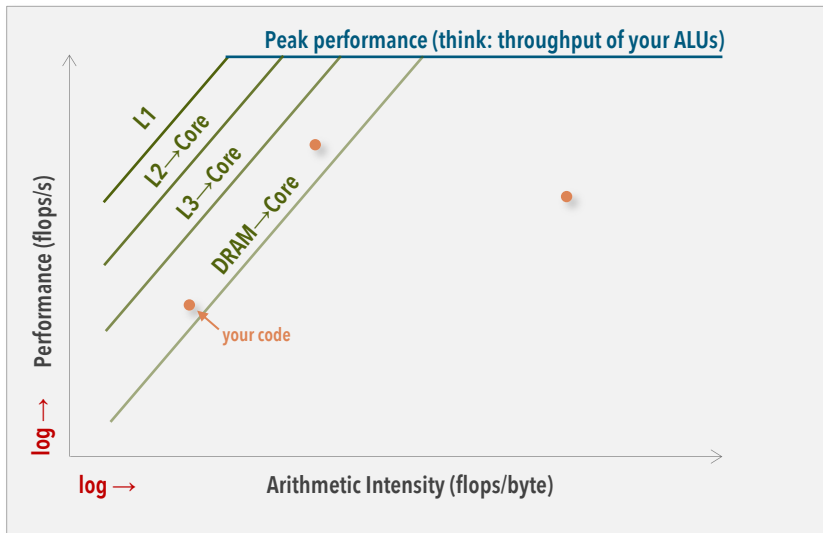
## Original Roofline Model<sup>2</sup>

- Several models, several intensities
- Several application "points"



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

# Implications ... bring cool features

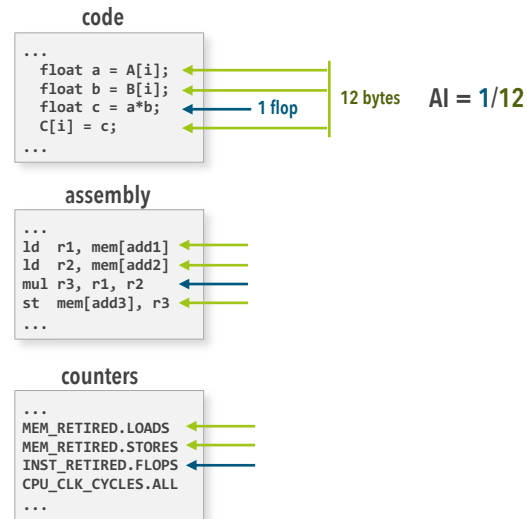


## Cache-aware Roofline Model

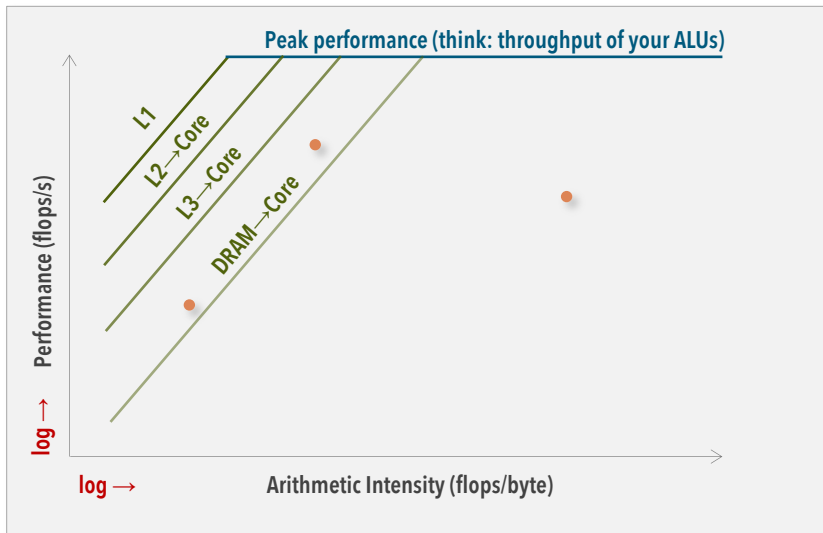
- Shows absolute architecture maximums (You can't break them! Can your application exploit them?)

## How to “plot” my code?

- CARM arithmetic intensity is exactly what you expect it to be!



# 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!

## Intel Advisor Roofline feature

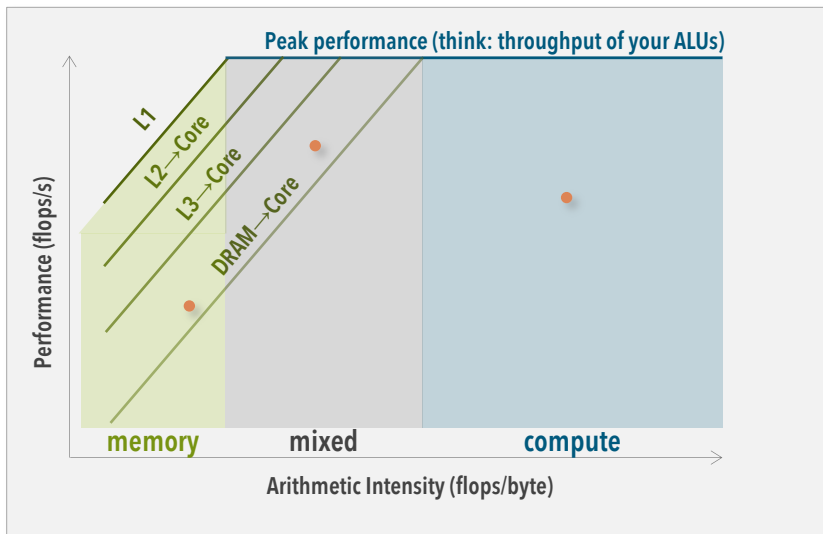
- CARM is there since 2017

## CARM Tool (in-house)

- Open-source at Champ Hub @ Github
- Support all major CPUs and GPUs
- PMU- and DBI-based application profiling



# 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 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

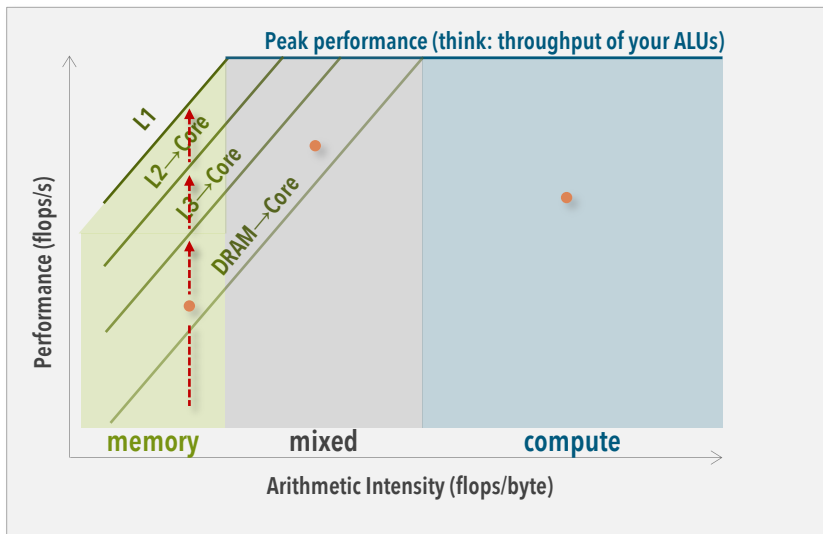
### ② Draw an imaginary vertical line

- 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...

