# Introduction to the Roofline Model

## Samuel Williams

**Computational Research Division**
**Lawrence Berkeley National Lab**
SWWilliams@lbl.gov

# Acknowledgements

# 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.

3

# 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)

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

BERKELEY LAB

# Roofline Model

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

  - Tracks rates not times

  - Augmented with Little's Law

    (concurrency = latency*bandwidth)

  - Independent of ISA and architecture

    (applies to CPUs, GPUs, Google TPUs[1], etc…)

- Three Components:

  - Machine Characterization

    (realistic performance potential of the system)

  - Application Execution Monitoring

  - Theoretical Application Bounds

    *(how well could my app perform with perfect*

    *compilers, caches, overlap, …)*
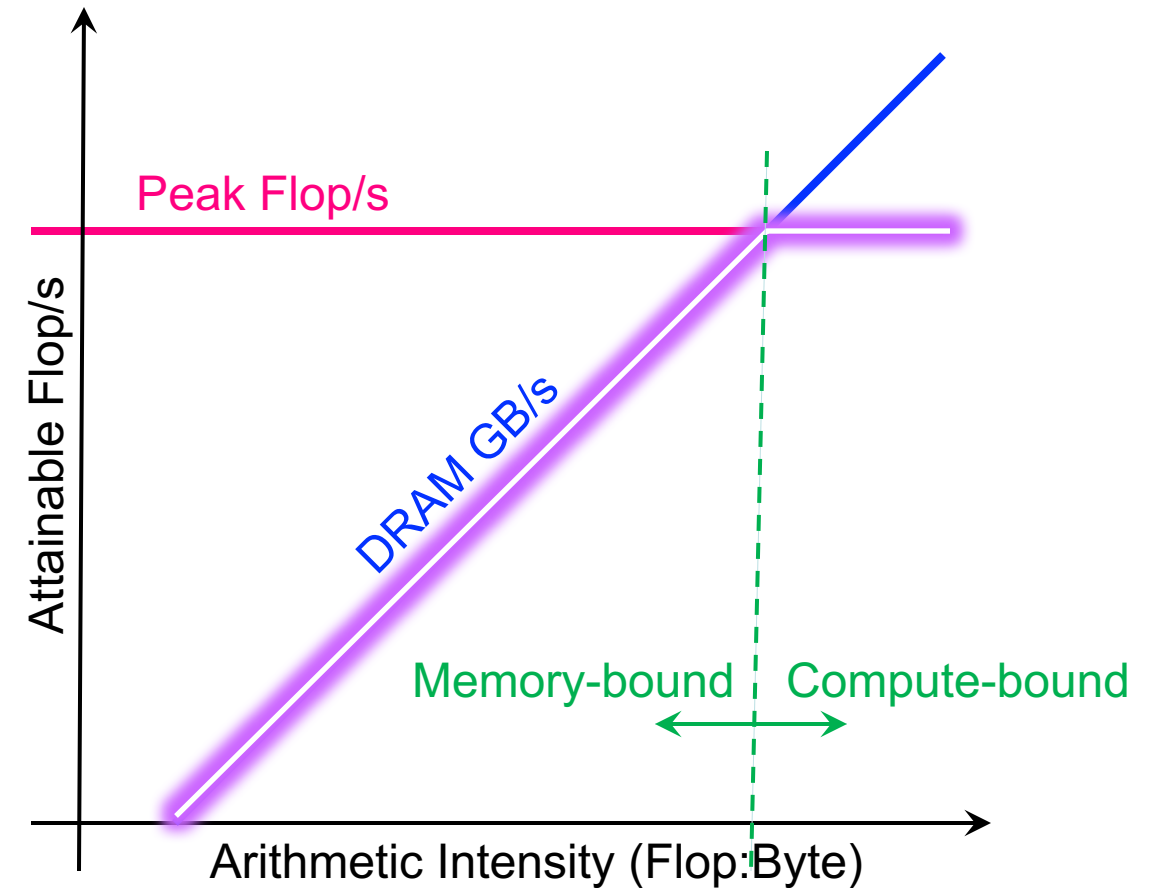


https://crd.lbl.gov/departments/computer-science/PAR/research/roofline

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

5

# (DRAM) Roofline

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

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

- Consider idealized processor/caches

- Plot the performance bound using Arithmetic Intensity (AI) as the x-axis…

  - AI = Flops / Bytes presented to DRAM
  - **Attainable Flop/s = min( peak Flop/s,  AI * peak GB/s )**
  - **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…)



Peak Flop/s

Attainable Flop/s

DRAM GB/s

Memory-bound | Compute-bound

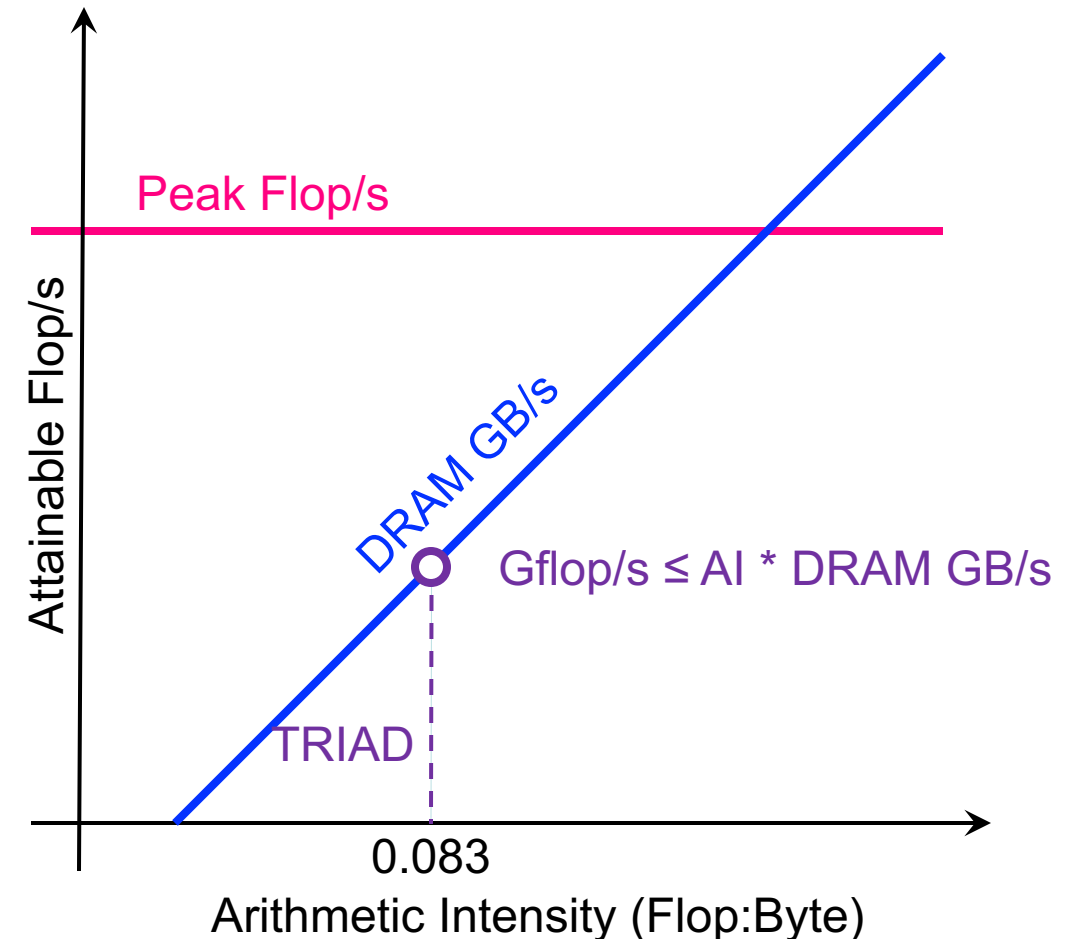Arithmetic Intensity (Flop:Byte)

