# Performance Tuning of Scientific Codes with the Roofline Model

| | | |
|---|---|---|
| 1:30pm | Introductions / Administration | all |
| 1:35pm | Roofline Introduction | Samuel Williams |
| 2:10pm | CARM / Energy / GPUs | Aleksandar Ilic |
| 2:40pm | Intel Advisor Installation | Zakhar Matveev |
| 3:00pm | coffee break | |
| 3:30pm | Introduction to Intel Advisor | Zakhar Matveev |
| 3:45pm | Hands-on with Intel Advisor | all |
| 4:30pm | HPC Application Studies | Charlene Yang |
| 4:55pm | closing remarks / Q&A | all |

# Introductions

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

## Aleksandar Ilic
**Institute of Systems and Computer**
**Engineering, Portugal**
aleksandar.ilic@inesc-id.pt

## Charlene Yang
**NERSC**
**Lawrence Berkeley National Lab**
CJYang@lbl.gov

## Zakhar Matveev
**Intel Corporation**
zakhar.a.matveev@intel.com

## Kiril Rogozhin
**Intel Corporation**
kirill.rogozhin@intel.com

# Materials:
# USB / Downloads

# Don't forget to take the Survey…

http://bit.ly/sc18-eval

# Introduction to the Roofline Model

## Samuel Williams

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

# Acknowledgements

# Background

# 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.

# Performance Models

- Many different components can contribute to kernel run time.
- Some are application-specific, and some architecture-specific.

| | |
|---:|:---|
| #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 |

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 |

LogP

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

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 bandwidth |
| Depth | OMP Overhead |
| MPI Message Size | Network Bandwidth |
| MPI Send:Wait ratio | Network Gap |
| #MPI Wait's | Network Latency |

**Roofline Model**

BERKELEY LAB

# Roofline Model:
## Arithmetic Intensity and Bandwidth

# Performance Models / Simulators

- Historically, many performance models and simulators tracked time 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)

- … resulted in a shift from a latency-limited computing regime to a **throughput-limited computing regime**

BERKELEY LAB

# 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.
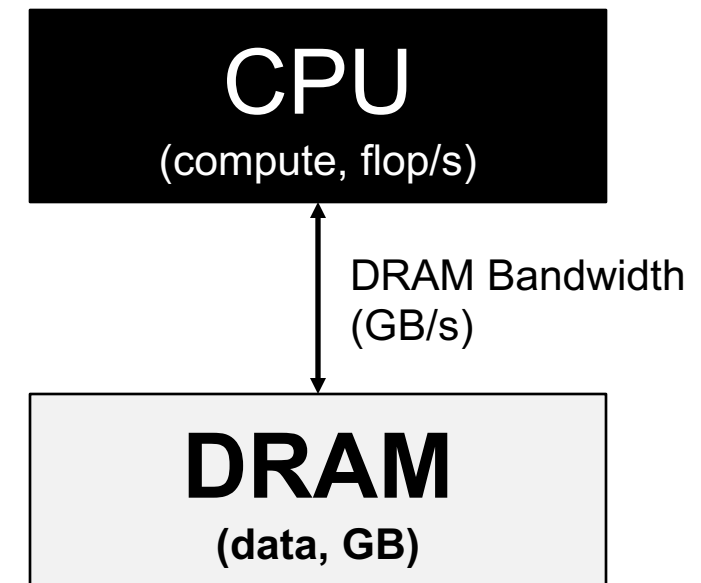
# (DRAM) Roofline

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

- However, finite reuse and bandwidth limit performance.

- Assuming perfect overlap of communication and computation…

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

DRAM Bandwidth
(GB/s)

**DRAM**
(data, GB)

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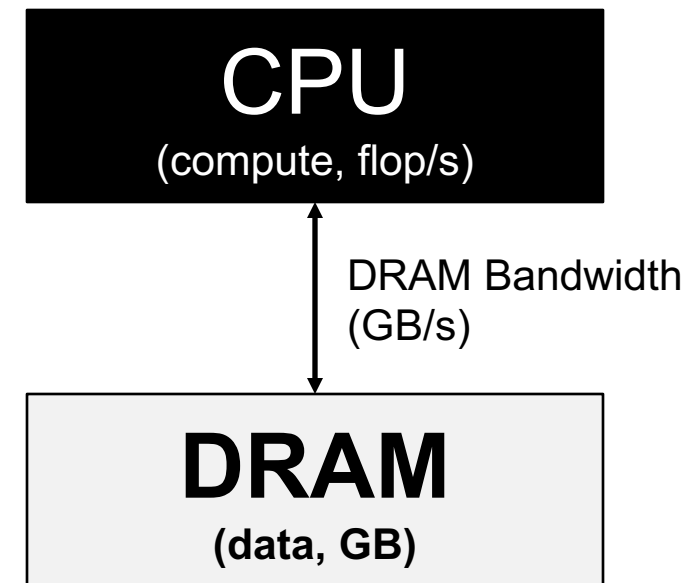
# (DRAM) Roofline

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

- However, finite reuse and bandwidth limit performance.

- Assuming perfect overlap of communication and computation…

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

DRAM Bandwidth
(GB/s)

**DRAM**
(data, GB)

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

BERKELEY LAB

# (DRAM) Roofline

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

- However, finite reuse and bandwidth limit performance.

- Assuming perfect overlap of communication and computation…

CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

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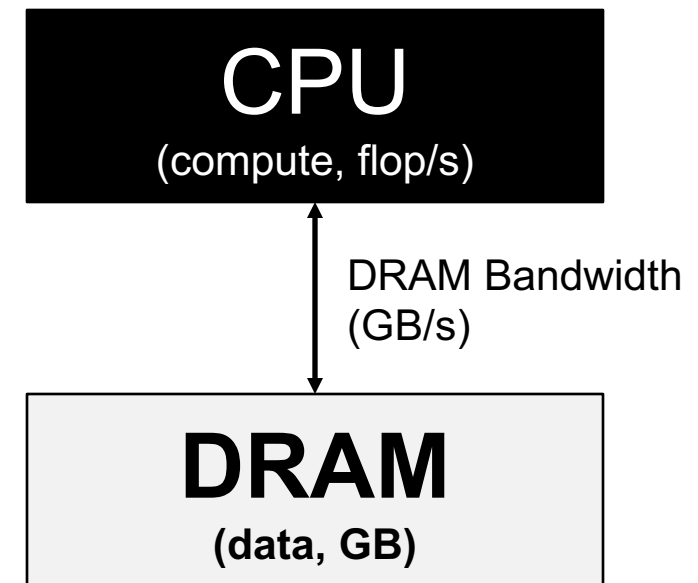
# (DRAM) Roofline

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

- However, finite reuse and bandwidth limit performance.

