# Automatic Cache Aware Roofline Model Building and Validation Using Topology Detection

Nicolas Denoyelle - Aleksandar Ilic - Brice Goglin - Leonel Sousa - Emmanuel Jeannot
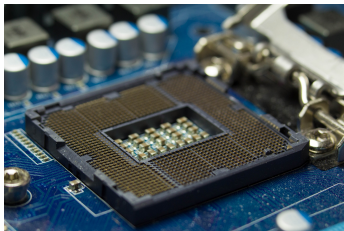
Inria (France) - INESC-ID (Portugal) - Inria (France) - INESC-ID (Portugal) - Inria (France)

October 4, 2016

# Hardware Complexity is Growing

- ↗ Cores: up to 72 on Intel KNL.
- ↘ Memory per Core.
- ↗ Cache Sharing .
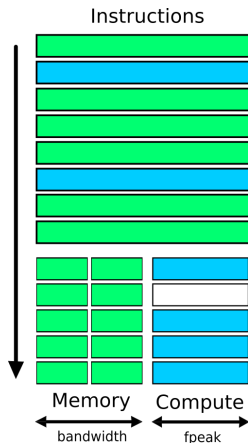- ↗ Cache Hierarchy: several levels, side cache *etc...*

The Cache Aware Roofline Model (CARM) simplifies this for you !

- What is the best you can get from the chip ?
  [Platform Model]

- Is it worth investing in code optimization ?
  [Application Model]
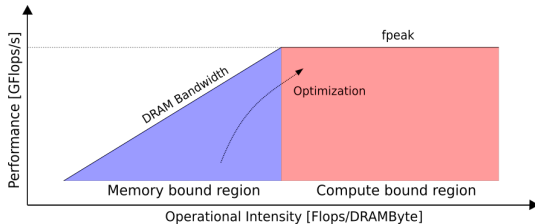
Improvement over the Original Roofline Model (ORM).

Instructions

Assumptions:

- The machine is bound whether by the memory unit or by the compute unit.
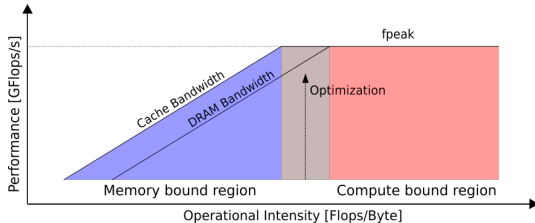- Both unit are not dependent on each other.

Memory   Compute

bandwidth   fpeak

$$ddot = ddot + b[i] * c[i] \qquad (1)$$

- memory: $b[i], c[i] = 8 + 8 Bytes$
- compute: $+, * = 2 flops$
- arithmetic intensity: $2^{-3} Flops/Bytes$
- performance: $Flops/Time$

# ORM versus CARM



Bytes transferred from DRAM to last level cache



Bytes issued from Cores.

# A Tool to Build and Validate the Model

## Build Model from Portable Blocks

- Hwloc gives the machine informations: cache structure and attributes, cores.
- The compiler gives architecture type and let us select the best instruction set.

## Micro-Benchmark Each Unit

- Bandwidth benchmark each cache level, and several stream types: Load, Store, 2Loads/Store.
- Floating Point peak benchmark: MUL, ADD, MAD, FMA.
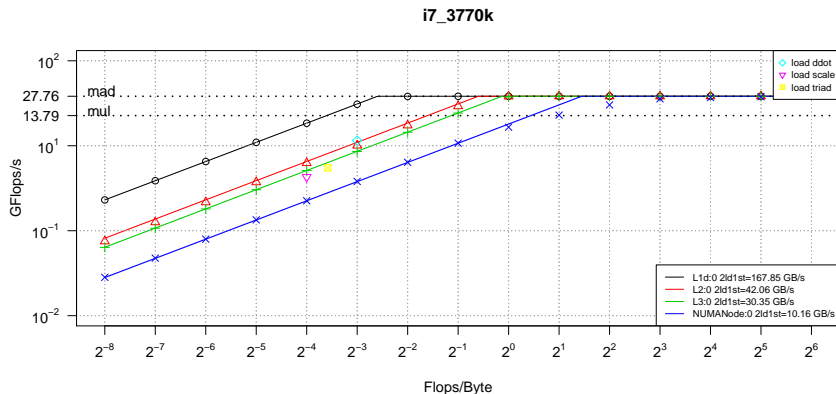- Micro Benchmarks of several arithmetic intensities for validation.

## Easy steps

- Build the model for the target platform.
- Collect CARM metrics from your application with our library (Hardware counters: PAPI).
- Plot this on a chart and get insights on your application.

## Quick Glance if the Model is Wrong

- Validation matches model ?
- Micro-Benchmarks matches theoretical throughput ?

i7_3770k

## Conclusion

We provide a Tool to build the CARM and model applications.
It is portable and easy to use. It reaches near theoretical architecture throughput.

## Future Work

The CARM gives insight over Caches (i.e temporal locality).
We will use it to model system heteogeneous memory, and give insights on data locality.

Thank you !

https://github.com/NicolasDenoyelle/LARM-Locality-Aware-Roofline-Model-

# Automatic Cache Aware Roofline Model Building and Validation Using Topology Detection

Nicolas Denoyelle - Aleksandar Ilic - Brice Goglin - Leonel Sousa - Emmanuel Jeannot

Inria (France) - INESC-ID (Portugal) - Inria (France) - INESC-ID (Portugal) - Inria (France)

October 4, 2016