# 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 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

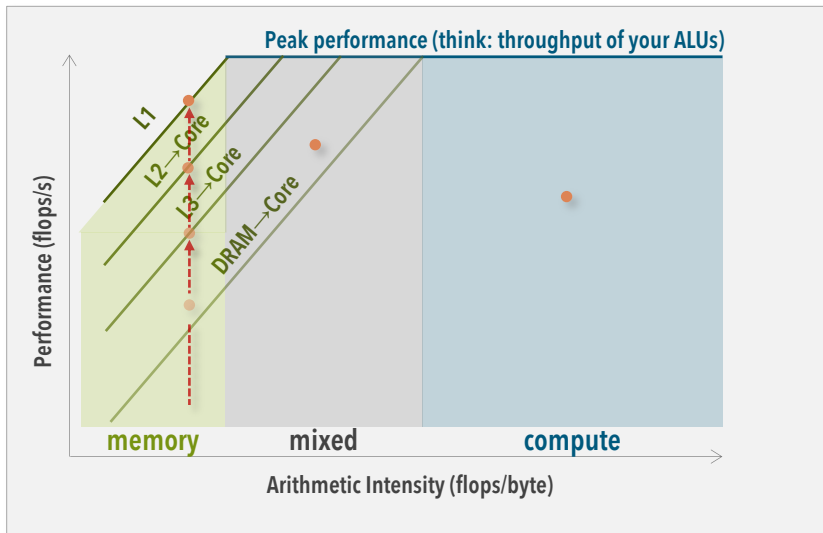
### ② Draw an imaginary vertical line

- 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...

# 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 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...



SOME THEORY



SOME EXAMPLE



SOME TOOLS



SOME ACTION



# Not all rooflines are the same

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

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

# CARM vs. ORM: Construction and Bandwidth

## CARM: Cache-aware Roofline

### ROOFS

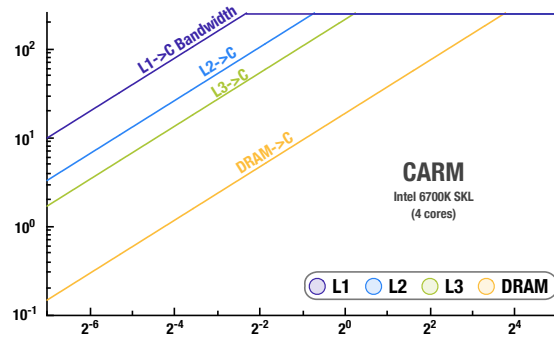
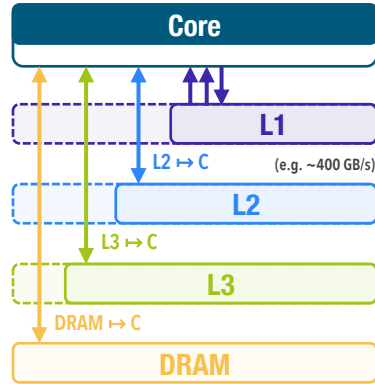
Micro-benchmarking (realistic)

### BANDWIDTH

As seen by the core

### ROOFLINE

Single-plot, single app point



## ORM: Original Roofline

### ROOFS

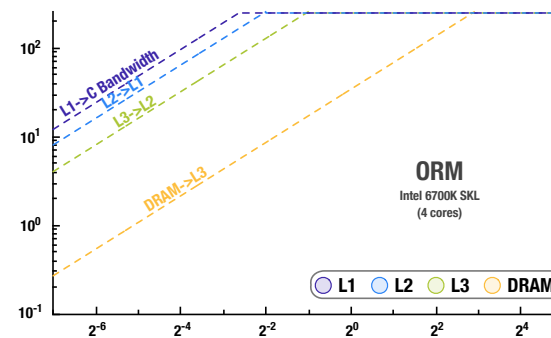
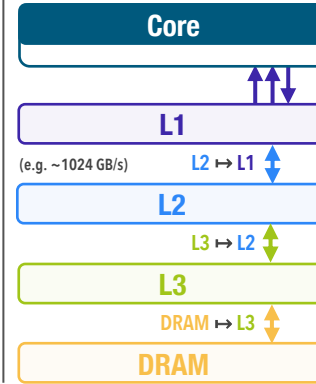
Data-sheets (theoretical)

### BANDWIDTH

Between two memory levels

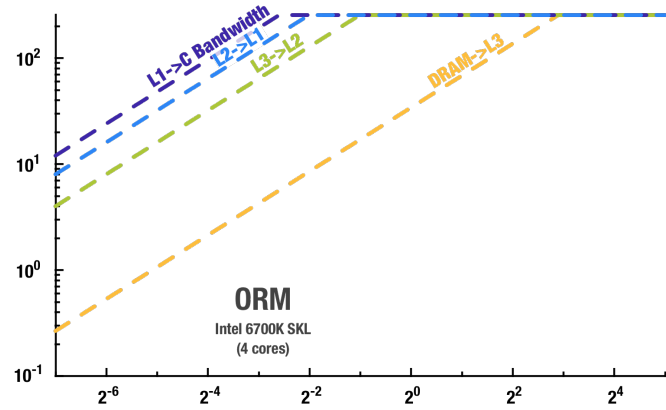
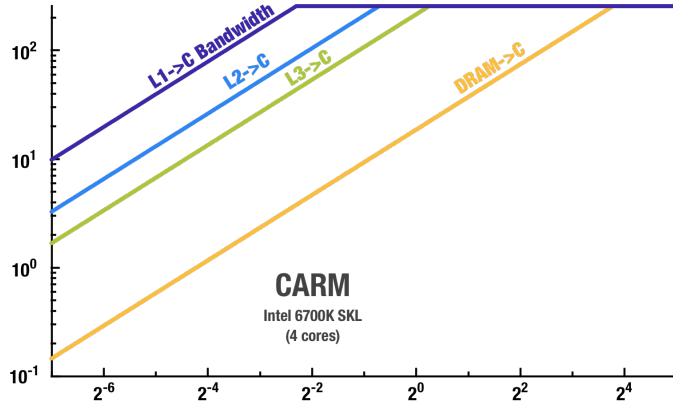
### ROOFLINE