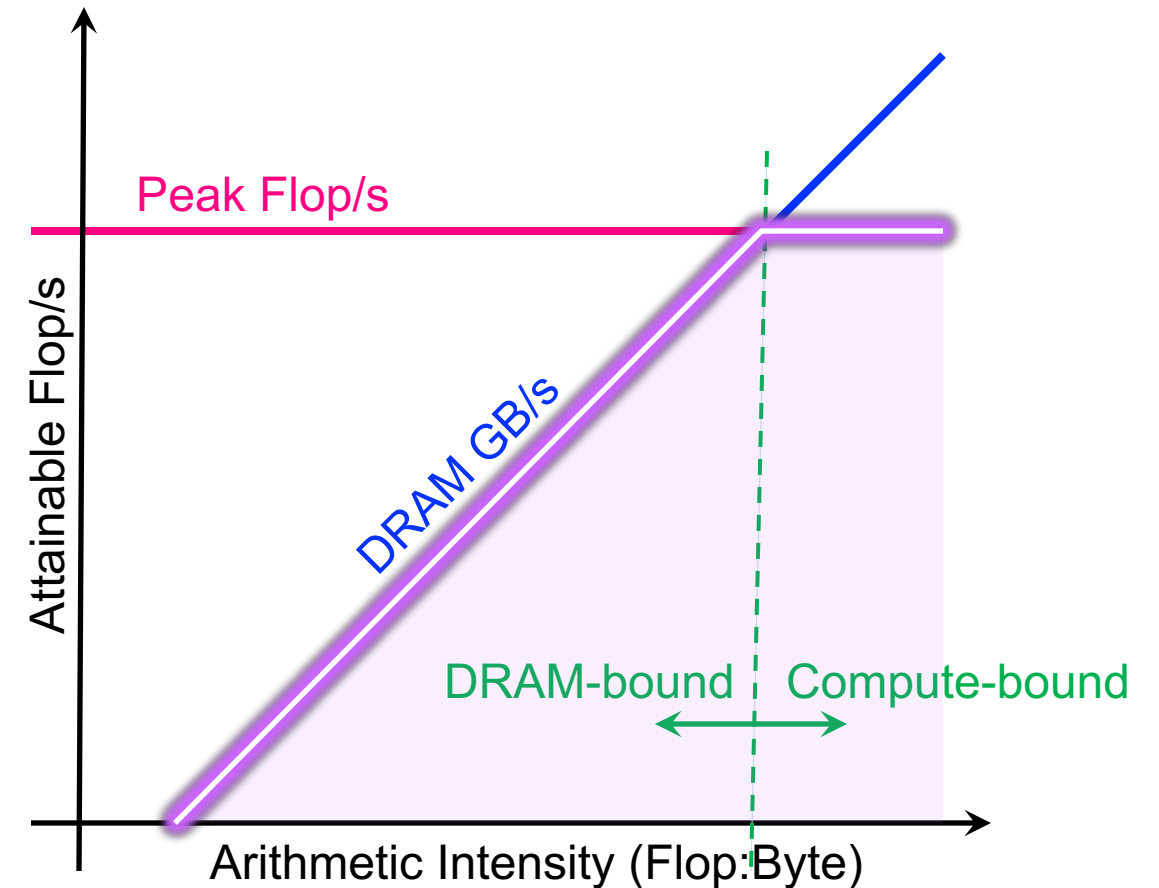- Assuming perfect overlap of communication and computation…

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

DRAM Bandwidth (GB/s)

**DRAM**
(data, GB)

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

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

BERKELEY LAB

# (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…)
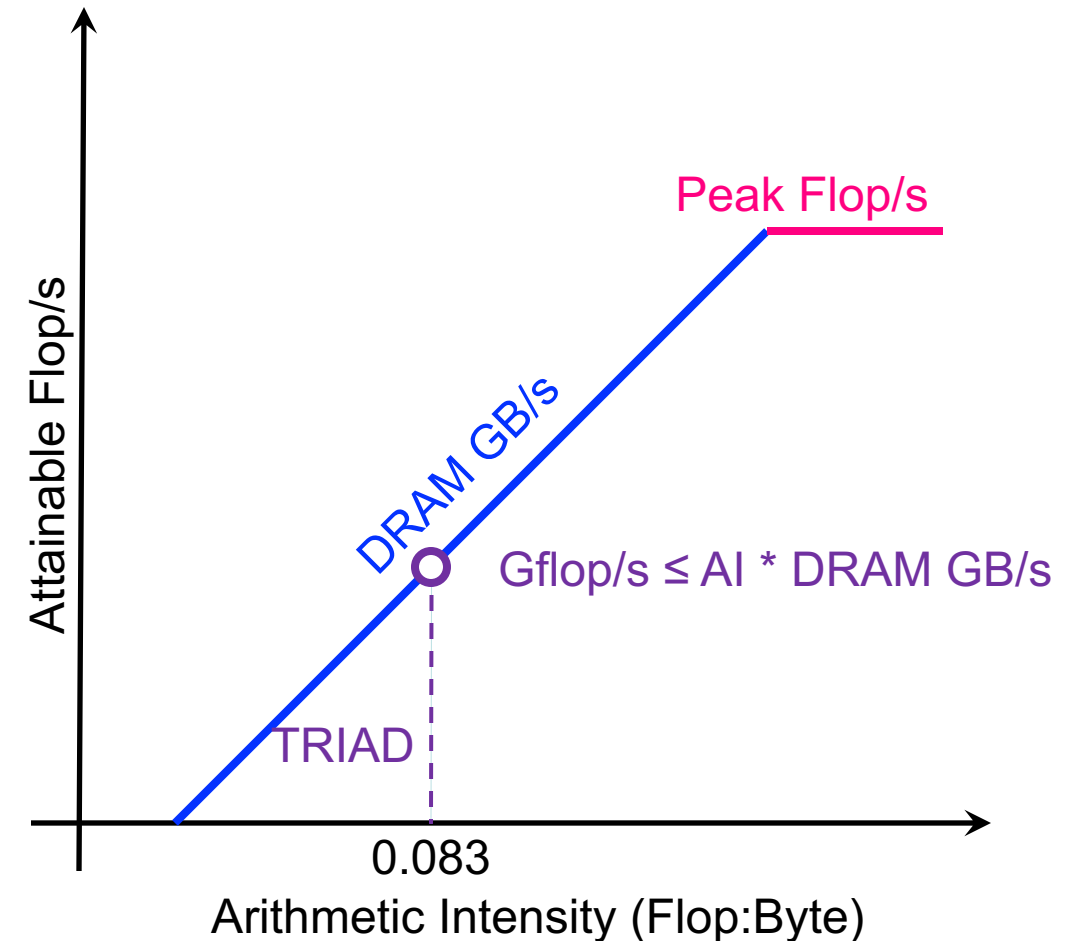
# Roofline Example #1

