# Roofline: A Throughput Oriented Performance Model

## Lenny Oliker

### Jack Deslippe, Tuomas Koskela, Samuel Williams
*Lawrence Berkeley National Laboratory, USA*

### Roman Belenov, Zakhar Matveev, Philippe Thierry
*Intel Corporation*

# Why Use Performance Models or Tools?

- Identify performance bottlenecks

- Motivate software optimizations

- **Determine when we're done optimizing**

  - Assess performance relative to machine capabilities

  - Motivate need for algorithmic changes

- Predict performance on future machines / architectures

  - Sets realistic expectations on performance for future procurements

  - Used for HW/SW Co-Design to ensure future architectures are well-suited for the computational needs of today's applications.

BERKELEY LAB

# Performance Models

- Many different components can contribute to kernel run time.

- Some are characteristics of the application, some are characteristics of the machine, and some are both (memory access pattern + caches).

| | |
|---|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| PCIe data movement | PCIe bandwidth |
| Depth | OMP Overhead |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

BERKELEY LAB

# **Performance Models**

- Can't think about all these terms all the time for every application…

**Computational Complexity**

| | |
|---|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| PCIe data movement | PCIe bandwidth |
| Depth | OMP Overhead |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

# Performance Models

- Because there are so many components, performance models often conceptualize the system as being dominated by one or more of these components.

| | |
|---|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| PCIe data movement | PCIe bandwidth |
| Depth | OMP Overhead |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

**Roofline Model**

Williams et al, "Roofline: An Insightful Visual Performance Model For Multicore Architectures", CACM, 2009.

BERKELEY LAB

# Performance Models

- Because there are so many components, performance models often conceptualize the system as being dominated by one or more of these components.

| | |
|---|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| **PCIe data movement** | **PCIe bandwidth** |
| **Depth** | **OMP Overhead** |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

**LogCA**

BERKELEY LAB

# Performance Models

- Because there are so many components, performance models often conceptualize the system as being dominated by one or more of these components.

| | |
|---:|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| PCIe data movement | PCIe bandwidth |
| Depth | OMP Overhead |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

**LogGP**

Alexandrov, et al, "LogGP: incorporating long messages into the LogP model - one step closer towards a realistic model for parallel computation", SPAA, 1995.

BERKELEY LAB

# Performance Models

- Because there are so many components, performance models often conceptualize the system as being dominated by one or more of these components.

| | |
|---|---|
| #FP operations | Flop/s |
| Cache data movement | Cache GB/s |
| DRAM data movement | DRAM GB/s |
| PCIe data movement | PCIe band... |
| Depth | OMP C... |
| MPI Message Size | Network B... |
| MPI Send:Wait ratio | Network Gap... |
| #MPI Wait's | Network Latency |

LogP

**Right model depends on app and problem size!**

Culler, et al, "LogP: a practical model of parallel computation", CACM, 1996.

BERKELEY LAB

# Performance Models / Simulators

- Historically, many performance models and simulators tracked latencies to predict performance (i.e. counting cycles)

- The last two decades saw a number of latency-hiding techniques…
  - Out-of-order execution (hardware discovers parallelism to hide latency)
  - HW stream prefetching (hardware speculatively loads data)
  - Massive thread parallelism (independent threads satisfy the latency-bandwidth product)

- Effective latency hiding has resulted in a shift from a latency-limited computing regime to a **throughput-limited computing regime**

# Roofline Model

- **Roofline Model** is a throughput-oriented performance model…

  - Tracks <u>rates</u> not times

  - Augmented with Little's Law
    (concurrency = latency*bandwidth)

  - Independent of ISA and architecture (applies to CPUs, GPUs, Google TPUs[1],  etc…)



https://crd.lbl.gov/departments/computer-science/PAR/research/roofline

[1]Jouppi et al, "In-Datacenter Performance Analysis of a Tensor Processing Unit", ISCA, 2017.

# (DRAM) Roofline

- One could hope to always attain peak performance (Flop/s)

- However, finite locality (reuse) and bandwidth limit performance.