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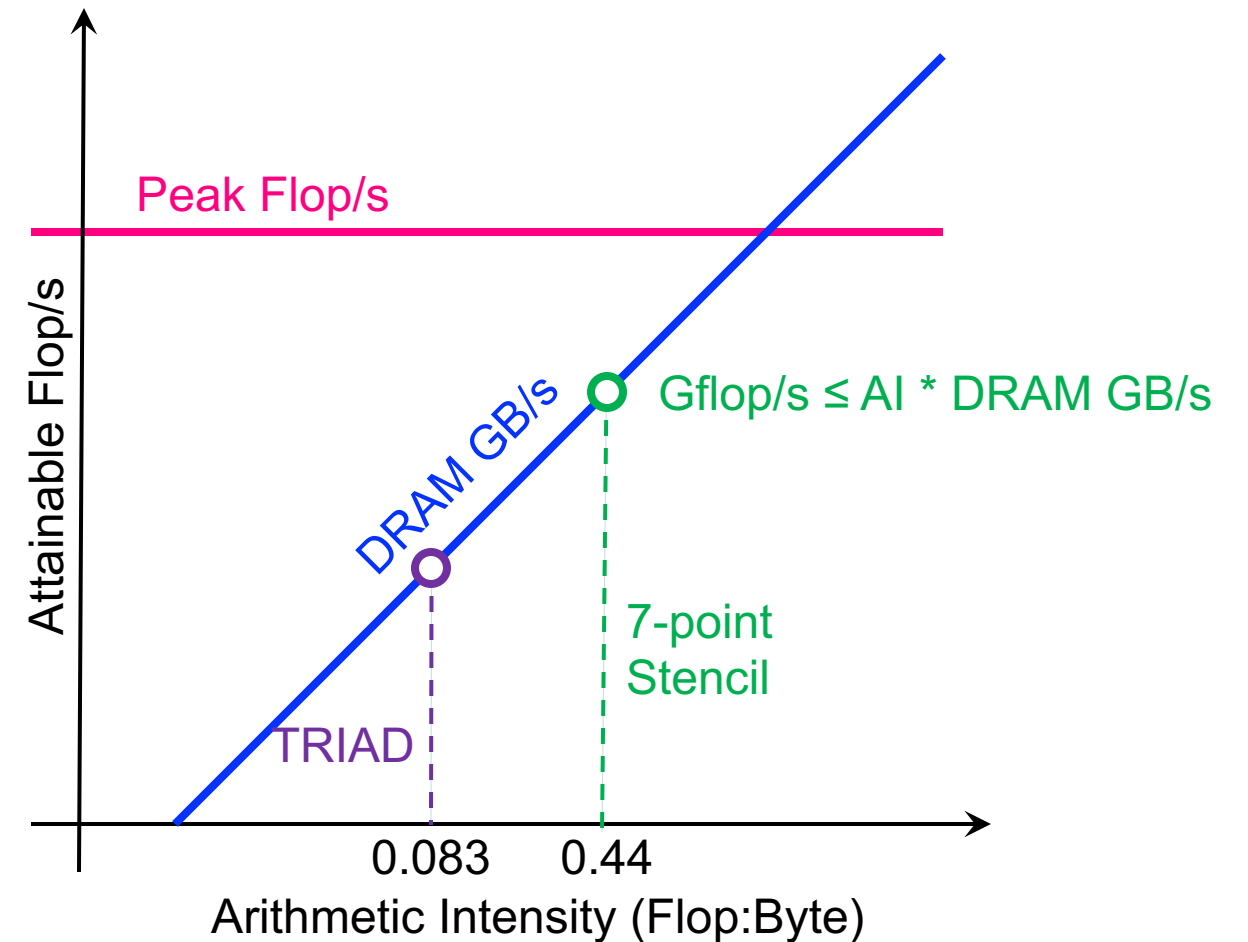
  - 2 flops per iteration
  - Transfer 24 bytes per iteration (read X[i], Y[i], write Z[i])
  - **AI = 0.083 flops per byte == Memory bound**



Peak Flop/s

Attainable Flop/s

DRAM GB/s

Gflop/s ≤ AI * DRAM GB/s

TRIAD

0.083

Arithmetic Intensity (Flop:Byte)