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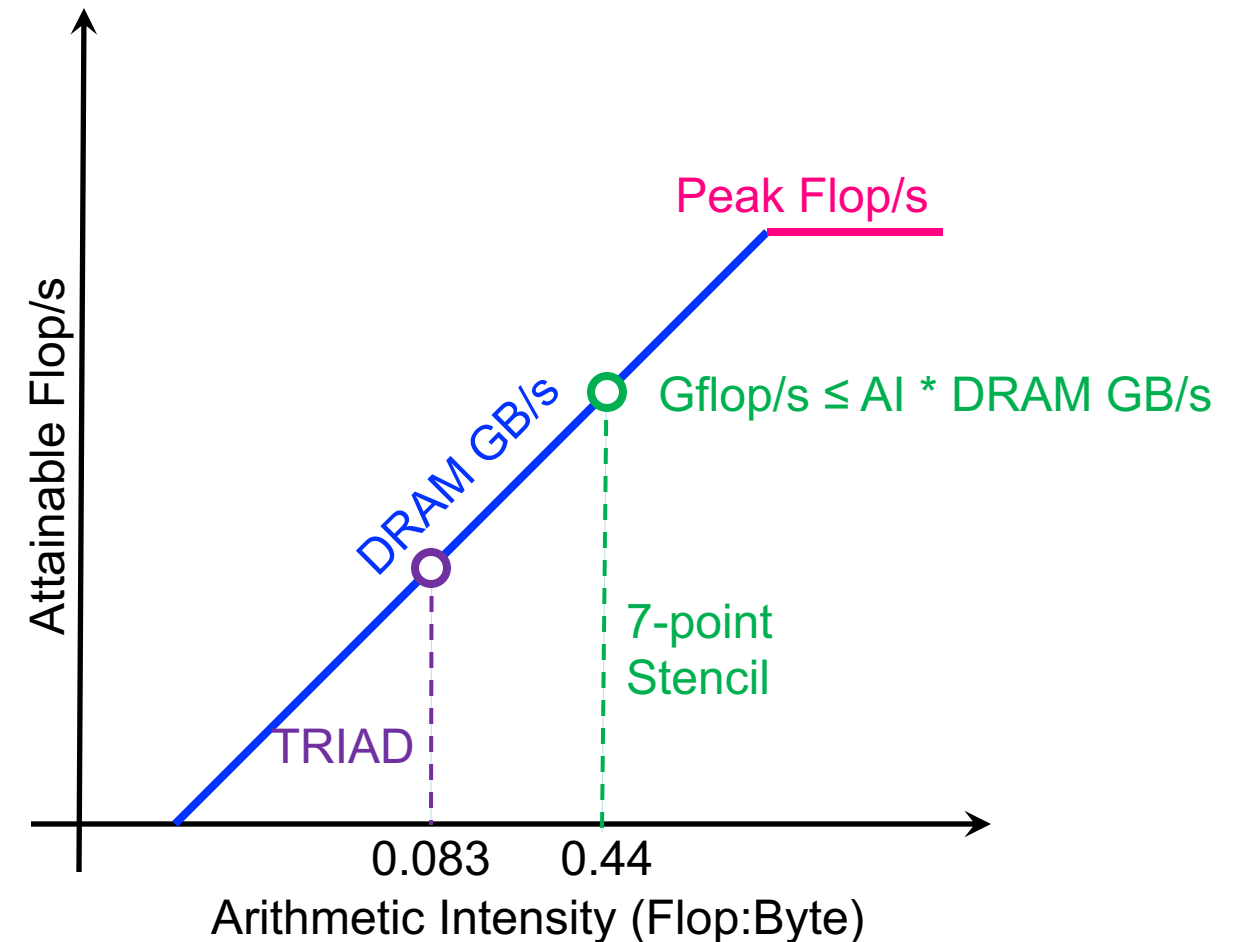
  - 2 flops per iteration
  - Transfer 24 bytes per iteration (read X[i], Y[i], write Z[i])
  - **AI = 0.083 flops per byte == Memory bound**