Multi-plot, multiple app points (one for each memory level)



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

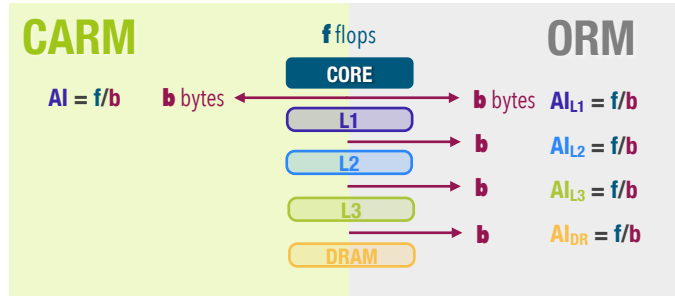


iteration: 1  
"cold caches" (misses in all caches)

```

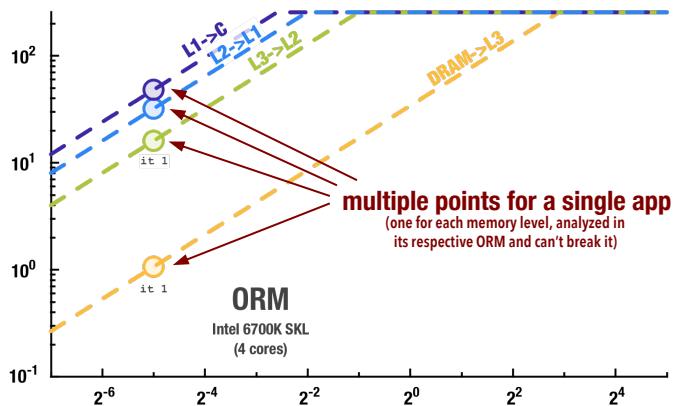
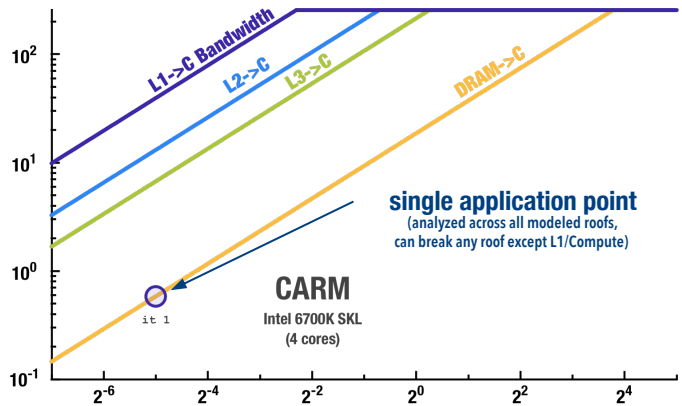
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

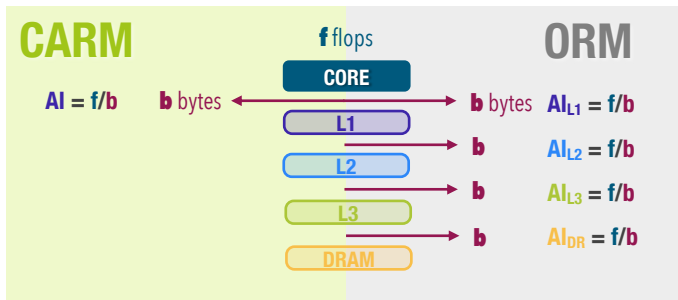


iteration: 1  
"cold caches" (misses in all caches)

```

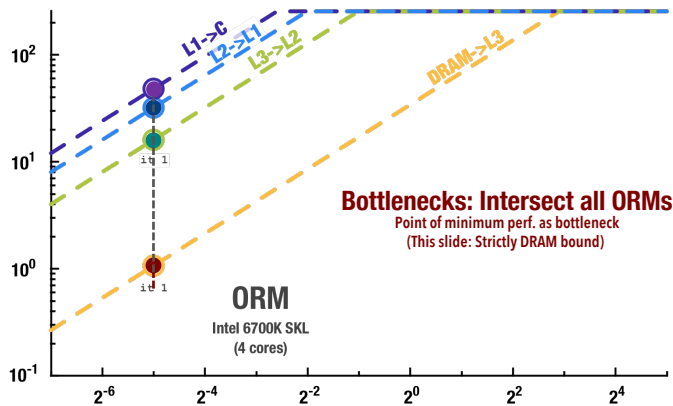
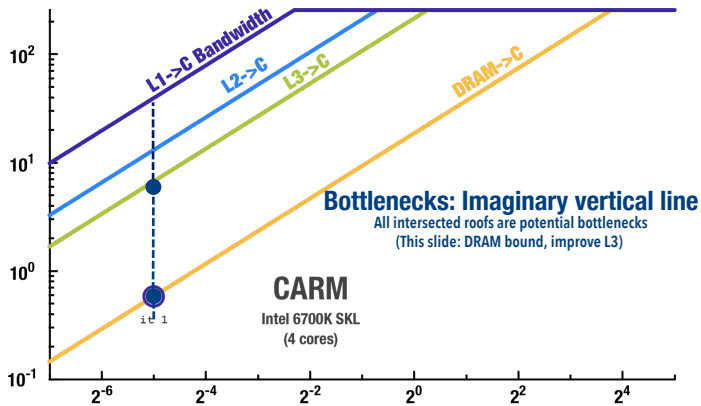
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ b bytes  
→ f flops