- Assume:
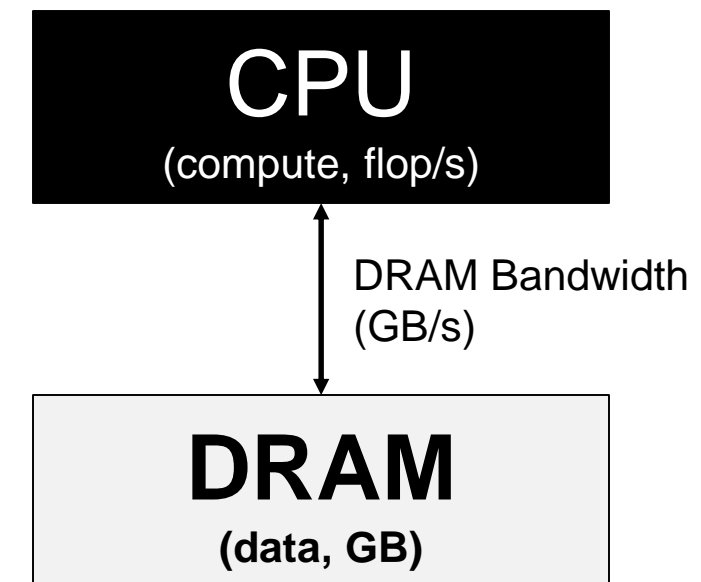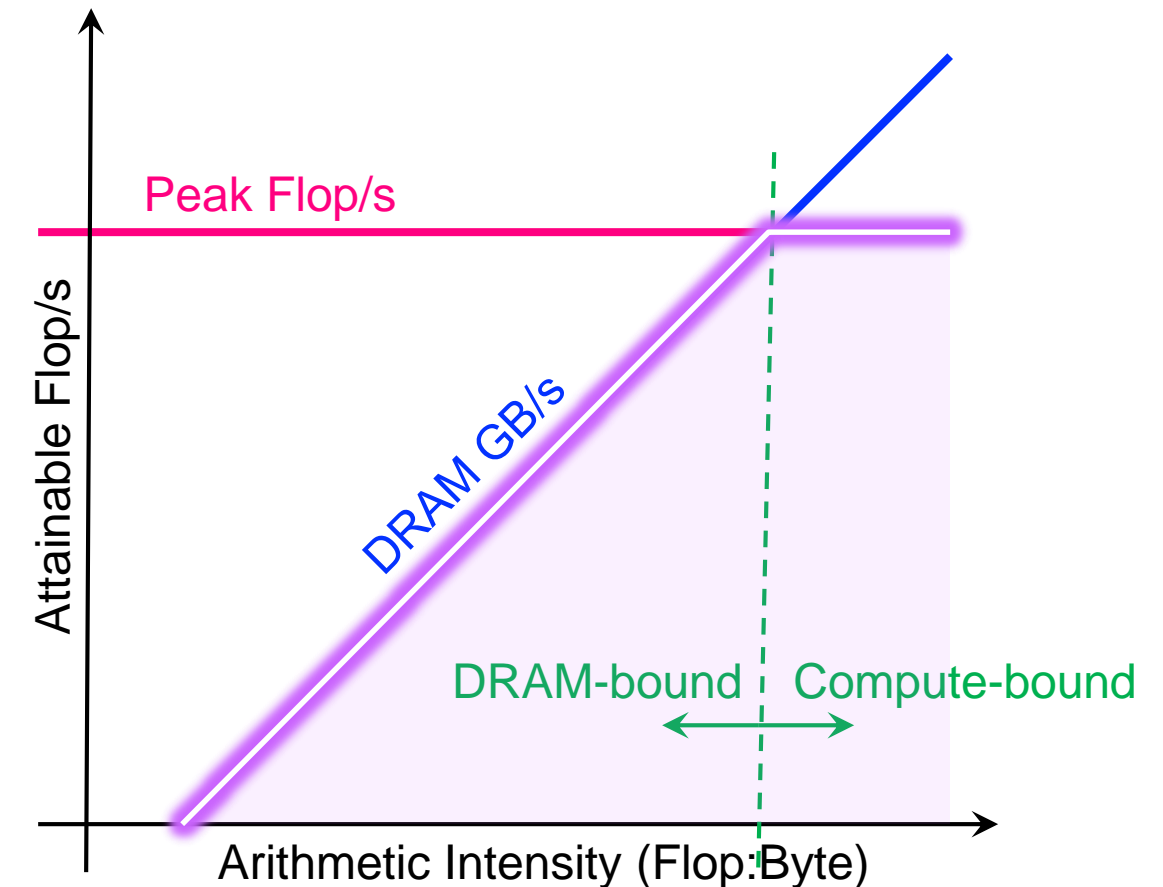  - Idealized processor/caches
  - Cold start (data in DRAM)

**CPU**
(compute, flop/s)

↕ DRAM Bandwidth (GB/s)

**DRAM**
**(data, GB)**

$$\text{Time} = \max \begin{cases} \text{\#FP ops / Peak GFlop/s} \\ \\ \text{\#Bytes / Peak GB/s} \end{cases}$$

BERKELEY LAB

# (DRAM) Roofline

- One could hope to always attain peak performance (Flop/s)

- However, finite locality (reuse) and bandwidth limit performance.

- Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)

CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

$$\frac{Time}{\#FP\ ops} = \max \begin{cases} 1\ /\ Peak\ GFlop/s \\ \#Bytes\ /\ \#FP\ ops\ /\ Peak\ GB/s \end{cases}$$

BERKELEY LAB

# (DRAM) Roofline

- One could hope to always attain peak performance (Flop/s)

- However, finite locality (reuse) and bandwidth limit performance.

- Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)

CPU
(compute, flop/s)

DRAM Bandwidth (GB/s)

DRAM
(data, GB)

$$\frac{\#FP\ ops}{Time} = \min \begin{cases} \text{Peak GFlop/s} \\ (\#FP\ ops\ /\ \#Bytes) * \text{Peak GB/s} \end{cases}$$

# (DRAM) Roofline

- One could hope to always attain peak performance (Flop/s)

- However, finite locality (reuse) and bandwidth limit performance.

- Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)

CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

$$
\text{GFlop/s} = \min \begin{cases} \textbf{Peak GFlop/s} \\ \\ \textbf{AI * Peak GB/s} \end{cases}
$$

*Note, Arithmetic Intensity (AI) = Flops / Bytes (as presented to DRAM )*

# (DRAM) Roofline

- Plot Roofline bound using Arithmetic Intensity as the x-axis

- **Log-log scale** makes it easy to doodle, extrapolate performance along Moore's Law, etc…

- Kernels with AI less than machine balance are ultimately DRAM bound (we'll refine this later…)

# Roofline Example #1

- **Typical machine balance is 5-10 flops per byte…**

  - 40-80 flops per double to exploit compute capability
  - Artifact of technology and money
  - **Unlikely to improve**

- **Consider STREAM Triad…**

```
#pragma omp parallel for
for(i=0;i<N;i++){
  Z[i] = X[i] + alpha*Y[i];
}
```

  - 2 flops per iteration
  - Transfer 24 bytes per iteration (read X[i], Y[i], write Z[i])
  - **AI = 0.083 flops per byte == Memory bound**

# Roofline Example #2

- Conversely, 7-point constant coefficient stencil…

  - 7 flops

  - 8 memory references (7 reads, 1 store) per point

  - Cache can filter all but 1 read and 1 write per point

  - **AI = 0.44 flops per byte == memory bound, but 5x the flop rate**

```
#pragma omp parallel for
for(k=1;k<dim+1;k++){
for(j=1;j<dim+1;j++){
for(i=1;i<dim+1;i++){
   new[k][j][i] = -6.0*old[k  ][j  ][i  ]
                      + old[k  ][j  ][i-1]
                      + old[k  ][j  ][i+1]
                      + old[k  ][j-1][i  ]
                      + old[k  ][j+1][i  ]
                      + old[k-1][j  ][i  ]
                      + old[k+1][j  ][i  ];
}}}
```

# Hierarchical Roofline

- Real processors have multiple levels of memory
  - Registers
  - L1, L2, L3 cache
  - MCDRAM/HBM (KNL/GPU device memory)
  - DDR (main memory)
  - NVRAM (non-volatile memory)

- Applications can have locality in each level
  - Unique data movements imply unique AI's
  - Moreover, each level will have a unique bandwidth

BERKELEY LAB

# Hierarchical Roofline

- Construct superposition of Rooflines…

  - Measure a bandwidth

  - Measure AI for each level of memory

  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…

  - **… performance is bound by the minimum**

# Hierarchical Roofline

- Construct superposition of Rooflines…
  - Measure a bandwidth
  - Measure AI for each level of memory
  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…
  - **… performance is bound by the minimum**



Peak Flop/s

L2 GB/s

MCDRAM cache GB/s

DDR GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

**DDR bottleneck pulls performance below MCDRAM Roofline**

# Hierarchical Roofline

- Construct superposition of Rooflines…

  - Measure a bandwidth

  - Measure AI for each level of memory

  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…

  - **… performance is bound by the minimum**

# Hierarchical Roofline

- Construct superposition of Rooflines…

  - Measure a bandwidth

  - Measure AI for each level of memory

  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…

  - **… performance is bound by the minimum**

# Data, Instruction, Thread-Level Parallelism…

- Modern CPUs use several techniques to increase per core Flop/s

## Fused Multiply Add

- w = x*y + z is a common idiom in linear algebra
- Ra~~ther~~ ~~performing~~ se~~parate~~ ~~multiply~~ ~~and~~ ~~add, use a~~ ~~fused multiply add (FMA)~~
- Th~~e FPU~~ chains the multiply and add in a single pipeline so that it can complete FMA/cycle

*Resurgence… Tensor Cores, QFMA, etc…*

## Vector Instructions

- Many HPC codes apply the same operation to a vector of elements
- Vendors provide vector instructions that apply the same operation to 2, 4, 8, 16 elements…
  x [0:7] *y [0:7] + z [0:7]
- Vector FPUs complete 8 vector operations/cycle

## Deep pipelines

- The hardware for a FMA is substantial.
- Breaking a single FMA up into several smaller operations and pipelining them allows vendors to increase GHz
- Little's Law applies… need FP_Latency * FP_bandwidth independent instructions

BERKELEY LAB

# Data, Instruction, Thread-Level Parallelism...

- If every instruction were an ADD (instead of FMA), **performance would drop by 2x on KNL or 4x on Haswell**

- Similarly, if one failed to vectorize, performance would drop by **another 8x on KNL and 4x on Haswell**

- Lack of threading (or load imbalance) will reduce performance by another 64x on KNL.