# 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++){
  int ijk = i + j*jStride + k*kStride;
  new[ijk] = -6.0*old[ijk         ]
                + old[ijk-1       ]
                + old[ijk+1       ]
                + old[ijk-jStride]
                + old[ijk+jStride]
                + old[ijk-kStride]
                + old[ijk+kStride];
}}}
```
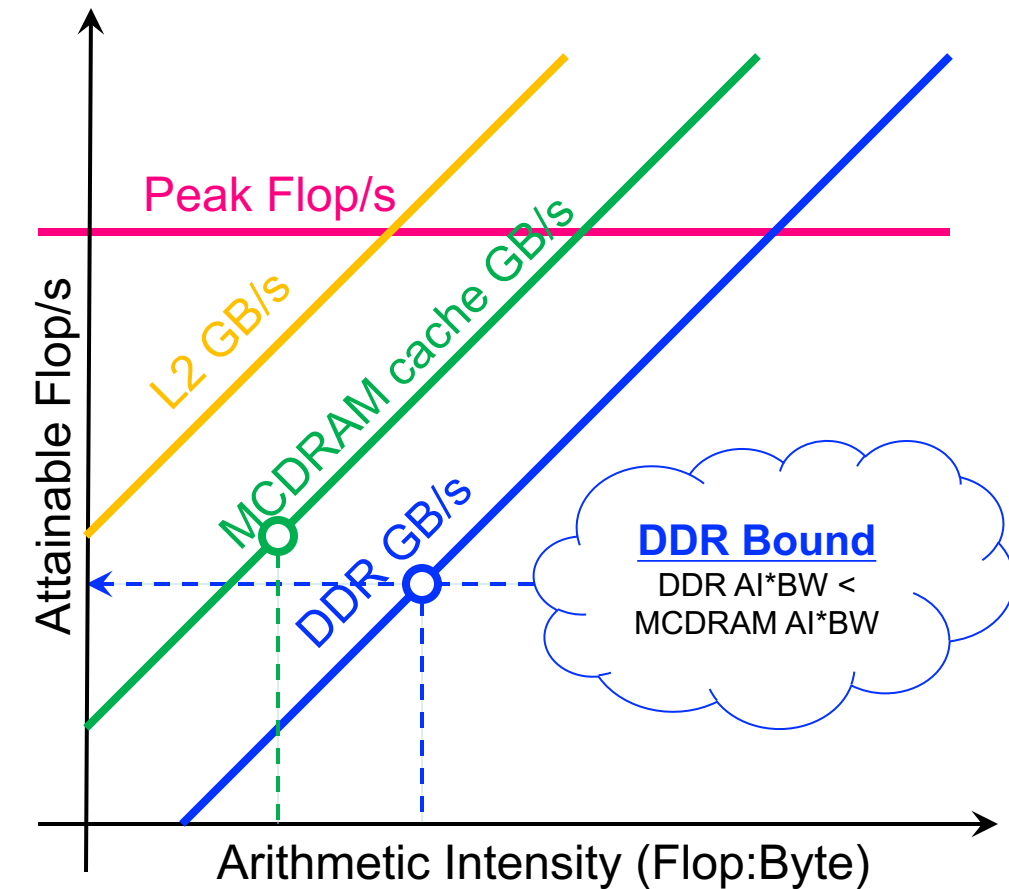
# Hierarchical Roofline

- 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
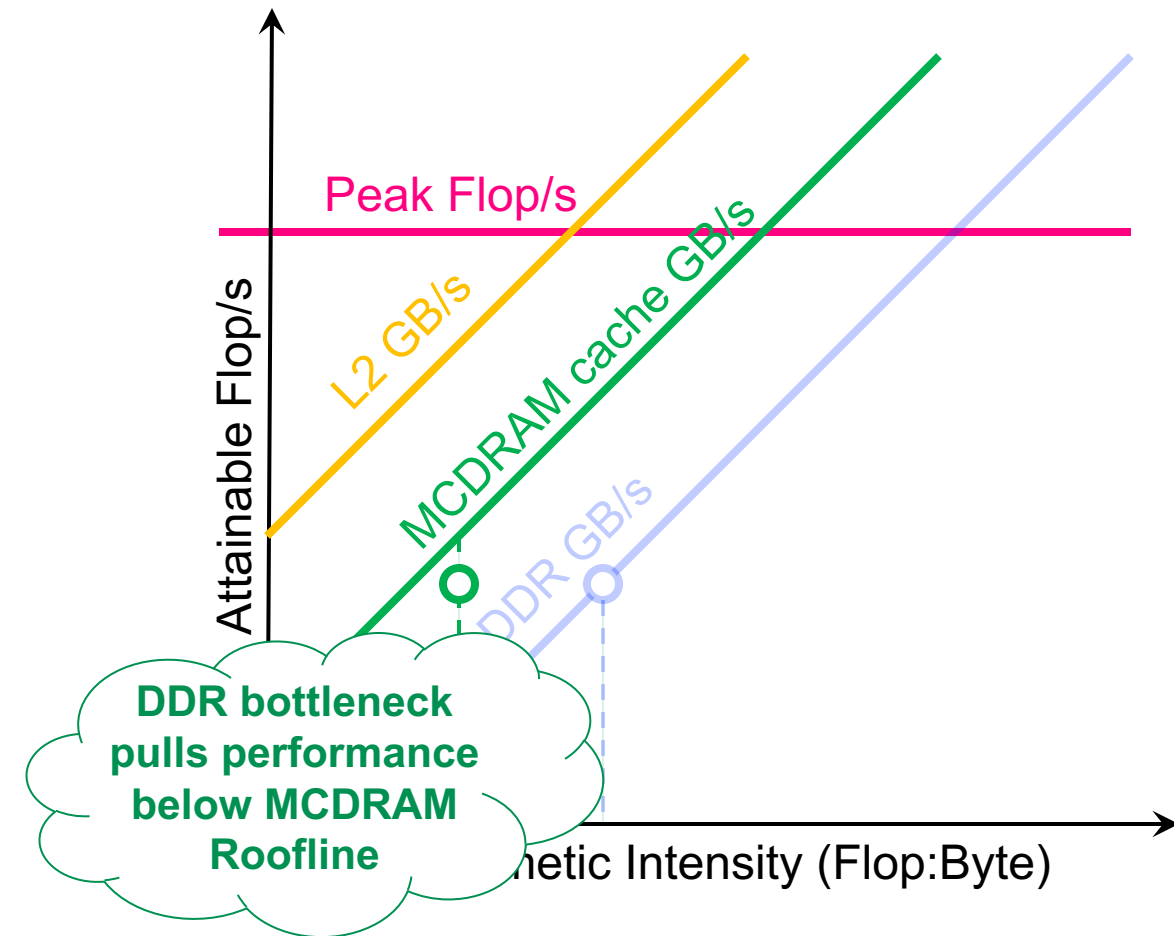  - Moreover, each level will have a unique bandwidth

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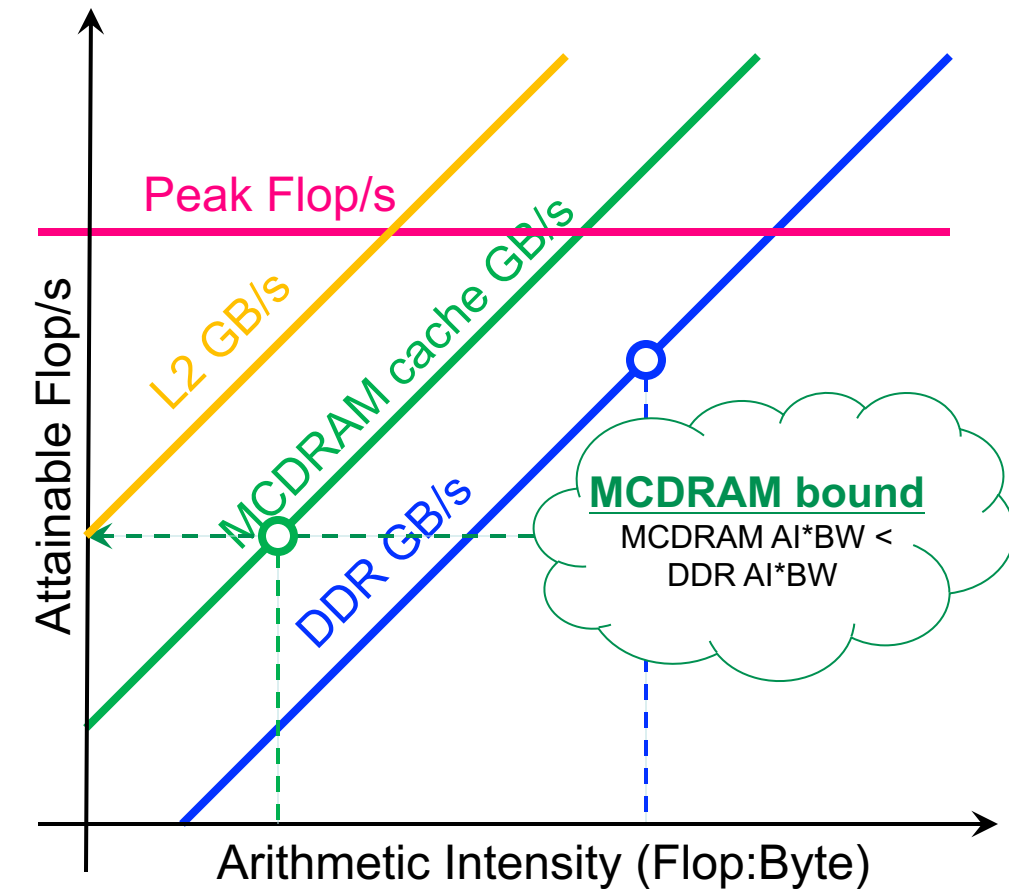
  - **… performance is bound by the minimum**

# 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)…

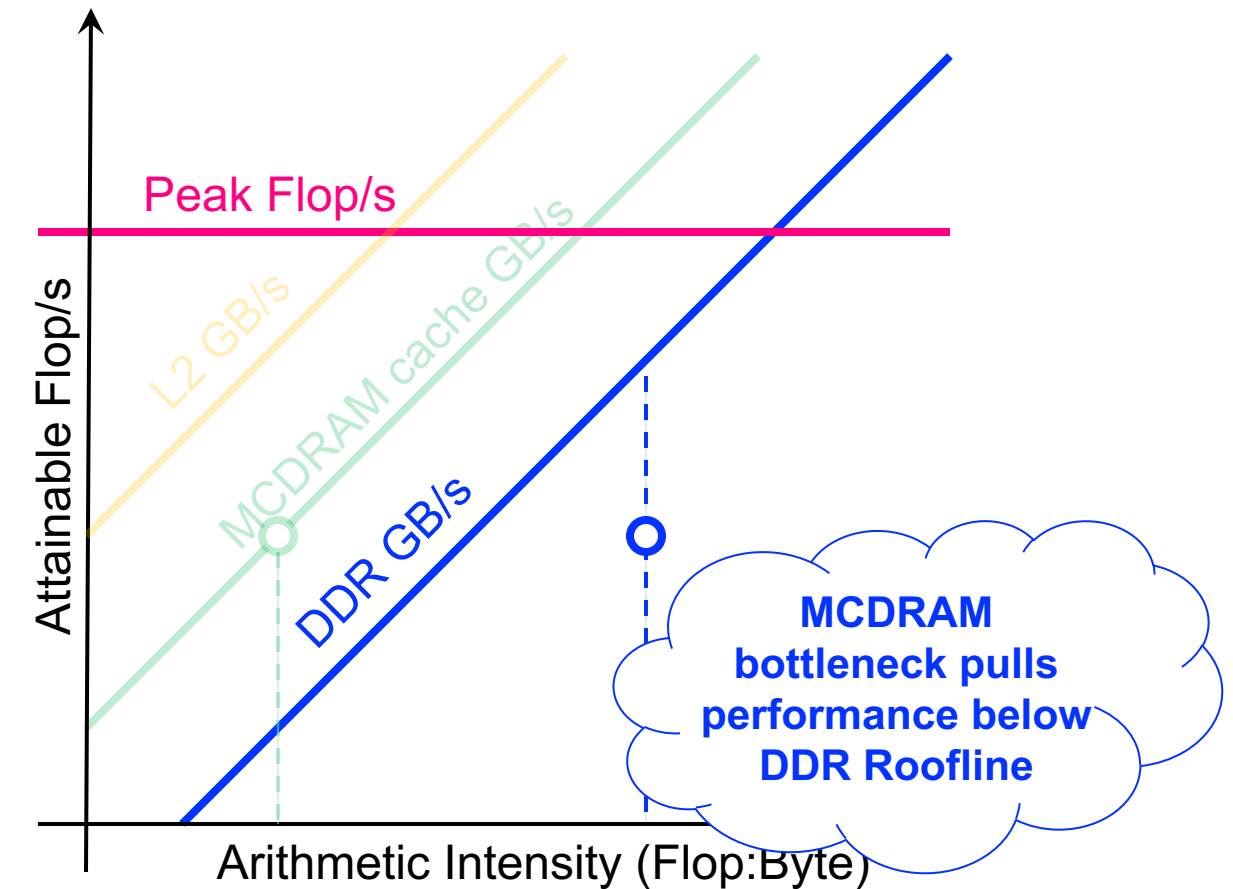  - **… performance is bound by the minimum**

# 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**
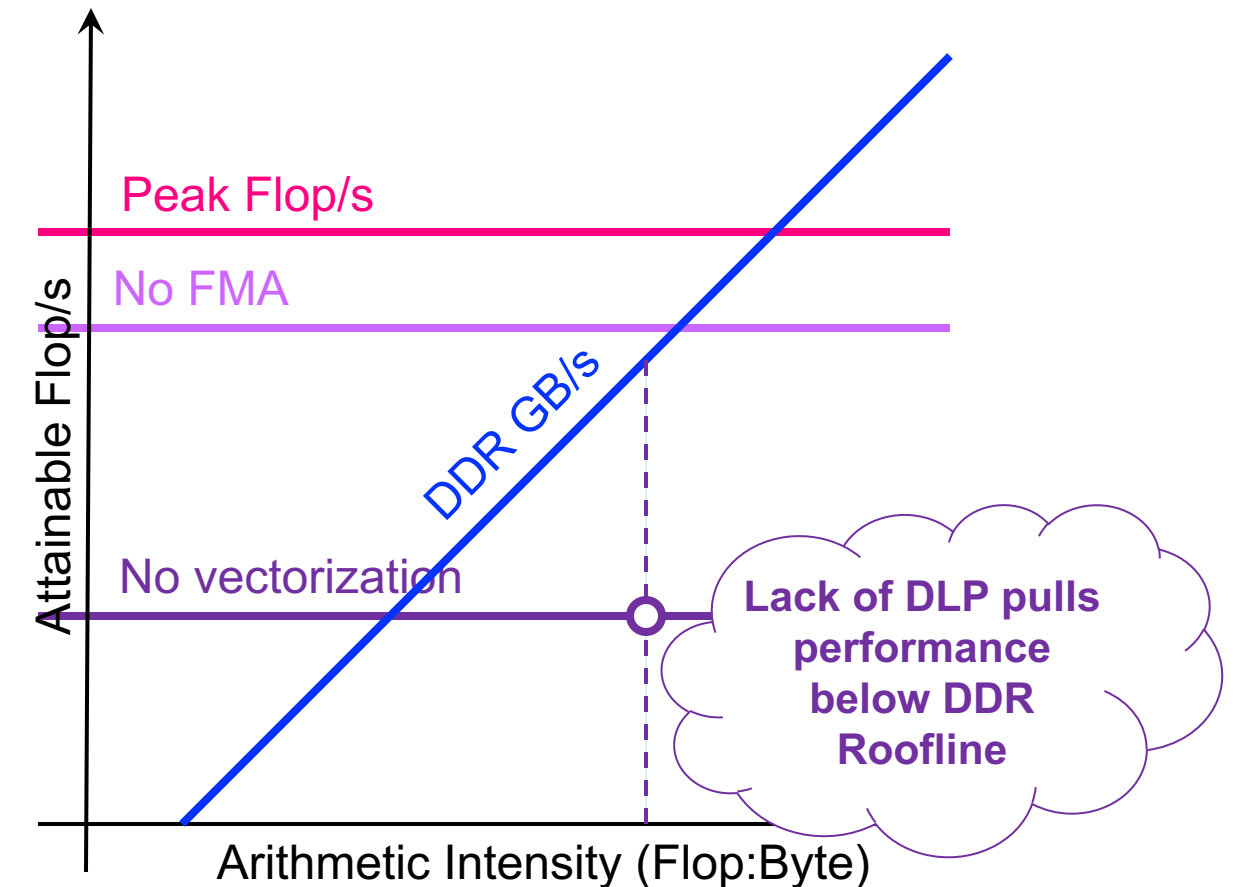


Peak Flop/s

L2 GB/s

MCDRAM cache GB/s

DDR GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

**MCDRAM bottleneck pulls performance below DDR Roofline**

BERKELEY LAB

# Data, Instruction, Thread-Level Parallelism…

- We have assumed one can attain peak flops with high locality.

- In reality, this is premised on sufficient…

  - Use special instructions (e.g. fused multiply-add)

  - Vectorization (16 flops per instruction)

  - unrolling, out-of-order execution (hide FPU latency)

  - OpenMP across multiple cores

- Without these, …

  - Peak performance is not attainable

  - Some kernels can transition from memory-bound to compute-bound

  - n.b. in reality, DRAM bandwidth is often tied to DLP and TLP (single core can't saturate BW w/scalar code)