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

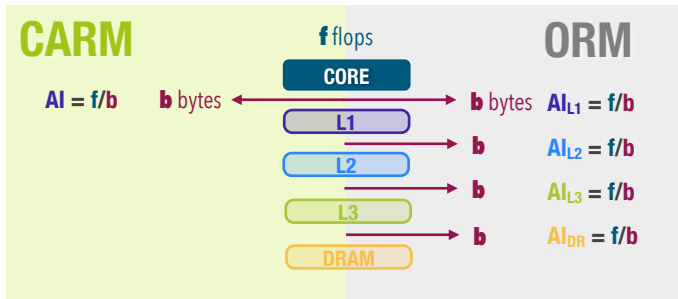


iteration: 1  
"cold caches" (misses in all caches)

```

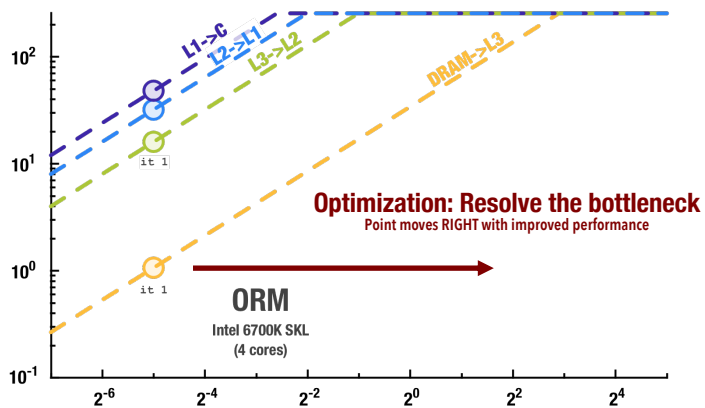
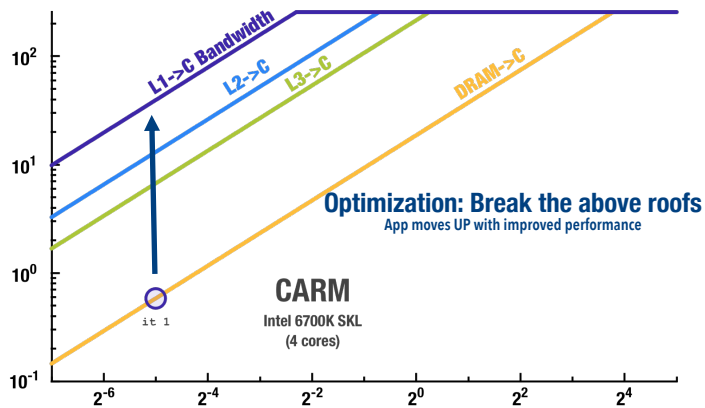
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

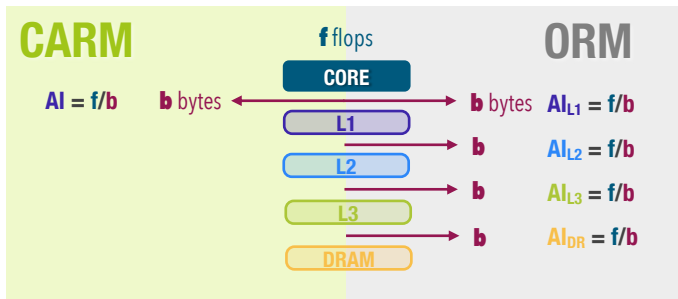


iteration: 1  
"cold caches" (misses in all caches)

```

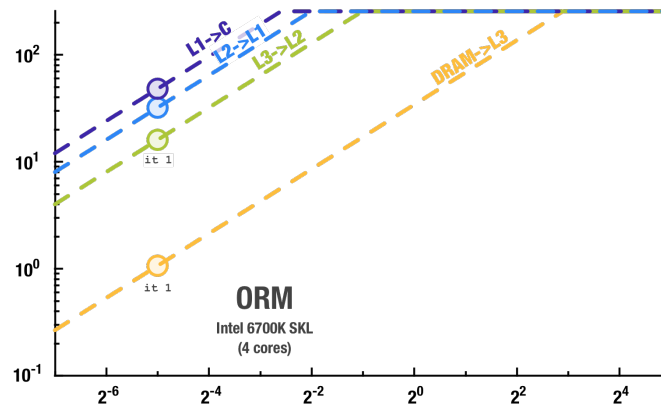
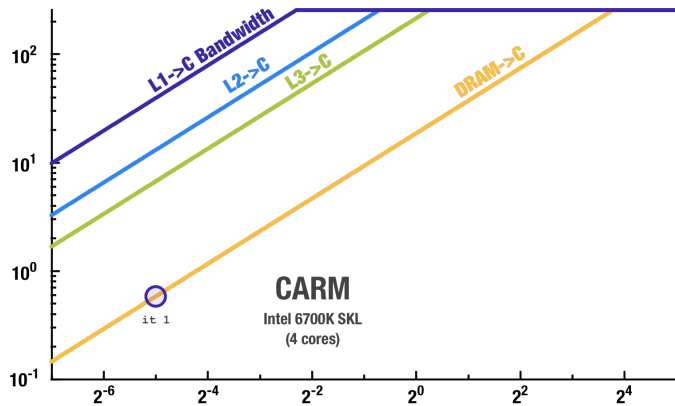
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC



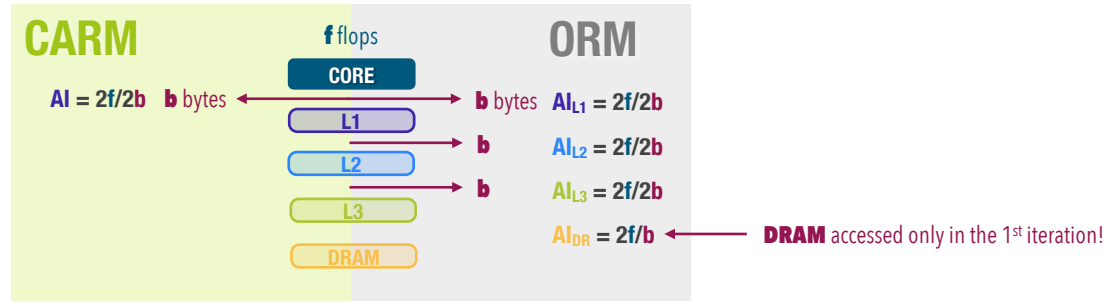
iterations: 1 and 2  
All data reused from L3