Peak Flop/s

DRAM GB/s

Gflop/s ≤ AI * DRAM GB/s

TRIAD

Attainable Flop/s

0.083

Arithmetic Intensity (Flop:Byte)
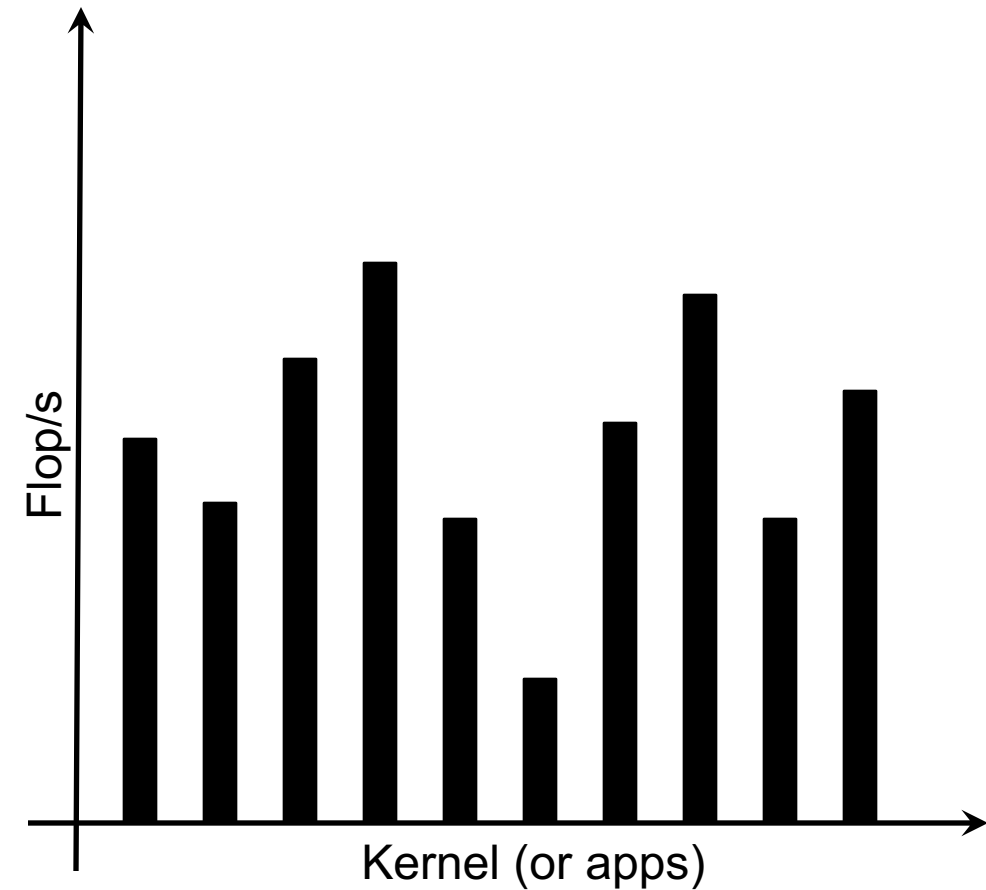
24

# 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  ];
}}}
```
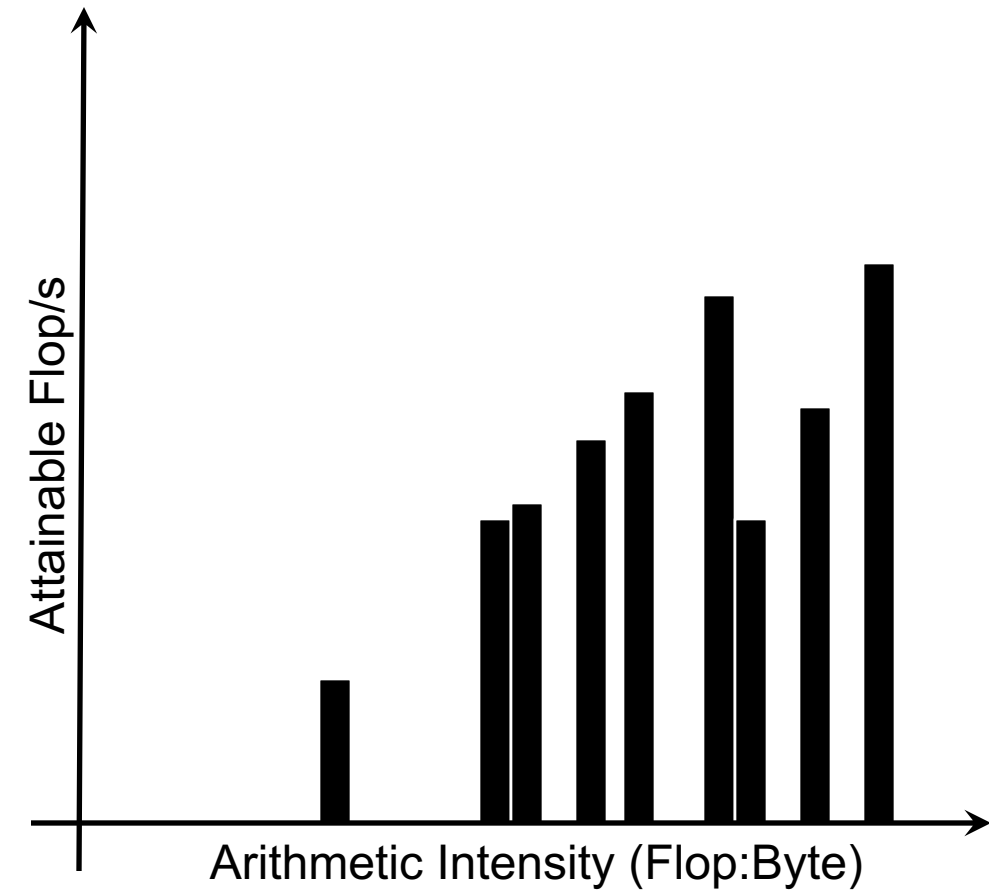


Peak Flop/s

Gflop/s ≤ AI * DRAM GB/s

DRAM GB/s

7-point Stencil

TRIAD

0.083    0.44

Arithmetic Intensity (Flop:Byte)

Attainable Flop/s

BERKELEY LAB

# Why is Roofline Useful?

- Imagine a mix of loop nests

- Flop/s alone may not be useful in deciding which to optimize first
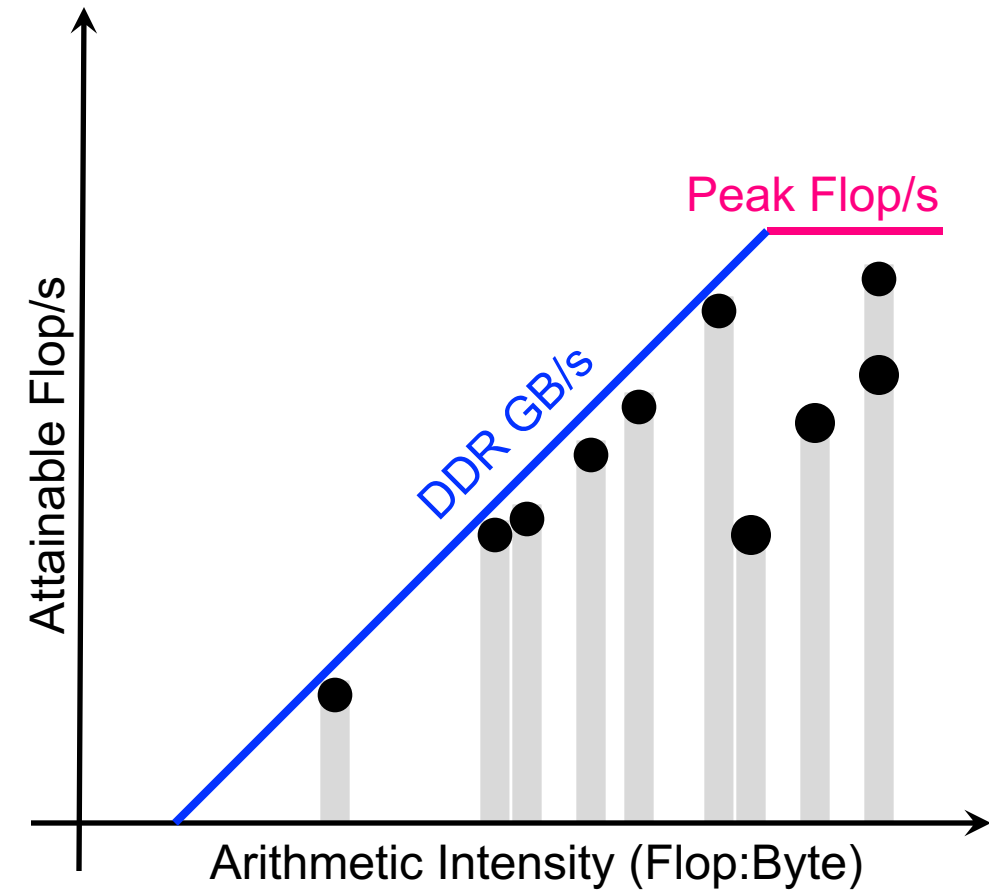


Flop/s

Kernel (or apps)

# Why is Roofline Useful?
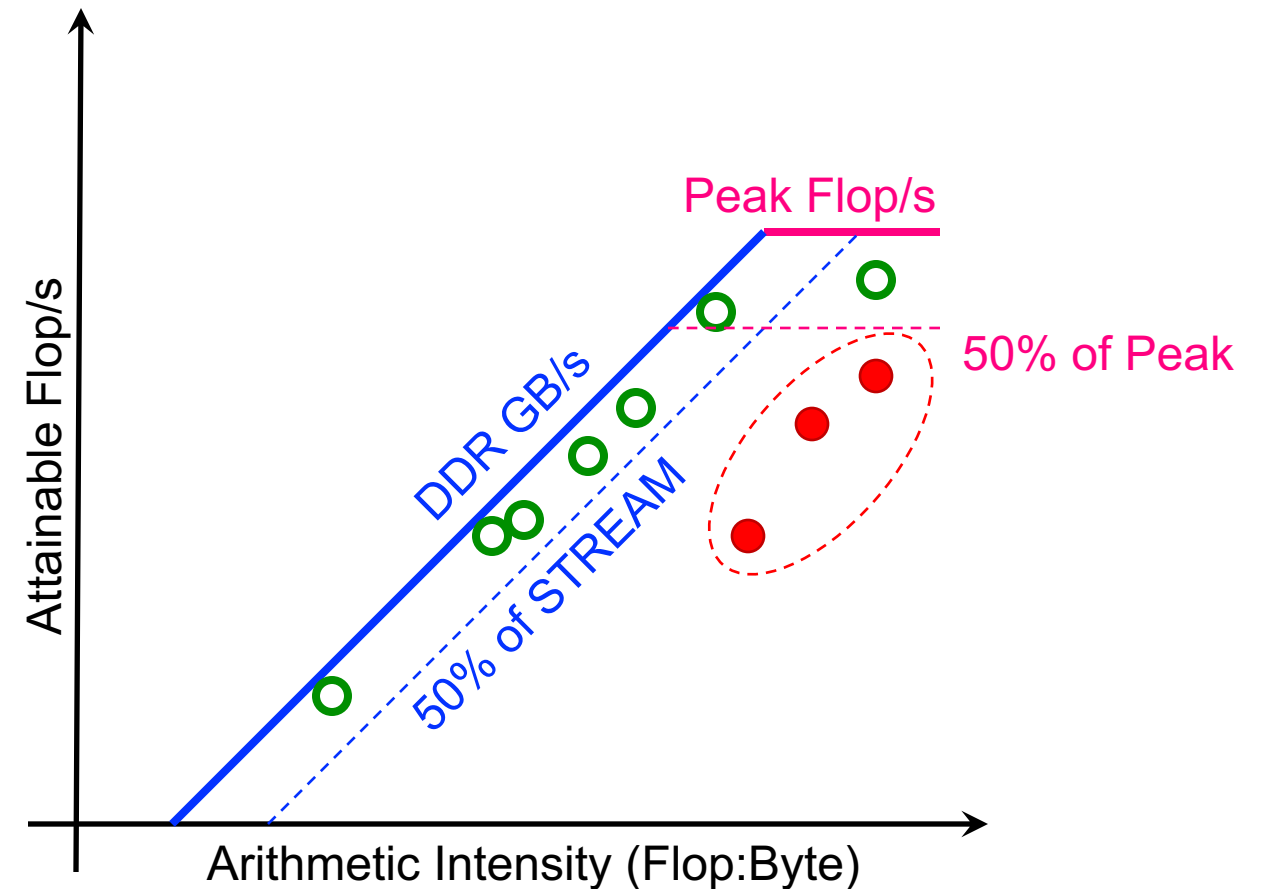
- We can sort kernels by AI …

# Why is Roofline Useful?

- We can sort kernels by AI …