Peak Flop/s

No FMA

DDR GB/s

No vectorization

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

Lack of DLP pulls performance below DDR Roofline

BERKELEY LAB

# Initial LBL/NERSC Roofline Efforts
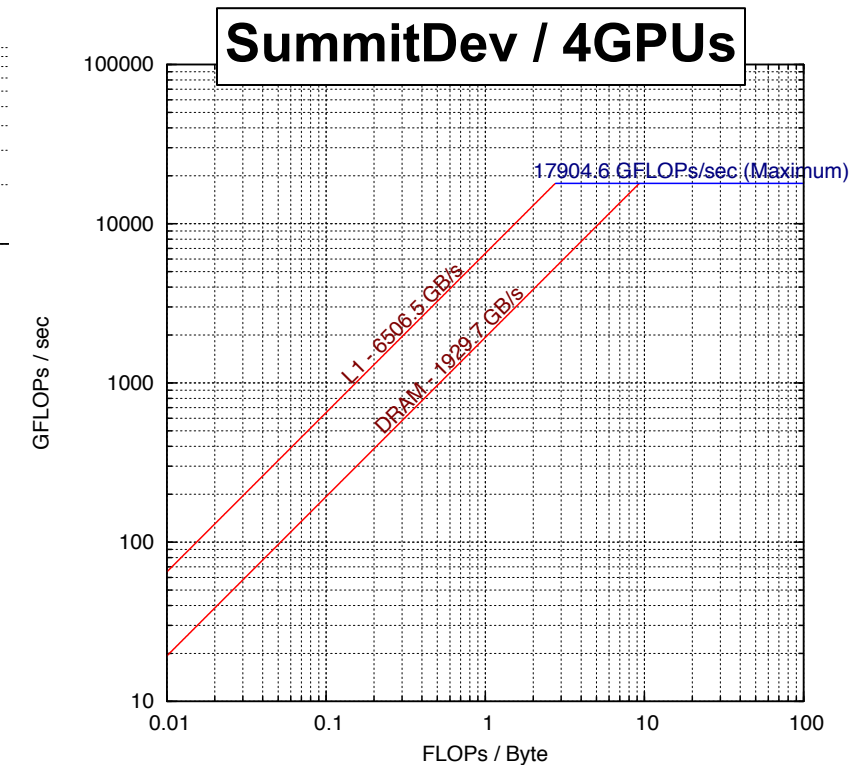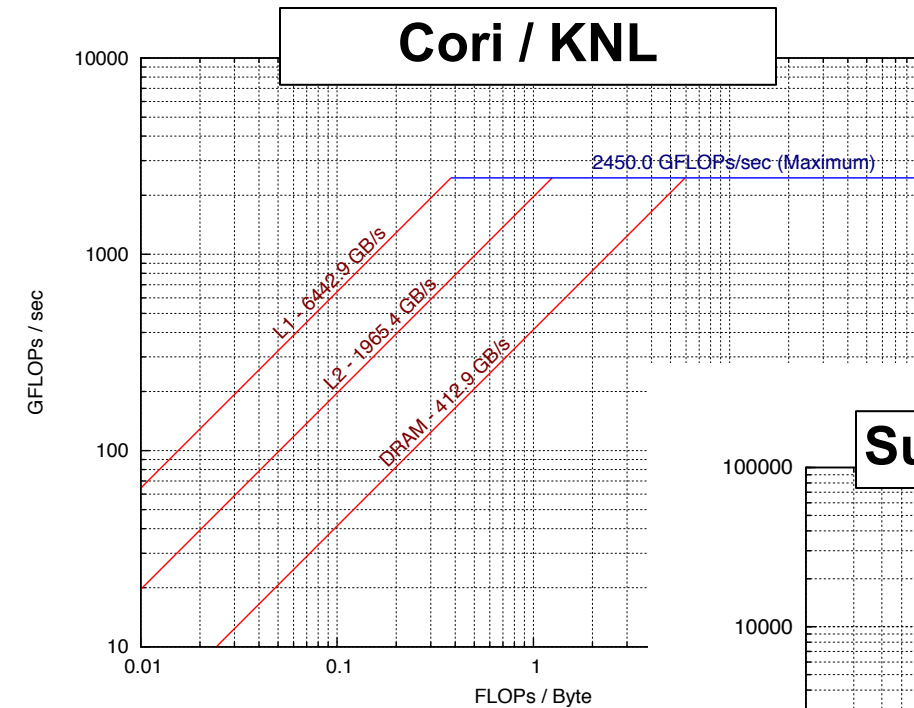
# Initial LBL Roofline Efforts / Goals

1. Node Characterization

2. Application Instrumentation/Characterization

3. Using Roofline to drive application performance analysis and optimization for KNL.

BERKELEY LAB

# Node Characterization?

- **"Marketing Numbers"** can be deceptive…
  - TurboMode / Underclock for AVX
  - Pin BW vs. real bandwidth
  - 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**
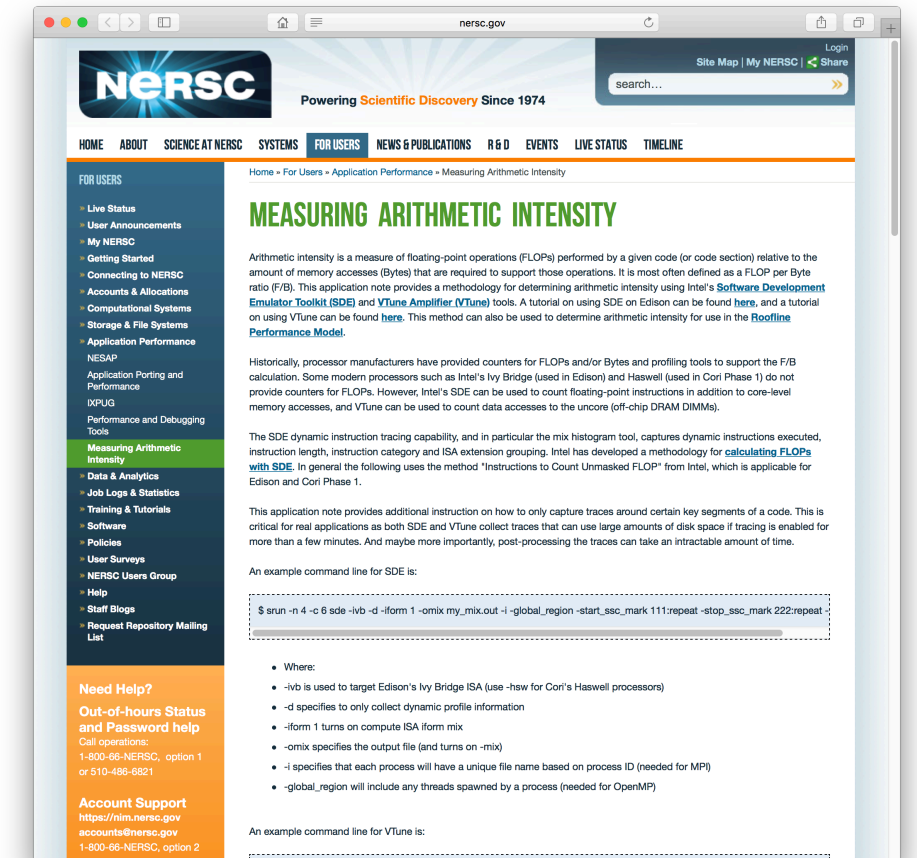


Cori / KNL



SummitDev / 4GPUs

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

BERKELEY LAB

# Instrumentation with Performance Counters?

- ***Characterizing applications with performance counters can be problematic…***

  ✘ Flop Counters can be broken/missing in production processors

  ✘ Vectorization/Masking can complicate counting Flop's

  ✘ Counting Loads and Stores doesn't capture cache reuse while counting cache misses doesn't account for prefetchers.

  ✘ DRAM counters (Uncore PMU) might be accurate, but…

    • are privileged and thus nominally inaccessible in user mode

    • may need vendor (e.g. Cray) and center (e.g. NERSC) approved OS/kernel changes

BERKELEY LAB

# 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)

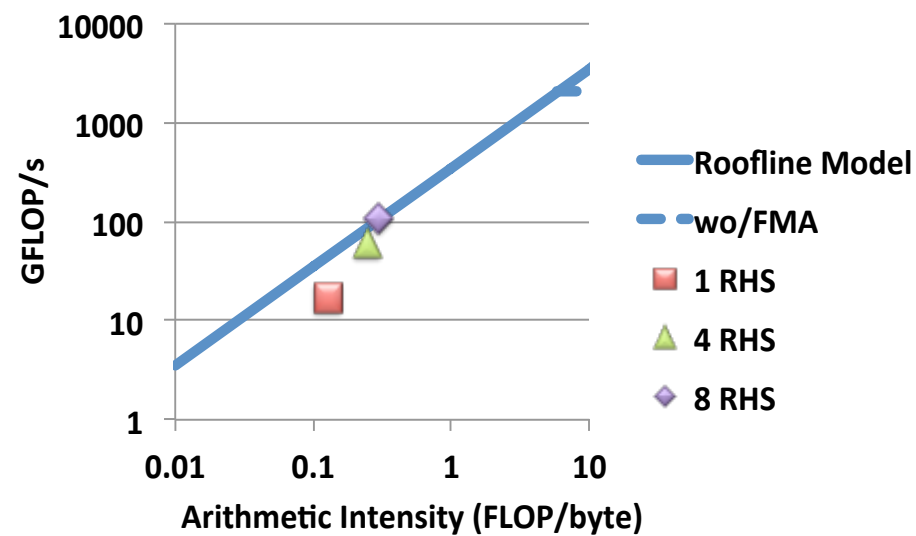- ➢ Used by NESAP (NERSC KNL application readiness project) to characterize apps on Cori…



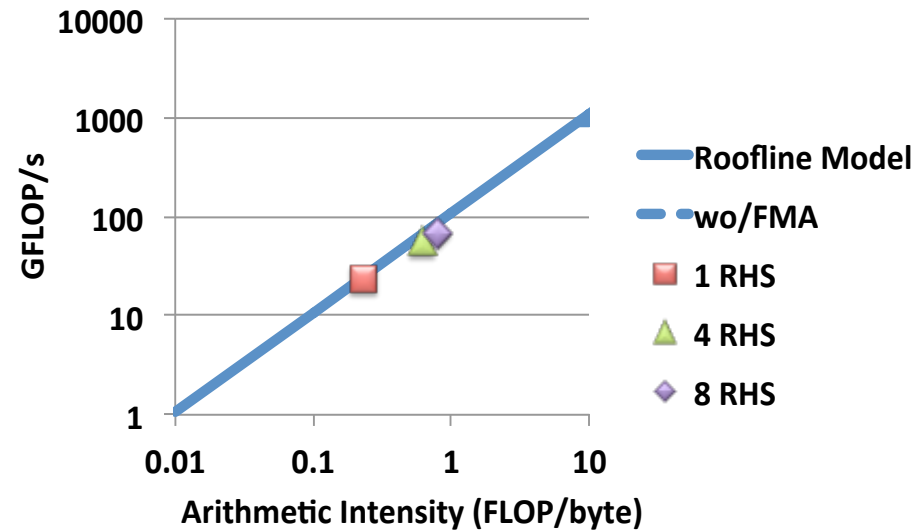http://www.nersc.gov/users/application-performance/measuring-arithmetic-intensity/