L3 Benchmark

```

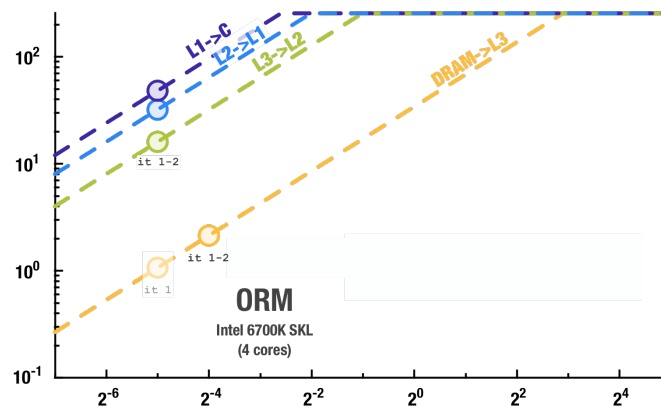
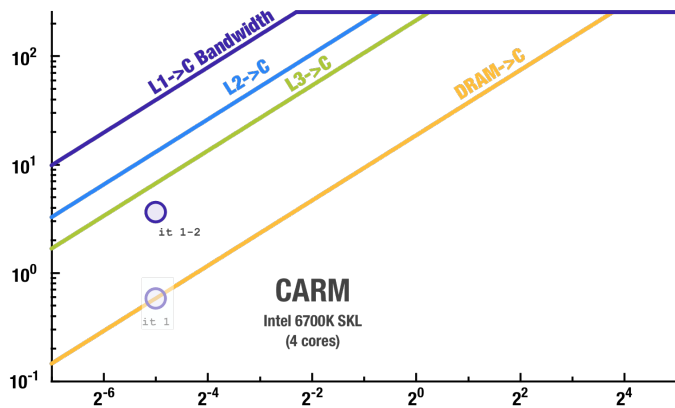
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC



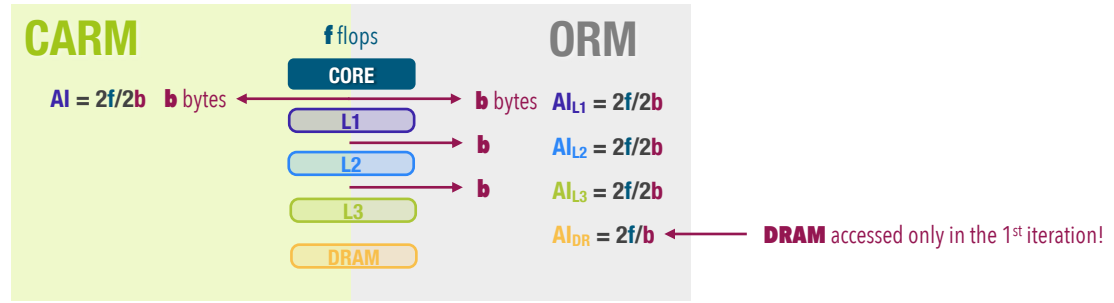
iterations: 1 and 2  
All data reused from L3

L3 Benchmark

```

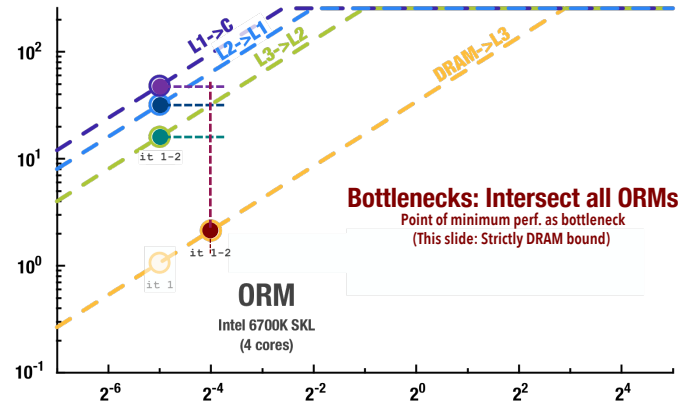
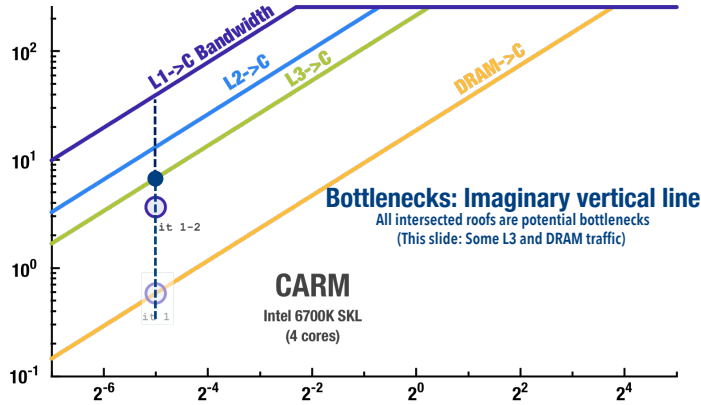
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

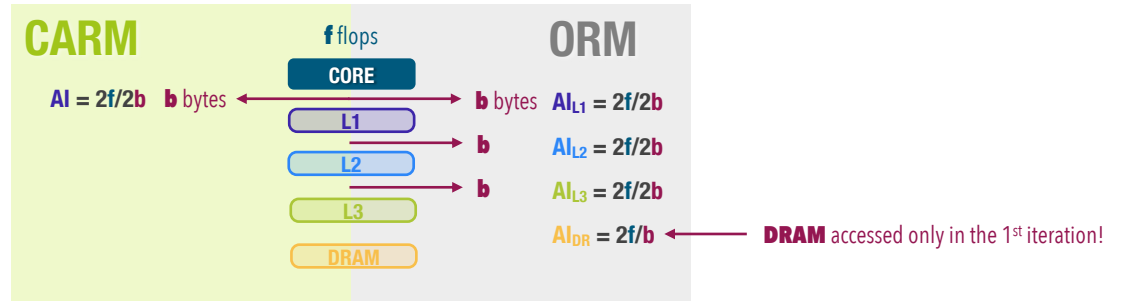


iterations: 1 and 2  
All data reused from L3

```

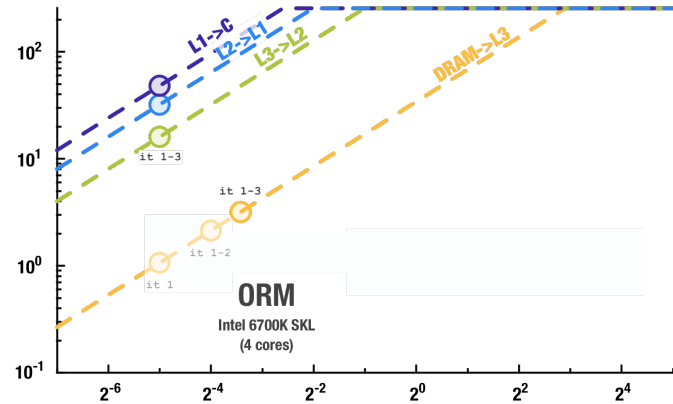
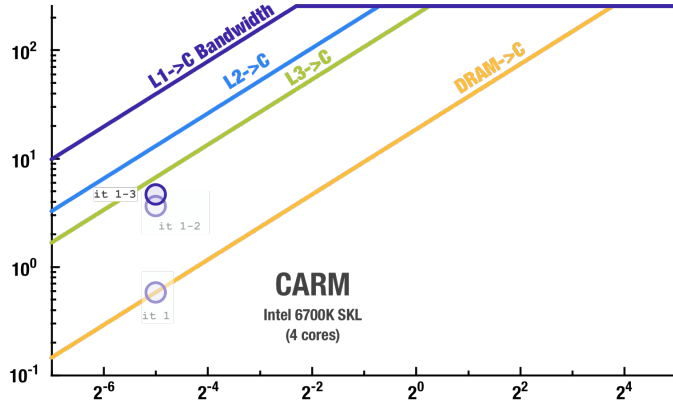
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b bytes**  
→ **f flops**



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

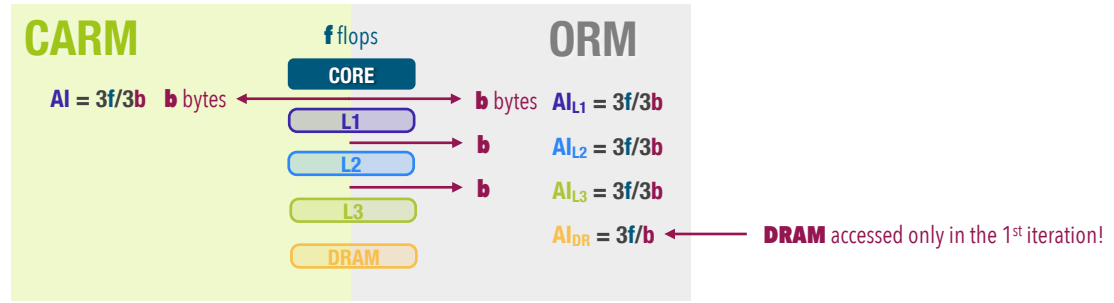


iterations: 1 to 3  
All data reused from L3