- … and compare performance relative to machine capabilities

# Why is Roofline Useful?

- **Kernels near the roofline are making good use of computational resources**

  ○ kernels can have low performance (Gflop/s), but make good use of a machine

  ○ kernels can have high performance (Gflop/s), but make poor use of a machine
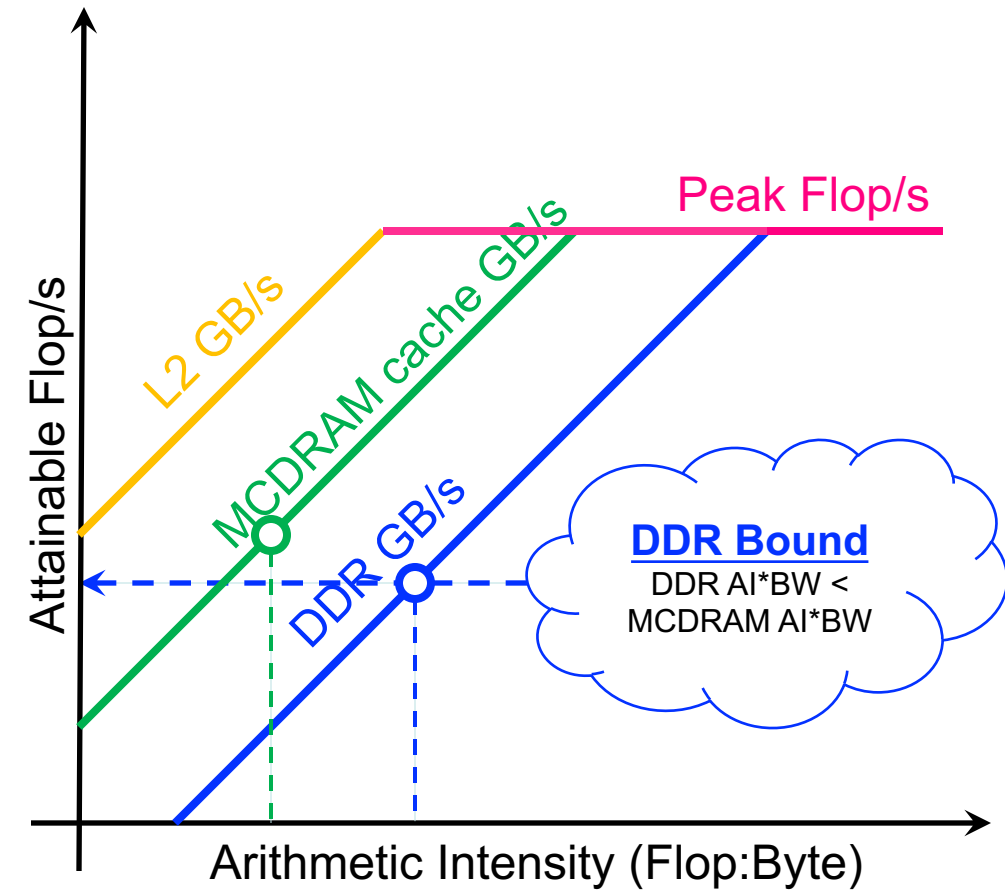
# Hierarchical Roofline

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

- **Applications 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 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 bandwidth

  - Measure AI for each level of memory

  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…
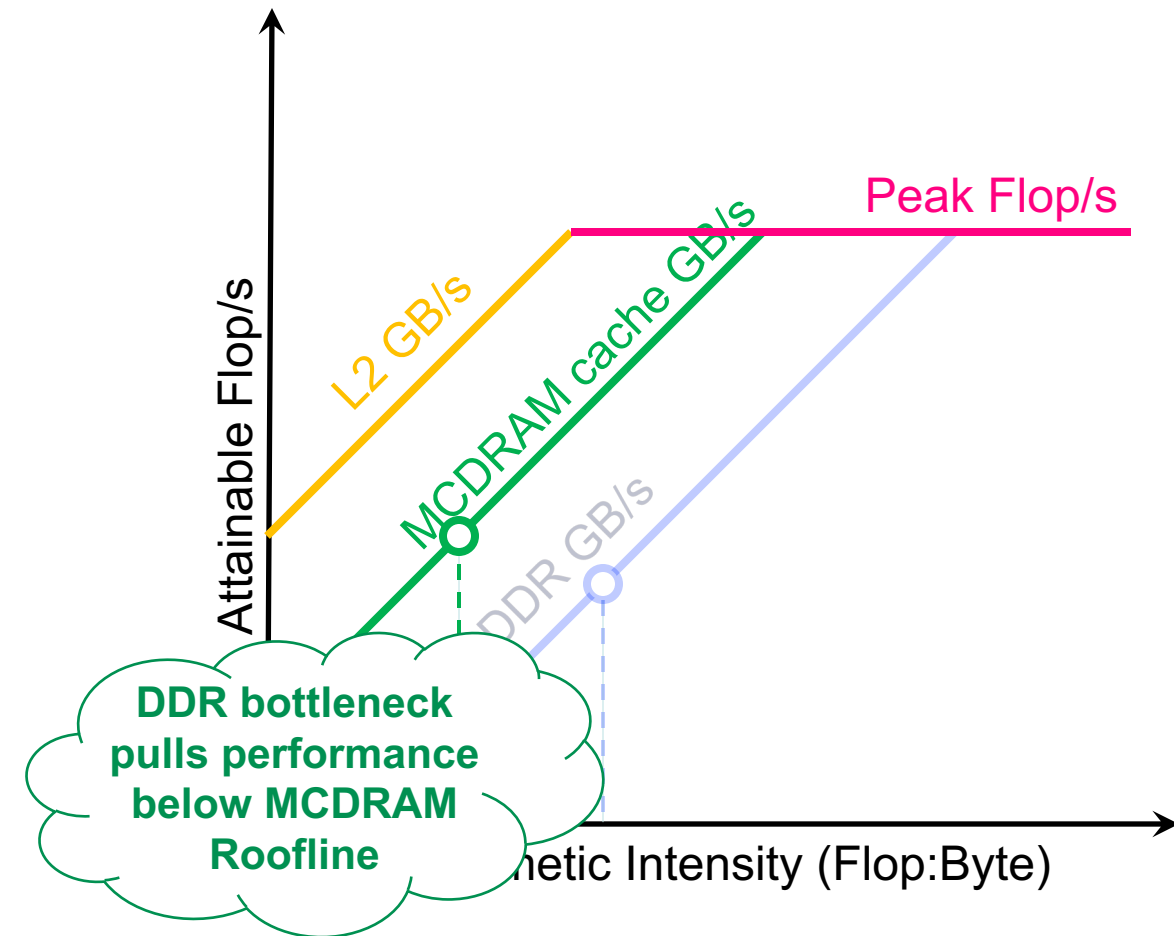
  - **… performance is bound by the minimum**

# Hierarchical Roofline

- Construct superposition of Rooflines…