# Superscalar vs. Instruction mix

- Define in-core ceilings based on instruction mix…

- e.g. Haswell

  - 4-issue superscalar

  - Only 2 FP data paths

  - Requires 50% of the instructions to be FP to get peak performance

# Superscalar vs. Instruction mix

- Define in-core ceilings based on instruction mix…

- e.g. Haswell
  - 4-issue superscalar
  - Only 2 FP data paths
  - Requires 50% of the instructions to be FP to get peak performance

- e.g. KNL
  - 2-issue superscalar
  - 2 FP data paths
  - Requires 100% of the instructions to be FP to get peak performance

# Superscalar vs. instruction mix

- Define in-core ceilings based on instruction mix…

- e.g. Haswell
  - 4-issue superscalar
  - Only 2 FP data paths
  - Requires 50% of the instructions to be FP to get peak performance

- e.g. KNL
  - 2-issue superscalar
  - 2 FP data paths
  - Requires 100% of the instructions to be FP to get peak performance

# Superscalar vs. instruction mix

- Define in-core ceilings based on instruction mix…

- e.g. Haswell
  - 4-issue superscalar
  - Only 2 FP data paths
  - Requires 50% of the instructions to be FP to get peak performance

- e.g. KNL
  - 2-issue superscalar
  - 2 FP data paths
  - Requires 100% of the instructions to be FP to get peak performance



Peak Flop/s

100% FP

50% FP (50% int)

25% FP (75% int)

DDR GB/s

Attainable Flop/s

Arithmetic Intensity

non-FP instructions can sap instruction issue bandwidth and pull performance below Roofline

# Divides and other Slow FP instructions

- FP Divides (sqrt, rsqrt, …) might support only limited pipelining

- As such, their throughput is substantially lower than FMA's

- If divides constitute even if 3% of the flop's come from divides, performance can be cut in half.

- *Penalty varies substantially between architectures and generations (e.g. IVB, HSW, KNL, …)*

# Roofline Model:
## Modeling Cache Effects

# Locality Walls

- Naively, we can bound AI using only compulsory cache misses



$$AI = \frac{\text{\#Flop's}}{\text{Compulsory Misses}}$$

# Locality Walls

- Naively, we can bound AI using only compulsory cache misses

- However, write allocate caches can lower AI



$$AI = \frac{\#Flop's}{Compulsory\ Misses + Write\ Allocates}$$

# Locality Walls

- Naively, we can bound AI using only compulsory cache misses

- However, write allocate caches can lower AI

- Cache capacity misses can have a huge penalty



$$AI = \frac{\#Flop's}{Compulsory\ Misses + Write\ Allocates + Capacity\ Misses}$$

# Locality Walls

- Naively, we can bound AI using only compulsory cache misses

- However, write allocate caches can lower AI

- Cache capacity misses can have a huge penalty

➢ **Compute bound became memory bound**

$$AI = \frac{\#Flop's}{Compulsory\ Misses + Write\ Allocates + Capacity\ Misses}$$

# General Strategy Guide

- Broadly speaking, there are three approaches to improving performance:

# General Strategy Guide

- Broadly speaking, there are three approaches to improving performance:

- **Maximize in-core performance (e.g. get compiler to vectorize)**

# General Strategy Guide

- Broadly speaking, there are three approaches to improving performance:

- Maximize in-core performance (e.g. get compiler to vectorize)

- **Maximize memory bandwidth (e.g. NUMA-aware allocation)**

# General Strategy Guide

- Broadly speaking, there are three approaches to improving performance:

- Maximize in-core performance (e.g. get compiler to vectorize)

- Maximize memory bandwidth (e.g. NUMA-aware allocation)

- **Minimize data movement (increase AI)**

**Constructing a Roofline Model requires answering some questions…**

# Questions can overwhelm users…

**Properties of the target machine**

**(Benchmarking)**

What is my machine's peak flop/s?

How important is FMA on my machine?

What is my machine's DDR GB/s?

L2 GB/s?

**Properties of an application's execution**

**(Instrumentation)**

How much data did my kernel actually move?

Did my kernel vectorize?

How many flop's did my kernel actually do?

How much did that divide hurt?

**Fundamental properties of the kernel constrained by hardware**

**(Theory)**

What is my kernel's compulsory AI? (communication lower bounds)

Can my kernel ever be vectorized?

# Node Characterization?

- **"Marketing Numbers"** can be deceptive…
  - Pin BW vs. real bandwidth
  - TurboMode / Underclock for AVX
  - compiler failings on high-AI loops.