BERKELEY LAB
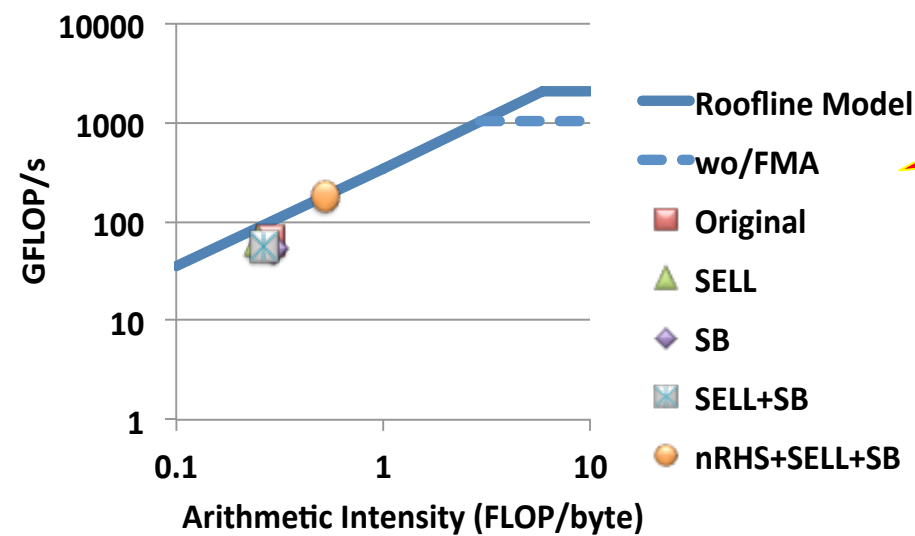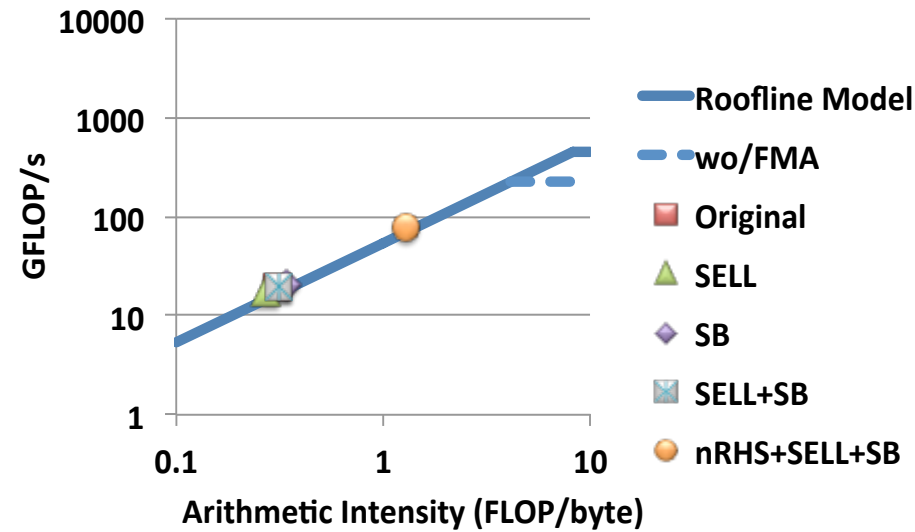
# Initial Roofline Analysis of NESAP Codes

## MFDn

## EMGeo

## PICSAR

**2P HSW**

**KNL**



DRAM-only Roofline is insufficient for PICSAR

BERKELEY LAB

# 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**



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

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

# Need an integrated solution…

- Having to compose VTune, SDE, and plotting tools…
  - ✓ worked correctly and benefited NESAP's application readiness
  - ✗ forced users to learn/run multiple tools and manually parse/graph the output
  - ✗ forced users to instrument routines of interest in their application
  - ✗ lacked integration with compiler/debugger/disassembly

- LIKWID was…
  - ✓ fast and easy to use
  - ✗ Suffered from the same limitations as VTune/SDE

- ERT…
  - ✓ Characterized flops, and bandwidths (cache, DRAM)
  - ✓ Interoperable with MPI, OpenMP, and CUDA
  - ✗ Required users to manually parse/incorporate the output

BERKELEY LAB

# Intel Advisor

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



Memory-bound, invest into cache blocking etc

Compute bound: invest into SIMD…

[1]Technology Preview, not in official product roadmap so far.
This version will be made available during the hands-on component of this tutorial.

BERKELEY LAB

# Hierarchical Roofline vs. Cache-Aware Roofline

*…understanding different Roofline formulations in Advisor*

# There are two Major Roofline Formulations:

- Hierarchical Roofline (original Roofline w/ DRAM, L3, L2, …)…