  - Measure bandwidth

  - Measure AI for each level of memory

  - Although an loop nest may have multiple AI's and multiple bounds (flops, L1, L2, … DRAM)…
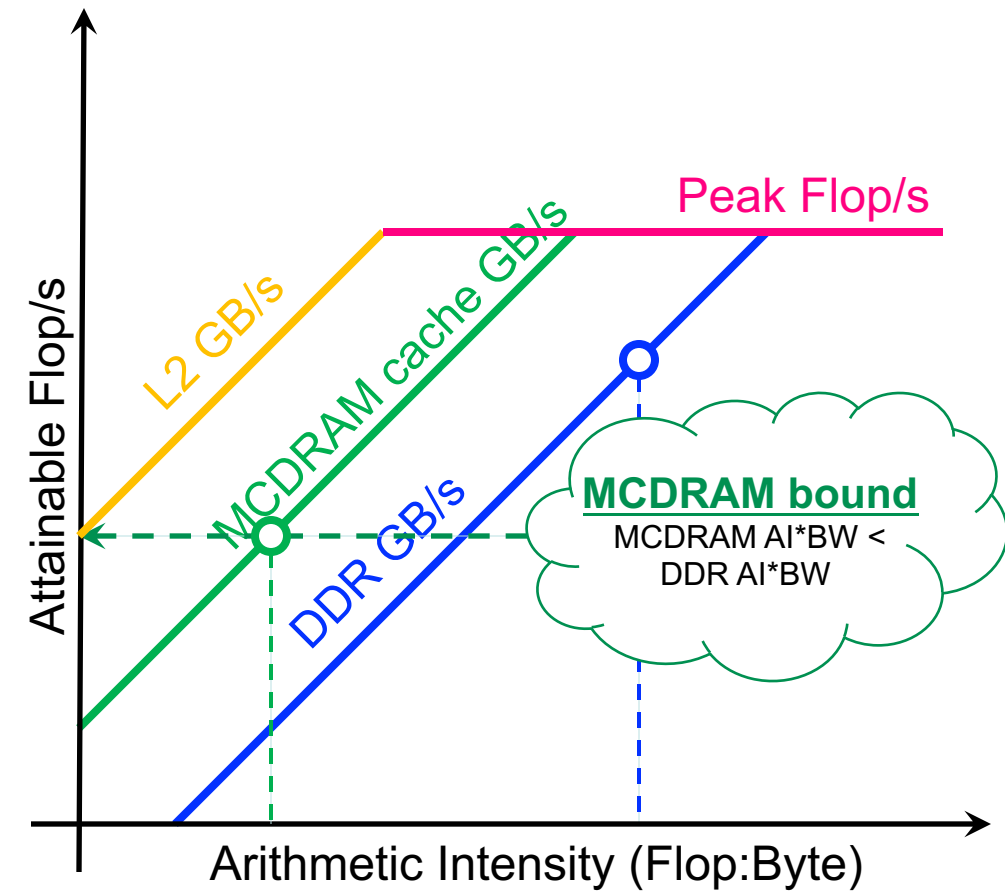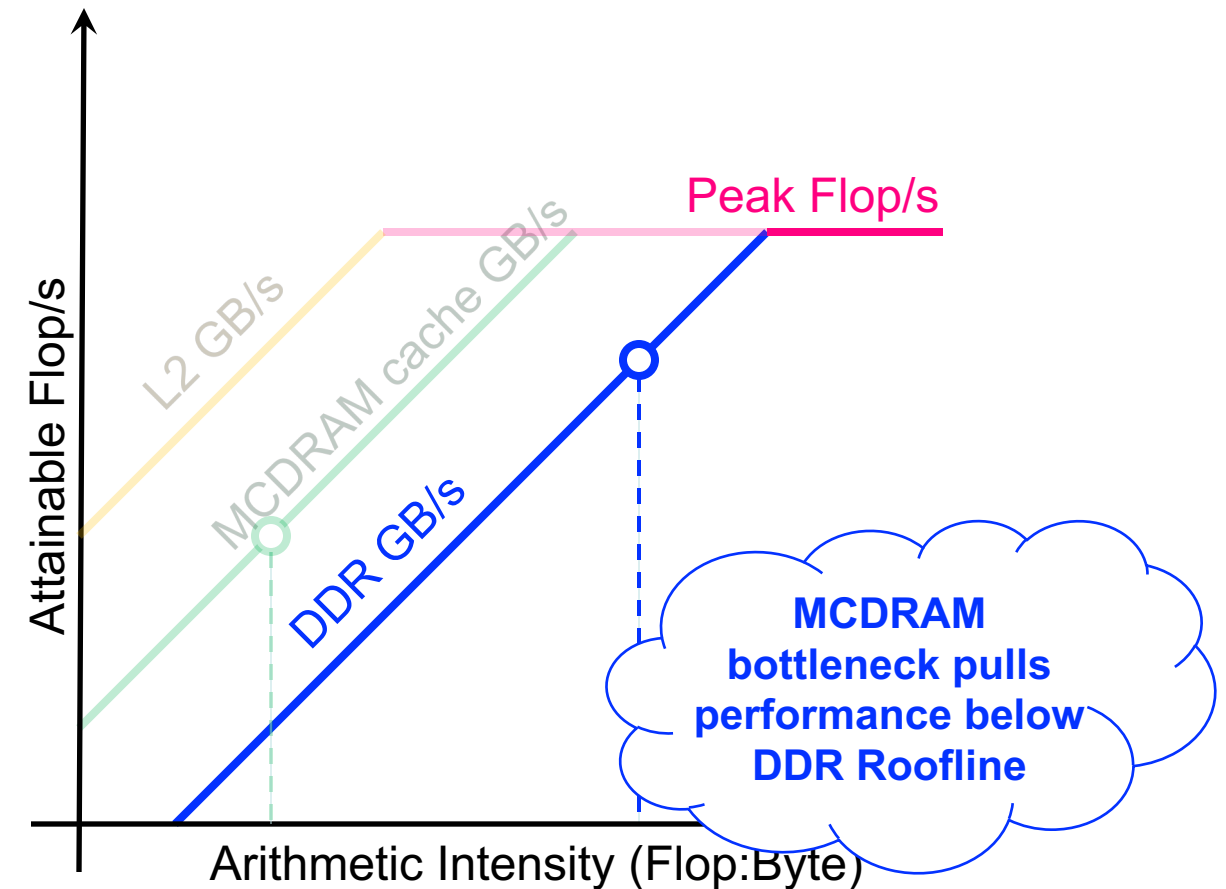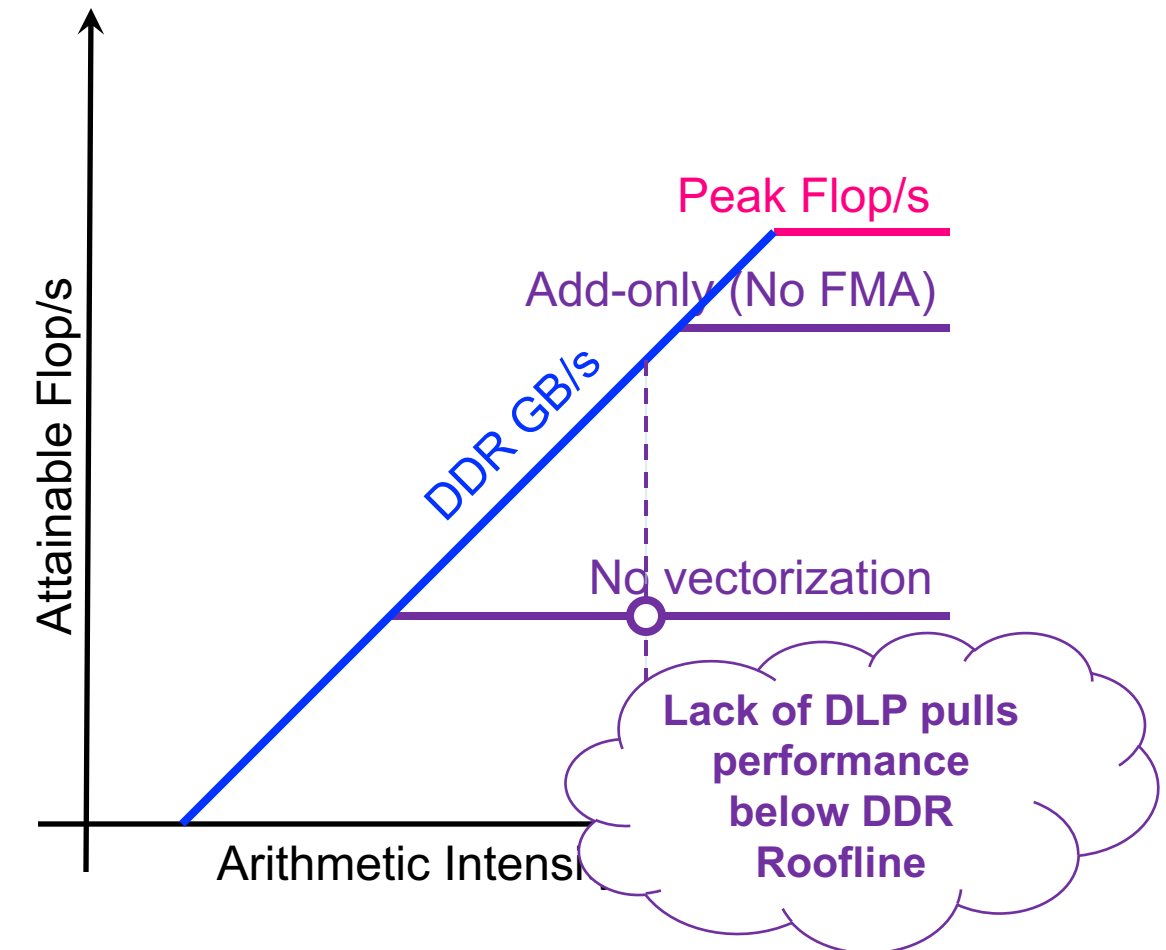
  - **… performance is bound by the minimum**

# Hierarchical Roofline

- Construct superposition of Rooflines…

  - Measure 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**



MCDRAM bottleneck pulls performance below DDR Roofline

# Data, Instruction, Thread-Level Parallelism…

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

- In reality, we must …
  - Use special instructions (e.g. FMA)
  - Vectorize loops (16 flops per instruction)
  - Hide FPU latency
    (unrolling, out-of-order execution)