- LBL developed the Empirical Roofline Toolkit (ERT)…
  - Characterize CPU/GPU systems
  - Peak Flop rates
  - Bandwidths for each level of memory
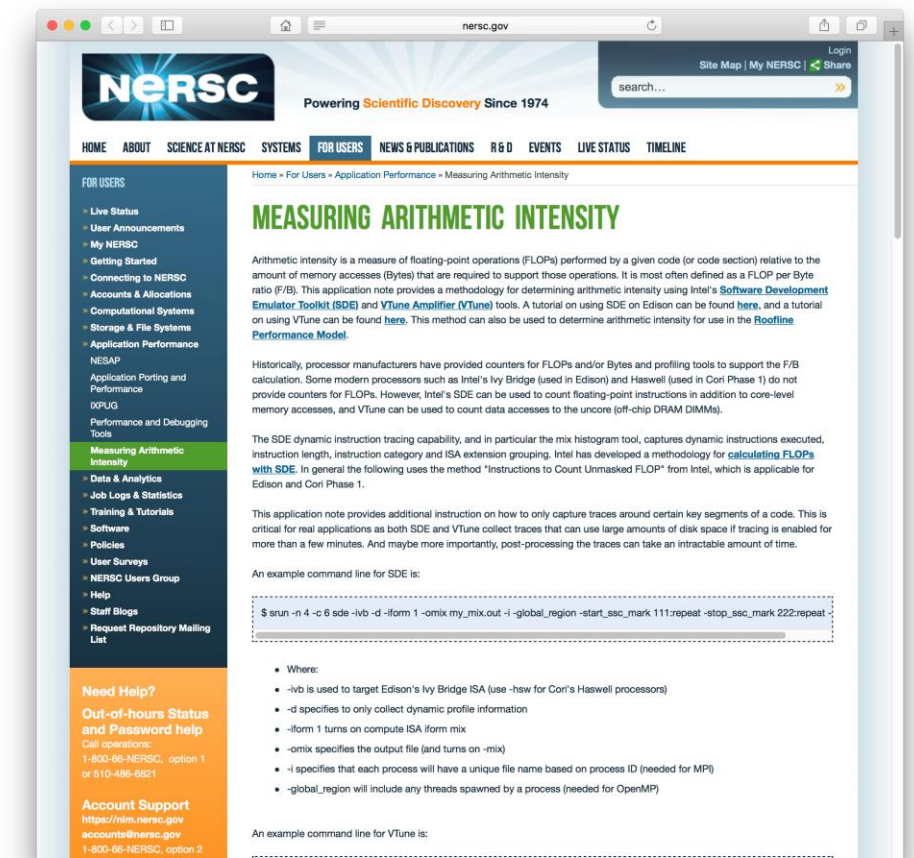  - **MPI+OpenMP/CUDA == multiple GPUs**





https://crd.lbl.gov/departments/computer-science/PAR/research/roofline/

# Instrumentation with Performance Counters?

- ***Characterizing applications with performance counters can be problematic…***
  - x Flop Counters can be broken/missing in production processors
  - x Vectorization/Masking can complicate counting Flop's
  - x Counting Loads and Stores doesn't capture cache reuse while counting cache misses doesn't account for prefetchers.
  - x DRAM counters (Uncore PMU) might be accurate, but…
    - x are privileged and thus nominally inaccessible in user mode
    - x may need vendor (e.g. Cray) and center (e.g. NERSC) approved OS/kernel changes

# Forced to Cobble Together Tools…

- Use tools known/observed to work on NERSC's Cori (KNL, HSW)…
  - Used **Intel SDE** (Pin binary instrumentation + emulation) to create software Flop counters
  - Used **Intel VTune** performance tool (NERSC/Cray approved) to access uncore counters
- ➢ Accurate measurement of Flop's (HSW) and DRAM data movement (HSW and KNL)
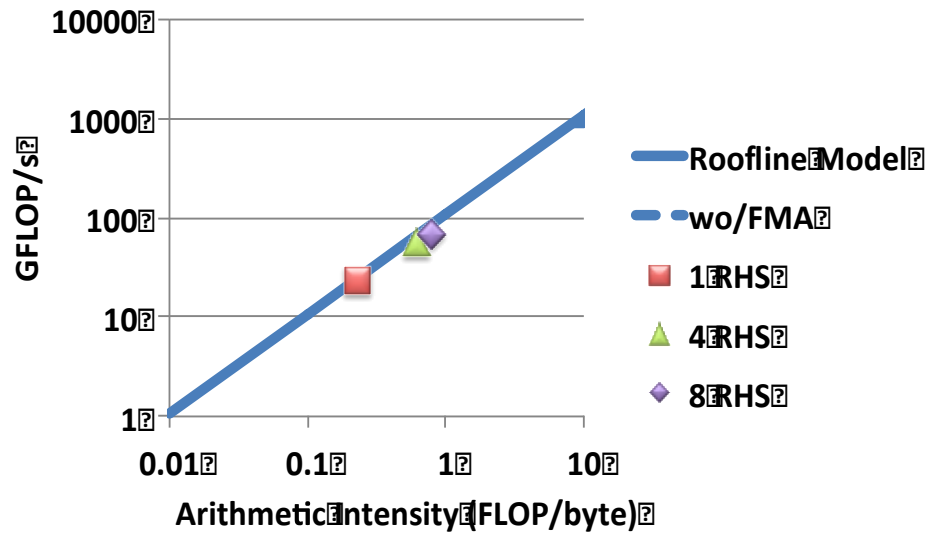- ➢ Used by NESAP (NERSC KNL application readiness project) to characterize apps on Cori…