```

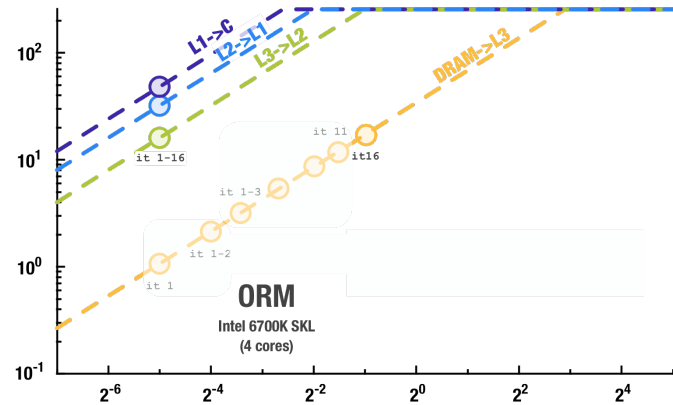
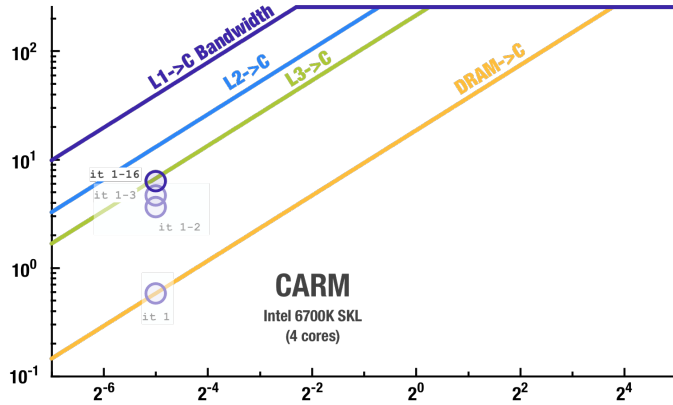
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b** bytes  
→ **f** flops



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

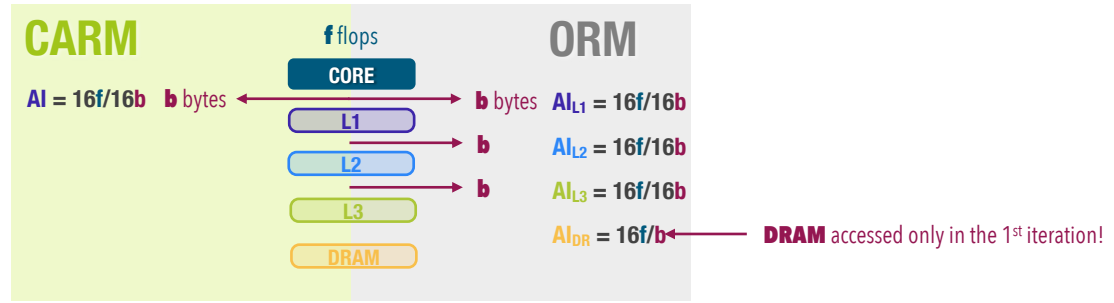


iterations: 1 to 16  
All data reused from L3

```

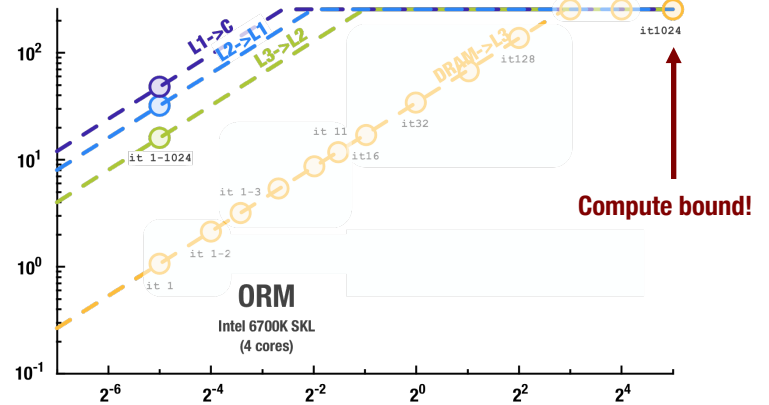
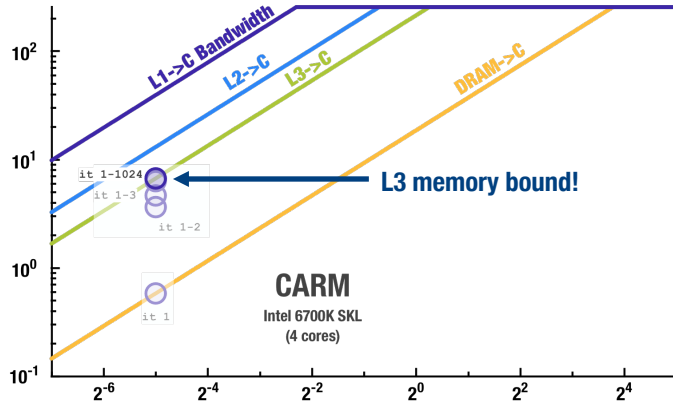
L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b** bytes  
→ **f** flops



# Example: L3 Benchmark

Memory-bound Application: Benchmark to obtain the sustainable bandwidth of LLC

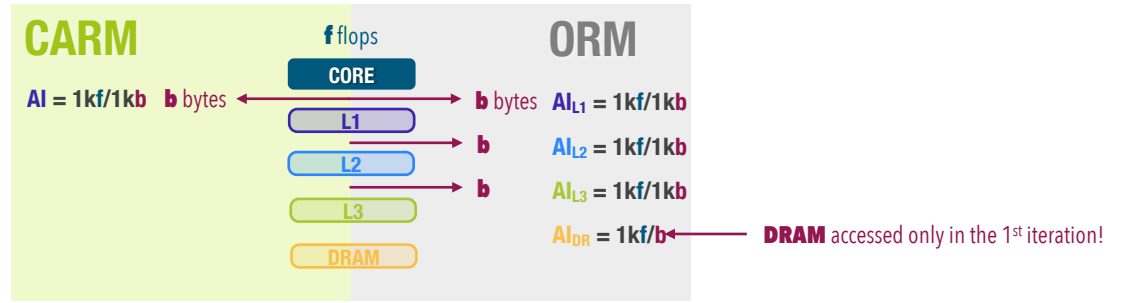


iterations: 1 to 1024  
All data reused from L3