  - Use all cores & sockets

- Without these, …
  - Peak performance is not attainable
  - Some kernels can transition from memory-bound to compute-bound



Peak Flop/s

Add-only (No FMA)

DDR GB/s

No vectorization

Attainable Flop/s

Arithmetic Intensity
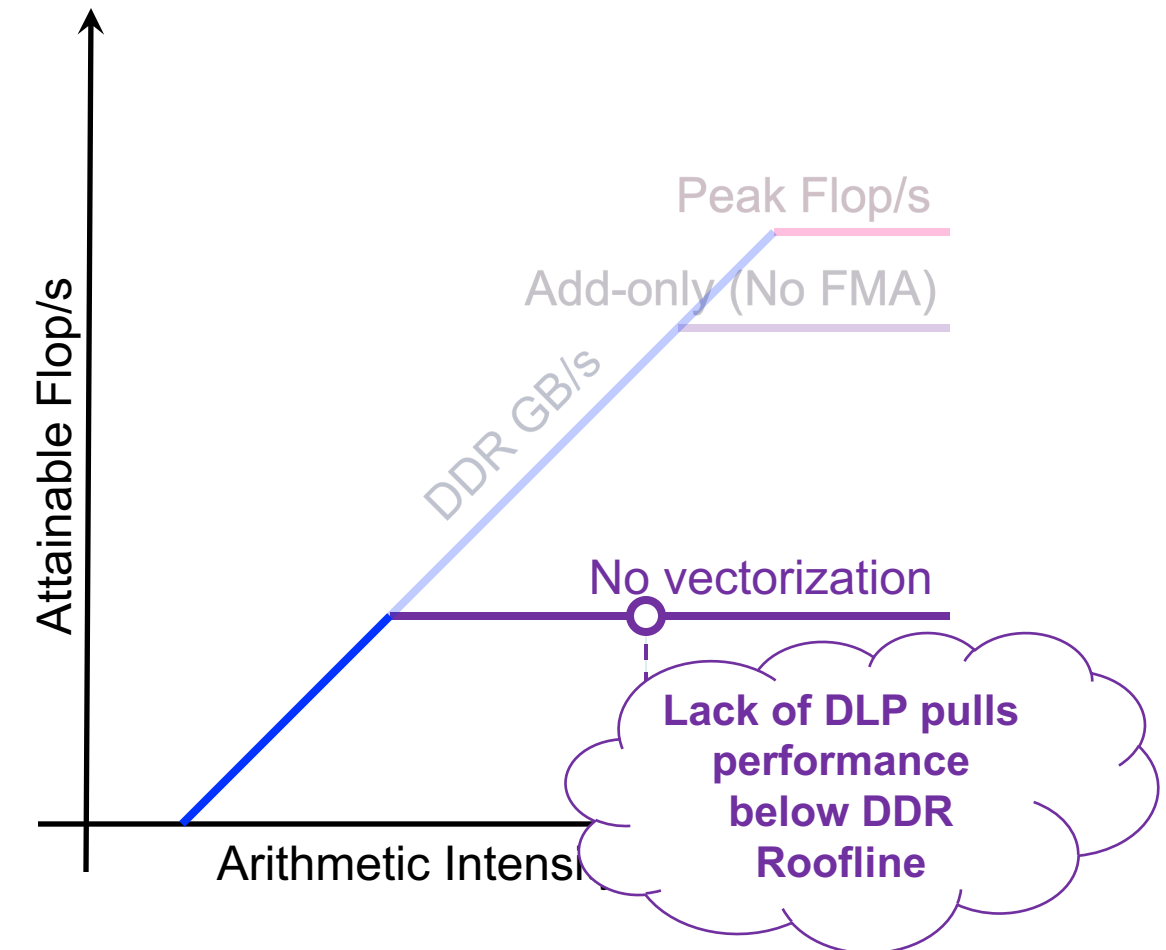
Lack of DLP pulls performance below DDR Roofline

# Data, Instruction, Thread-Level Parallelism…

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

- In reality, we must …
  - Use special instructions (e.g. FMA)
  - Vectorize loops (16 flops per instruction)
  - Hide FPU latency
    (unrolling, out-of-order execution)
  - Use all cores & sockets

- Without these, …
  - Peak performance is not attainable
  - Some kernels can transition from memory-bound to compute-bound



Peak Flop/s

Add-only (No FMA)