http://www.nersc.gov/users/application-performance/measuring-arithmetic-intensity/
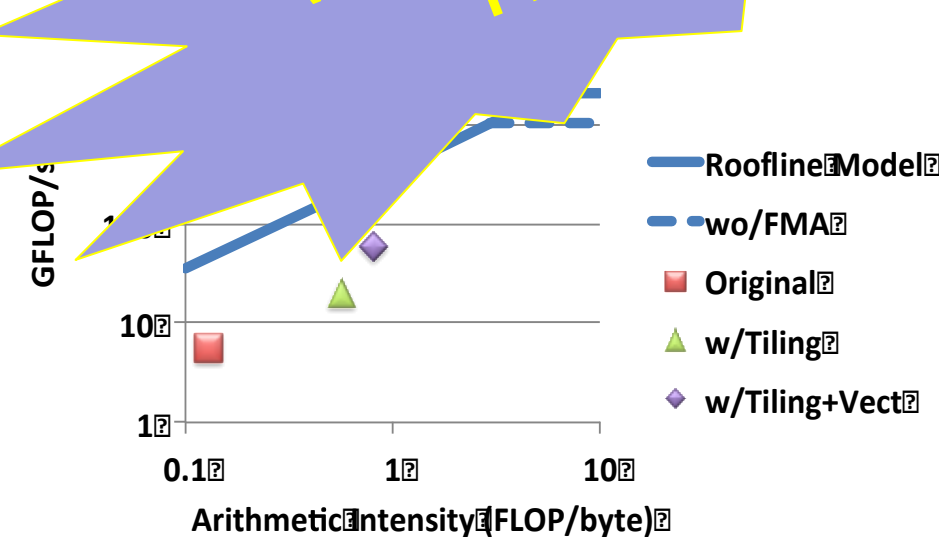
# Initial Roofline Analysis of NESAP Codes



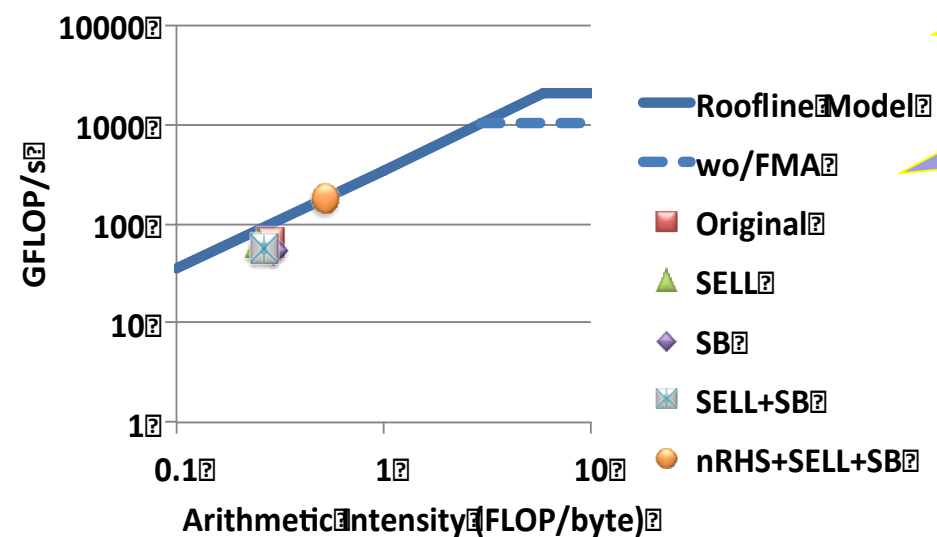Nuclear physics — MFDn (2P HSW, KNL)

Geophysical imaging — EMGeo (2P HSW, KNL)