```

L3 Benchmark
for it=1..N do
  mops (2LD+ST)
  mops (2LD+ST)
  mops (2LD+ST)
  flops (FMAD)
end for
  
```

→ **b** bytes  
→ **f** flops





SOME THEORY



SOME EXAMPLE



SOME TOOLS



SOME ACTION



# Don't do it manually!

# State of the art – Roofline Tools

	<b>Intel Advisor</b>	<b>AMD uProf</b>	<b>ERT</b>
<b>Supported Architectures</b>	Intel	AMD	Intel   AMD
<b>Supported Roofline</b>	CARM	ORM	? CARM ?
<b>Application Analysis</b>	DBI	PMUs	No support
<b>Open-Source</b>	No	No	Yes

# State of the art – Roofline Tools

	Intel Advisor	AMD uProf	ERT	CARM Tool
<b>Supported Architectures</b>	Intel	AMD	Intel   AMD	Intel   AMD ARM   RISC-V
<b>Supported Roofline</b>	CARM	ORM	? CARM ?	CARM
<b>Application Analysis</b>	DBI	PMUs	No support	DBI   PMUs
<b>Open-Source</b>	No	No	Yes	Yes

## SUPPORTED ISA (SIMD extensions)

**intel**  
**AMD**

SCALAR | SSE | AVX2 | AVX512

**arm**

SCALAR | NEON | SVE

**RISC-V**

SCALAR | RVV 0.7.1 | RVV 1.0



DEMO (SOON)

- **Interface: Two types**

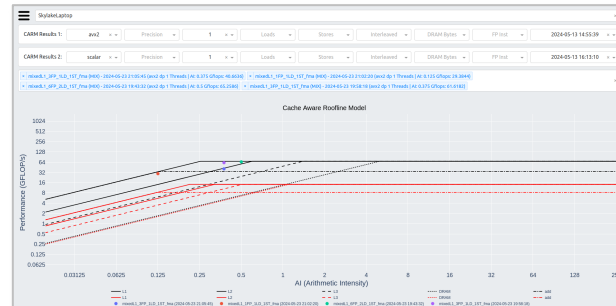
- Command Line Interface (CLI)
- Graphical User Interface (GUI)

- **Automatic Benchmarking**

- CARM Benchmarks
- Memory / FP / Mixed Benchmarks

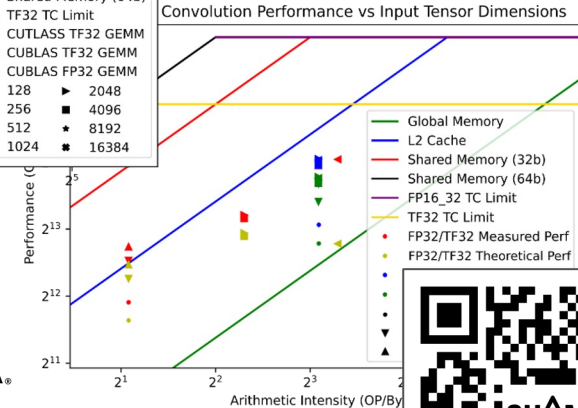
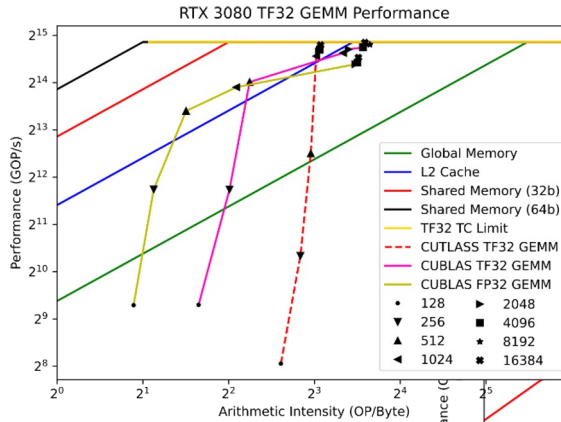
- **Application Analysis**

- PMU: PAPI
- DBI: DynamoRIO, Intel SDE



# The CARM Tool :: GPU Extension

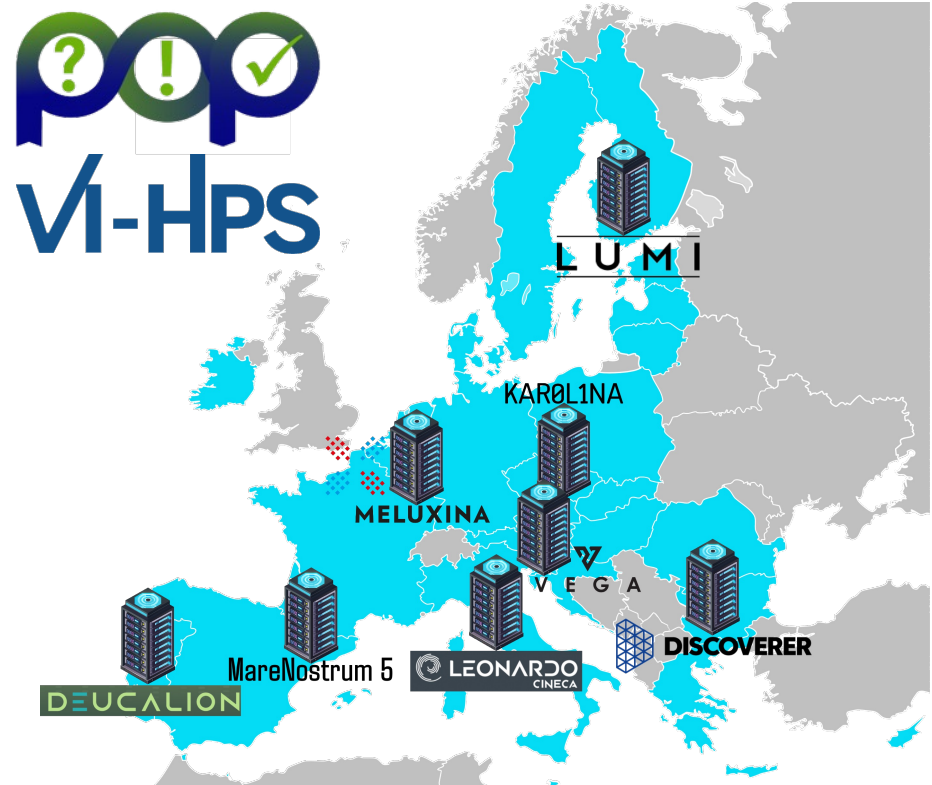
- **Nvidia GPU + Tensor CARM**
  - Compute: CUDA + Tensor Cores
  - Memory: Global, L2 cache, Shared
- **AMD GPU Extension**
  - Targeting MI GPUs via HIP
  - AMD Roc Prof for profiling
- **GPU App Profiling: Nsight Compute**
  - Entire app or specific GPU kernels
  - Potential shown in profiling convolutions etc.





# The CARM Tool :: Integration into HPC Centres

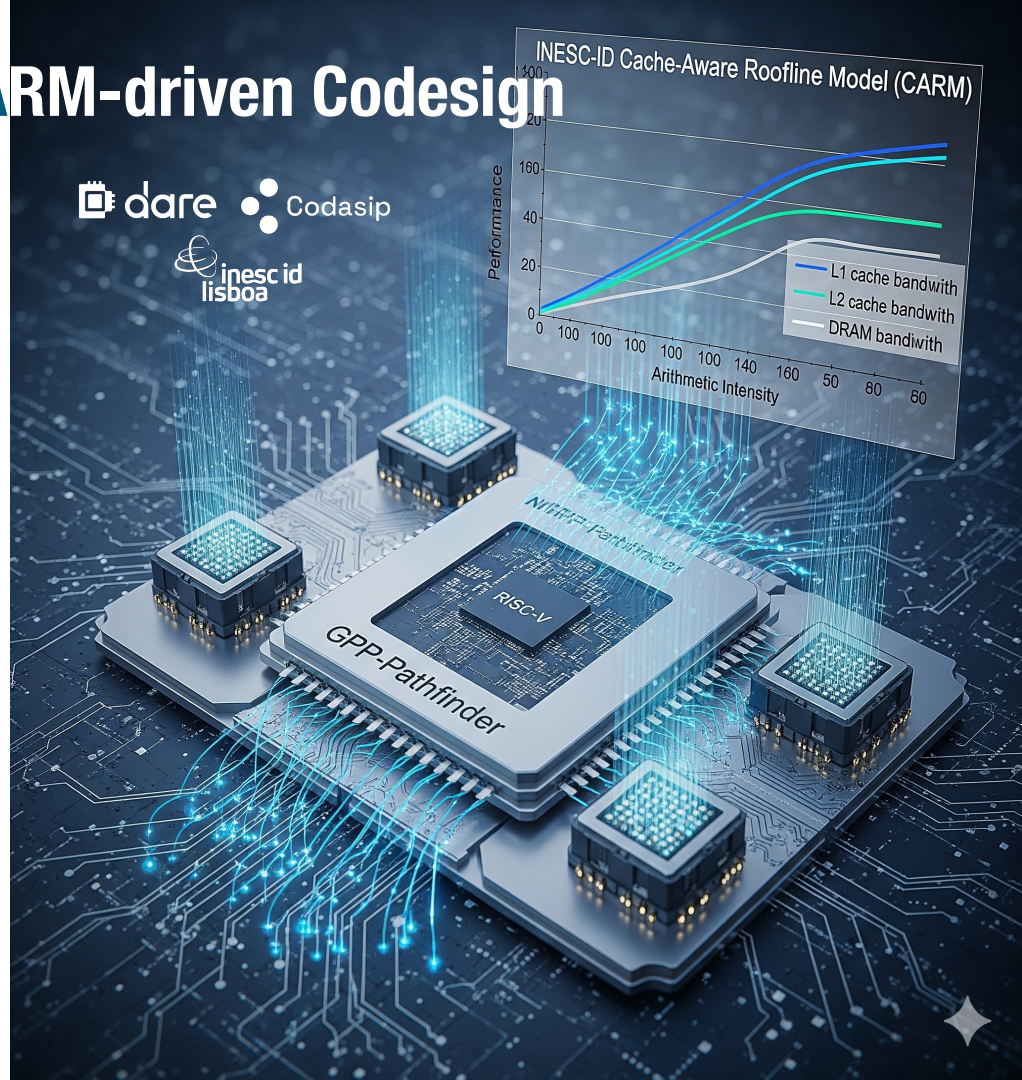