  - **Williams, et al, "Roofline: An Insightful Visual Performance Model for Multicore Architectures", CACM, 2009**

  - **Chapter 4 of "Auto-tuning Performance on Multicore Computers", 2008**

  - Defines multiple bandwidth ceilings and multiple AI's per kernel

  - Performance bound is the minimum of flops and the memory intercepts (superposition of original, single-metric Rooflines)

- Cache-Aware Roofline

  - **Ilic et al, "Cache-aware Roofline model: Upgrading the loft", IEEE Computer Architecture Letters, 2014**

  - Defines multiple bandwidth ceilings, but uses a single AI (flop:L1 bytes)

  - As one looses cache locality (capacity, conflict, …) performance falls from one BW ceiling to a lower one at constant AI

- Why Does this matter?

  - Some tools use the Hierarchical Roofline, some use cache-aware **== Users need to understand the differences**

  - Cache-Aware Roofline model was integrated into production Intel Advisor

  - Evaluation version of Hierarchical Roofline[1] (cache simulator) has also been integrated into Intel Advisor

  - **You will be allowed to explore both in the hand—on component of this tutorial**

[1]Technology Preview, not in official product roadmap so far.
This version will be made available during the hands-on component of this tutorial.

# Hierarchical Roofline

- Captures cache effects

- AI is Flop:Bytes after being *filtered by lower cache levels*

- Multiple Arithmetic Intensities

  (one per level of memory)

- AI *dependent* on problem size

  (capacity misses reduce AI)

- Memory/Cache/Locality effects are *observed as decreased AI*

- Requires *performance counters or cache simulator* to correctly measure AI

# Cache-Aware Roofline

- Captures cache effects