Particle/laser interaction — PICSAR

DRAM-only Roofline was insufficient for PICSAR

# Evaluation of LIKWID
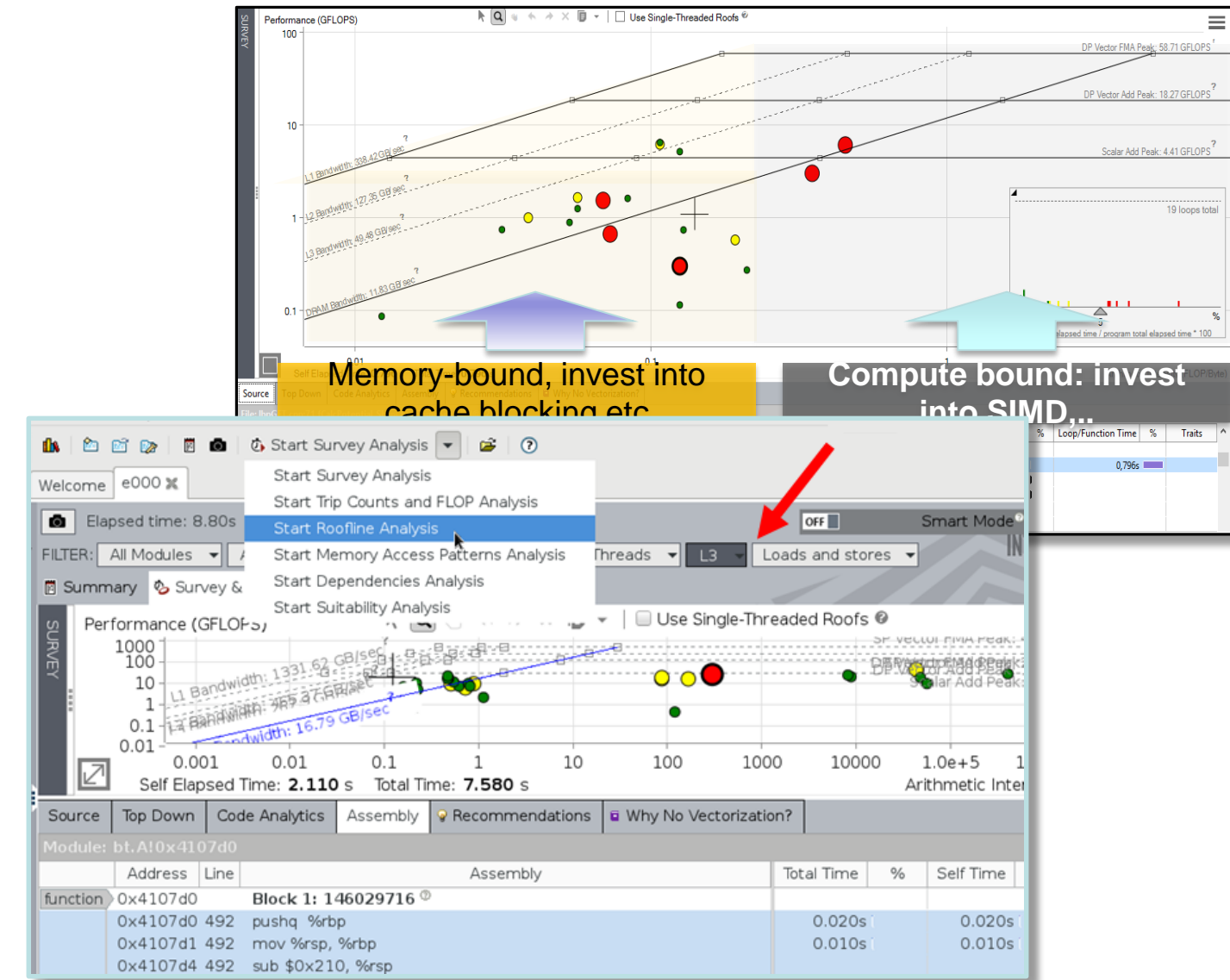
- **LIKWID provides easy to use wrappers for measuring performance counters…**
  - ✓ **Works on NERSC production systems**
  - ✓ Minimal overhead (<1%)
  - ✓ Scalable in distributed memory (MPI-friendly)
  - ✓ Fast, high-level characterization
  - ✗ **No detailed timing breakdown or optimization advice**
  - ✗ **Limited by quality of hardware performance counter implementation (garbage in/garbage out)**
- ➤ **Useful tool that complements other tools**

https://github.com/RRZE-HPC/likwid

## AMReX Application Characterization
### (2Px16c HSW == Cori Phase 1)



Legend:
- L2
- L3
- DRAM
- Roofline

Y-axis: Bandwidth(GB/s) — 8, 16, 32, 64, 128, 256, 512, 1024

X-axis categories: HPGMG (32Px1T), HPGMG (4Px8T), Combustor (32Px1T), Combustor (4Px8T), MFIX (32Px1T), Nyx (32Px1T), Nyx (4Px8T), PeleLM (32Px1T), WarpX (32Px1T), WarpX (4Px8T)