- **CARM-based application profiling**
  - Via the CARM Tool
  - Via Extrae traces and Paraver integration
- **POP3: Deployment to HPC Centres**
  - Karolina automated deployment complete
  - Available on EuroHPC centres soon
- **CARM Dissemination Initiatives**
  - POP3 Webinars
  - VI-HPS Tools Guide and Blog
  - <https://www.vi-hps.org/tools/carm.html>



# DARE Project :: CARM-driven Codesign

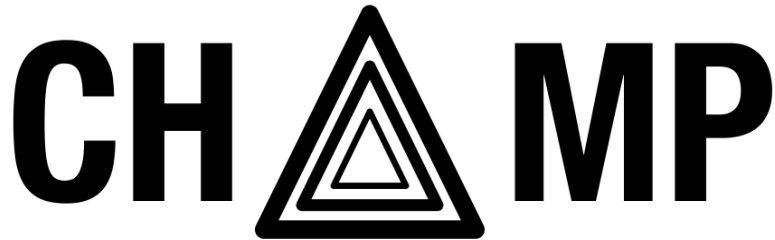


- **DARE SGA1 Project**
  - Driving Europe's next-generation HPC systems
  - Focus on performance, efficiency, and sovereignty
- **Codasip: RISC-V GPP Pathfinding**
  - Low-level architecture-aware support for CARM
  - Help optimizing the apps running on the GPP core
- **Insights about characteristics of the architecture**
- **Collaborate with the DARE SW TA**



# CARM Tool – Available on Github

- The CARM Tool is open source and available on Github
  - <https://github.com/champ-hub/carm-roofline>
- The first of many tools to be developed in the scope of CHAMP hub



**Heterogeneous Computing and Performance Modeling Hub**  
*Hub para Computação Heterogénea e Modelação de Performance*





SOME THEORY



SOME EXAMPLE



SOME TOOLS



SOME ACTION



# Finally, something to really see!



DEFINING TECHNOLOGY

Thank you!



The DARE SGA1 project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 101202459. The JU receives support from the European Union's Horizon Europe research and innovation programme and Spain, Germany, Czechia, Italy, Netherlands, Belgium, Finland, Greece, Croatia, Portugal, Poland, Sweden, France and Austria. Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or of the granting authority. Neither the European Union nor the granting authority can be held responsible for