- AI is Flop:Bytes **as *presented to the L1 cache (plus non-temporal stores)***

- Single Arithmetic Intensity

- AI *independent* of problem size

- Memory/Cache/Locality effects are *observed as decreased performance*

- Requires static analysis or *binary instrumentation* to measure AI

BERKELEY LAB

# Example: STREAM

- ## L1 AI…

  - 2 flops
  - 2 x 8B load (old)
  - 1 x 8B store (new)
  - = 0.08 flops per byte

- ## No cache reuse…

  - Iteration i doesn't touch any data associated with iteration i+delta for any delta.

- ## … leads to a DRAM AI equal to the L1 AI

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

BERKELEY LAB

# Example: STREAM

## Hierarchical Roofline



Performance is bound to the minimum of the two Intercepts…

$AI_{L1}$ * L1 GB/s
$AI_{DRAM}$ * DRAM GB/s

Multiple AI's….
1) Flop:DRAM bytes
2) Flop:L1 bytes (same)

Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

0.083

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Observed performance is correlated with DRAM bandwidth

Single AI based on flop:L1 bytes

Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

0.083

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

# Example: 7-point Stencil (Small Problem)

- ## L1 AI…

  - 7 flops

  - 7 x 8B load (old)

  - 1 x 8B store (new)

  - = 0.11 flops per byte

  - some compilers may do register shuffles to reduce the number of loads.

- ## Moderate cache reuse…

  - old[ijk] is reused on subsequent iterations of i,j,k

  - old[ijk-1] is reused on subsequent iterations of i.

  - old[ijk-jStride] is reused on subsequent iterations of j.

  - old[ijk-kStride] is reused on subsequent iterations of k.
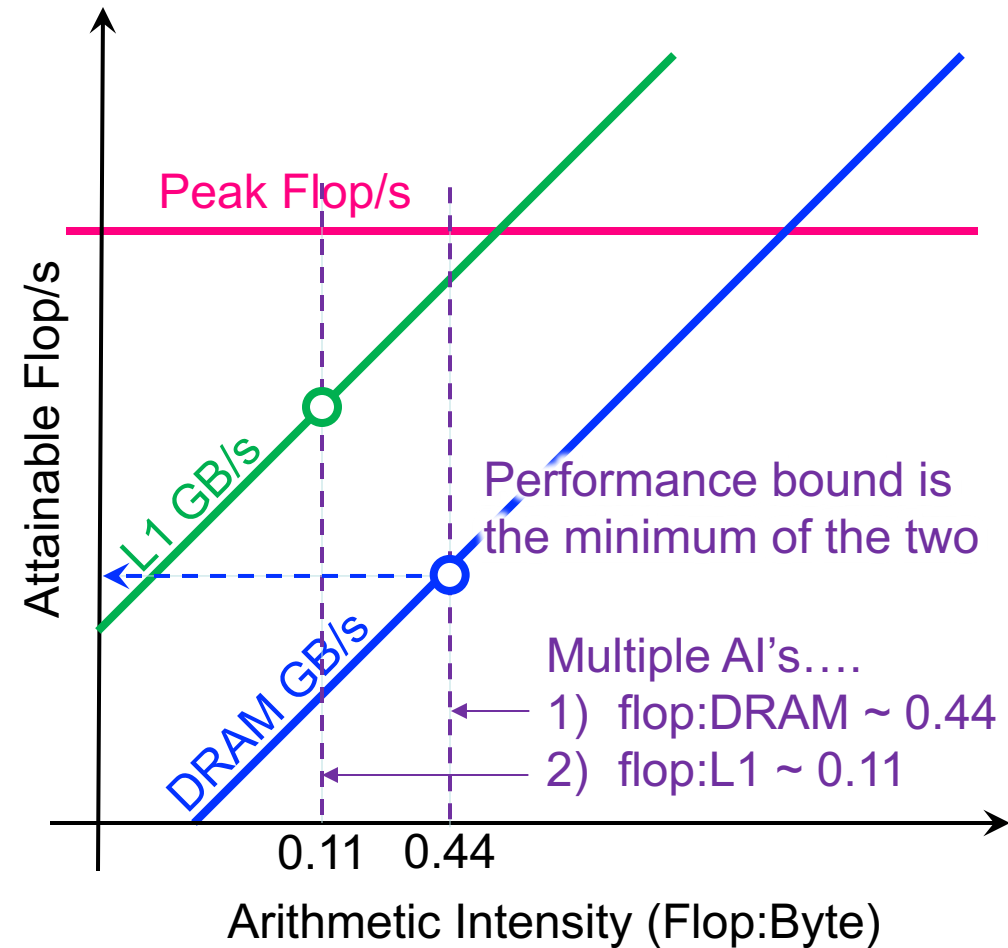
- ## … leads to DRAM AI larger than the L1 AI