# Intel Advisor

- **Includes Roofline Automation…**
  - ✓ Automatically instruments applications (one dot per loop nest/function)
  - ✓ Computes FLOPS and AI for each function
  - ✓ AVX-512 support that incorporates masks
  - ✓ **Integrated Cache Simulator[1] (hierarchical roofline / multiple AI's)**
  - ✓ Automatically benchmarks target system (calculates ceilings)
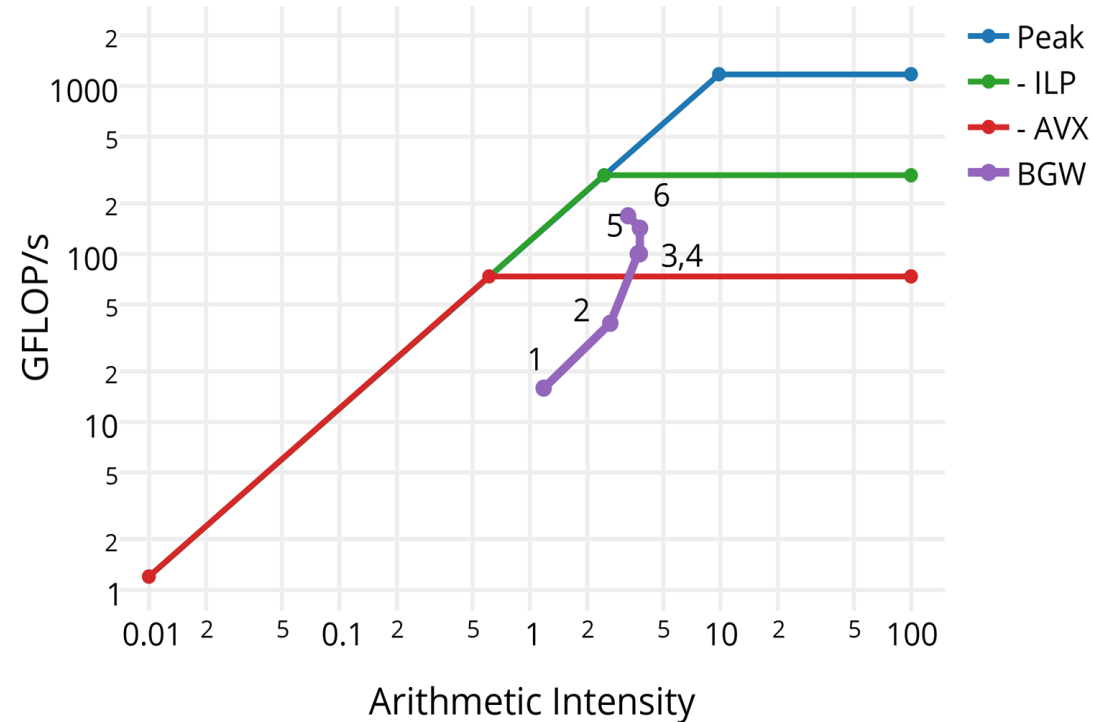  - ✓ Full integration with existing Advisor capabilities

  http://www.nersc.gov/users/training/events/roofline-training-1182017-1192017
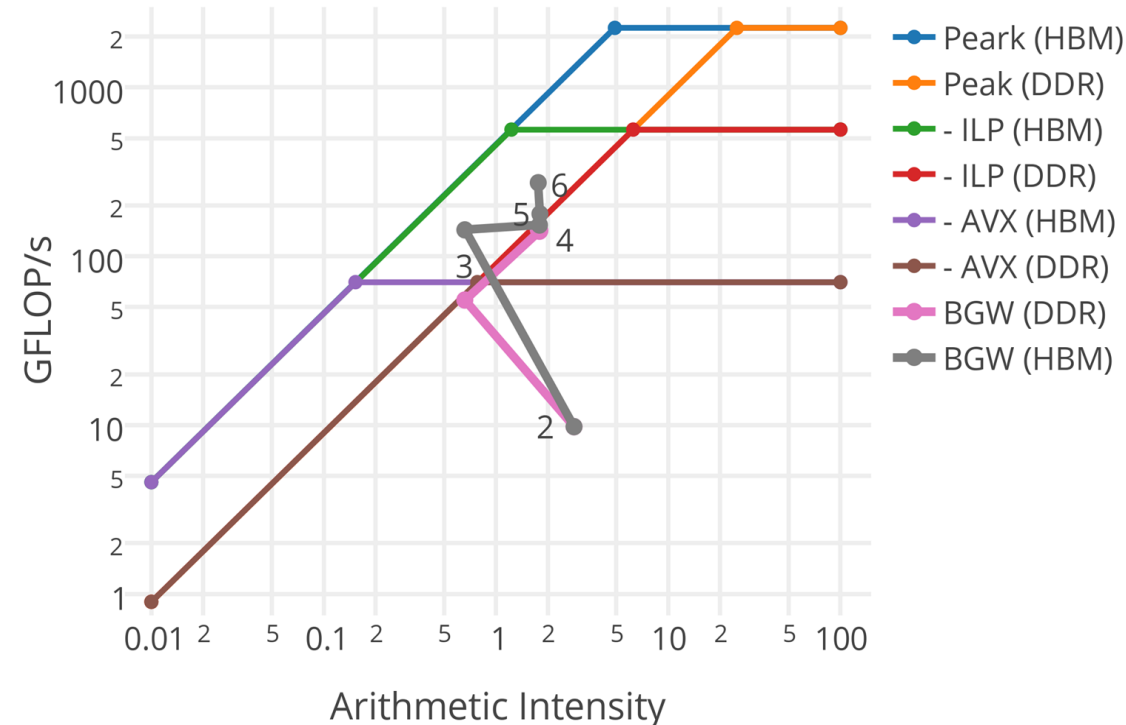


[1]Technology Preview, not in official product roadmap so far.

# Tracking Optimization Behavior

Haswell Roofline Optimization Path



KNL Roofline Optimization Path



BerkeleyGW: Optimization process for Kernel-C (Sigma code):

1. Refactor (3 Loops for MPI, OpenMP, Vectors)
2. Add OpenMP
3. Initial Vectorization (loop reordering, conditional removal)
4. Cache-Blocking
5. Improved Vectorization (Divides)
6. Hyper-threading

*BerkeleyGW is a material science application that is dominated by dense linear algebra, including distributed matrix multiplication, inversion, diagonalization, and contraction and fast fourier transforms (FFT).*

# Acknowledgements