```
#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++){
   int ijk = i + j*jStride + k*kStride;
   new[ijk] = -6.0*old[ijk        ]
                 + old[ijk-1      ]
                 + old[ijk+1      ]
                 + old[ijk-jStride]
                 + old[ijk+jStride]
                 + old[ijk-kStride]
                 + old[ijk+kStride];
}}}
```
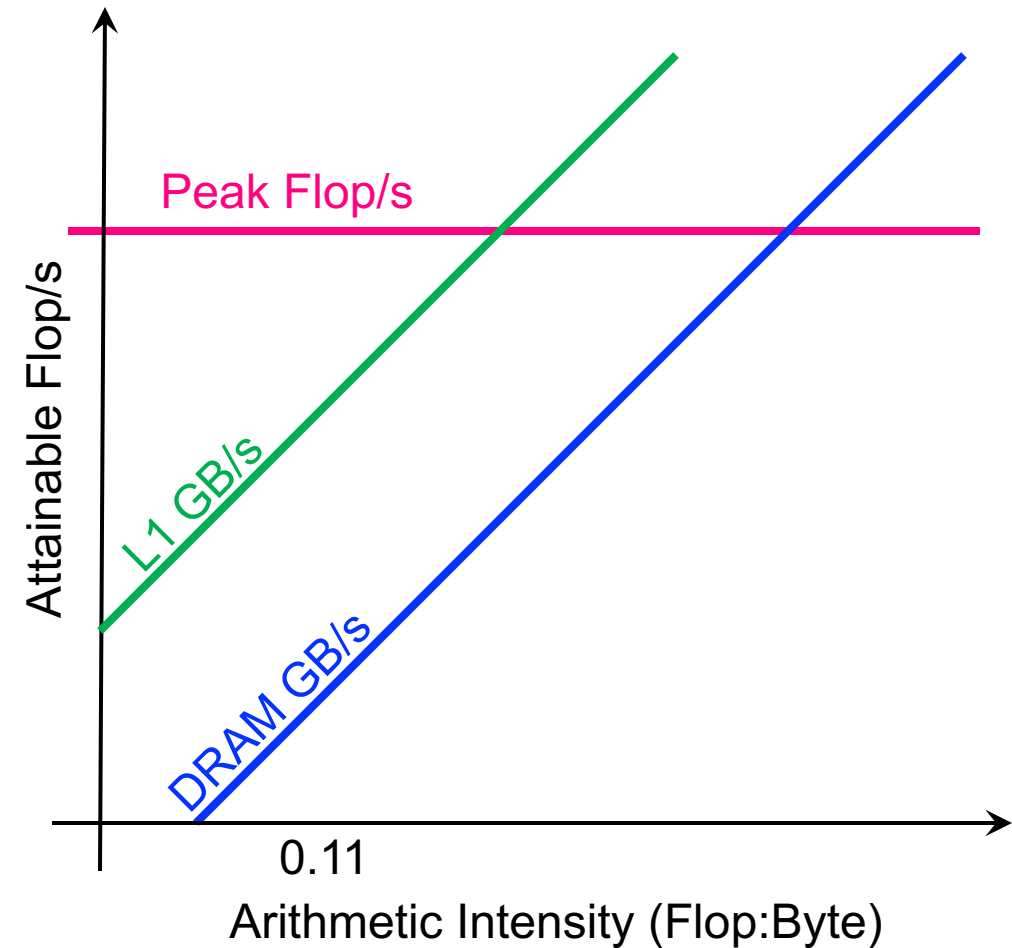
BERKELEY LAB

# Example: 7-point Stencil (Small Problem)

## Hierarchical Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

Performance bound is the minimum of the two

Multiple AI's....
1) flop:DRAM ~ 0.44
2) flop:L1 ~ 0.11

0.11  0.44

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

0.11

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

# Example: 7-point Stencil (Small Problem)

## Hierarchical Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Performance bound is
the minimum of the two

Multiple AI's….
1) flop:DRAM ~ 0.44
2) flop:L1 ~ 0.11

0.11   0.44

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Observed performance
is between L1 and DRAM lines
(== some cache locality)

Single AI based on flop:L1 bytes

0.11

Arithmetic Intensity (Flop:Byte)

# Example: 7-point Stencil (Large Problem)

## Hierarchical Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Capacity misses reduce DRAM AI and performance

Multiple AI's….
1) flop:DRAM ~ 0.20
2) flop:L1 ~ 0.11

0.11    0.20

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Observed performance is closer to DRAM line (== less cache locality)

Single AI based on flop:L1 bytes

0.11

Arithmetic Intensity (Flop:Byte)

# Example: 7-point Stencil (Observed Perf.)

## Hierarchical Roofline



Peak Flop/s

L1 GB/s

Attainable Flop/s

DRAM GB/s

Actual **observed** performance is tied to the bottlenecked resource and can be well below a cache Roofline (e.g. L1).

0.11   0.20

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

L1 GB/s

Attainable Flop/s

DRAM GB/s

Observed performance is closer to DRAM line (== less cache locality)

Single AI based on flop:L1 bytes

0.11

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

# Example: 7-point Stencil (Observed Perf.)

## Hierarchical Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Actual **observed** performance is tied to the bottlenecked resource and can be well below a cache Roofline (e.g. L1).

0.11  0.20

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

Attainable Flop/s

L1 GB/s

DRAM GB/s

Observed performance is closer to DRAM line (== less cache locality)

Single AI based on flop:L1 bytes

0.11

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

Questions?