DDR GB/s

Attainable Flop/s

No vectorization

Arithmetic Intensity

Lack of DLP pulls performance below DDR Roofline

BERKELEY LAB

# Roofline-Driven Performance Optimization

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



(Graph with y-axis "Attainable Flop/s" and x-axis "Arithmetic Intensity (Flop:Byte)". Lines labeled: "Peak Flop/s", "No FMA", "DDR GB/s", "DDR GB/s (NUMA)", "No vectorization".)

# Roofline-Driven Performance Optimization

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

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

# Roofline-Driven Performance Optimization

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

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

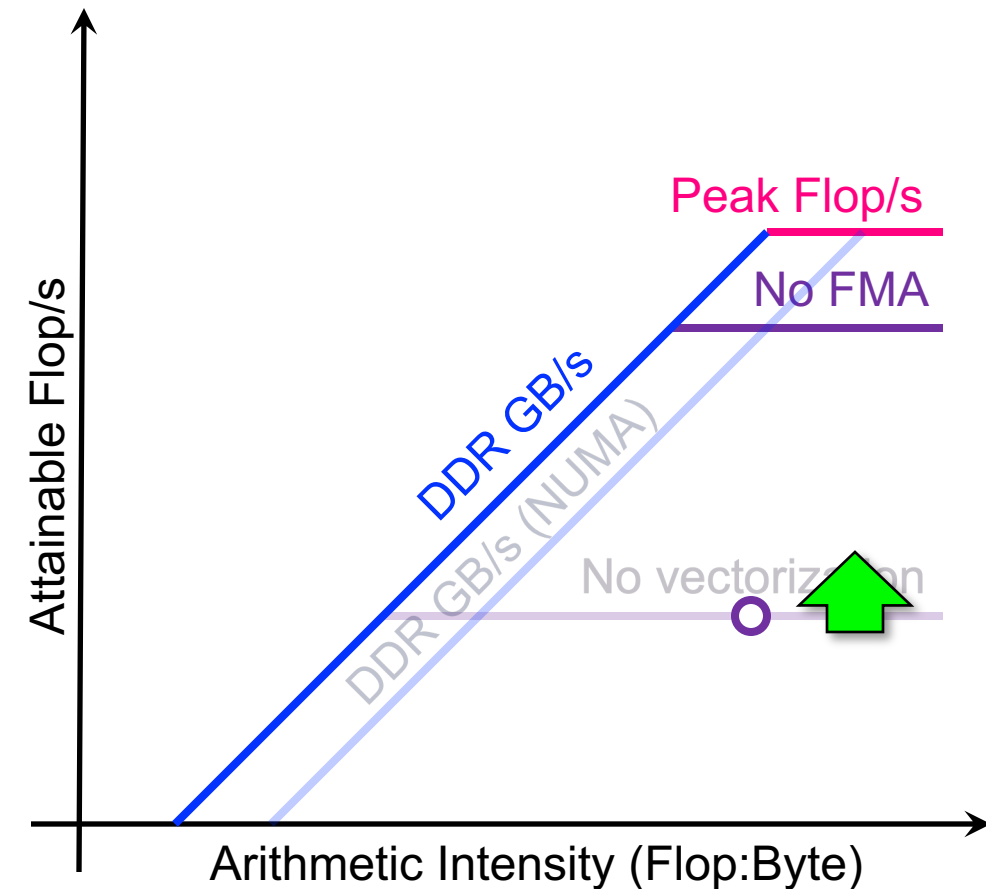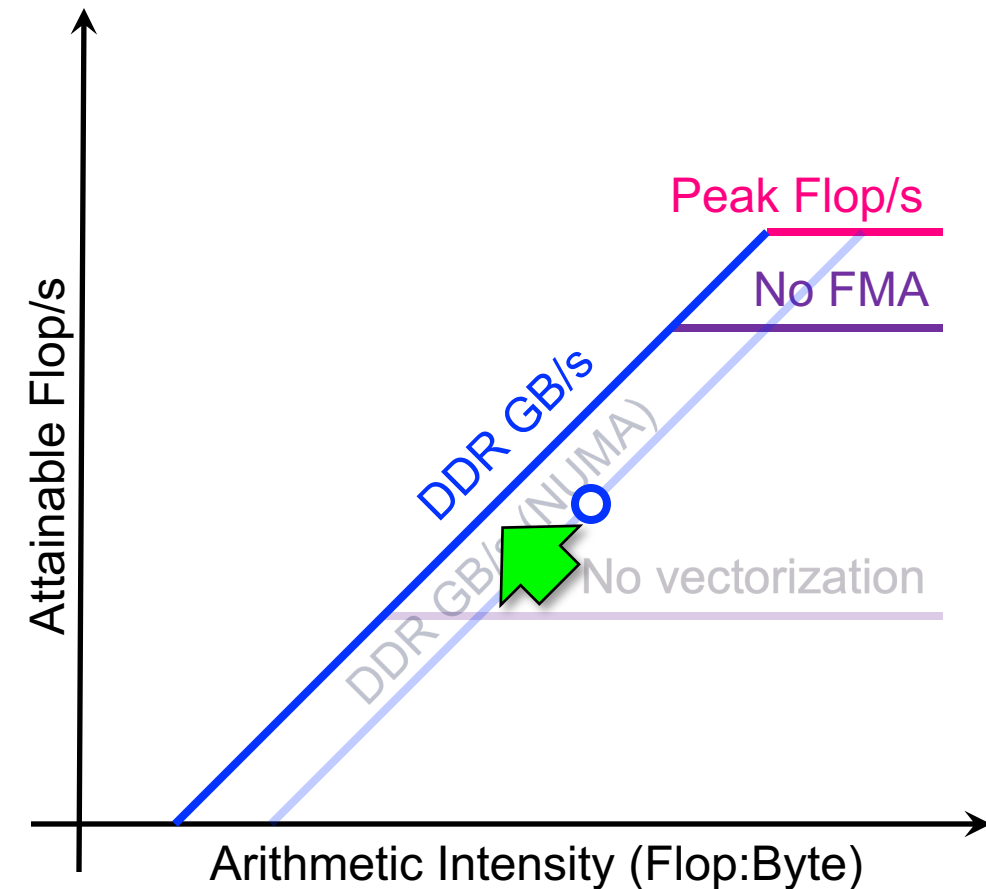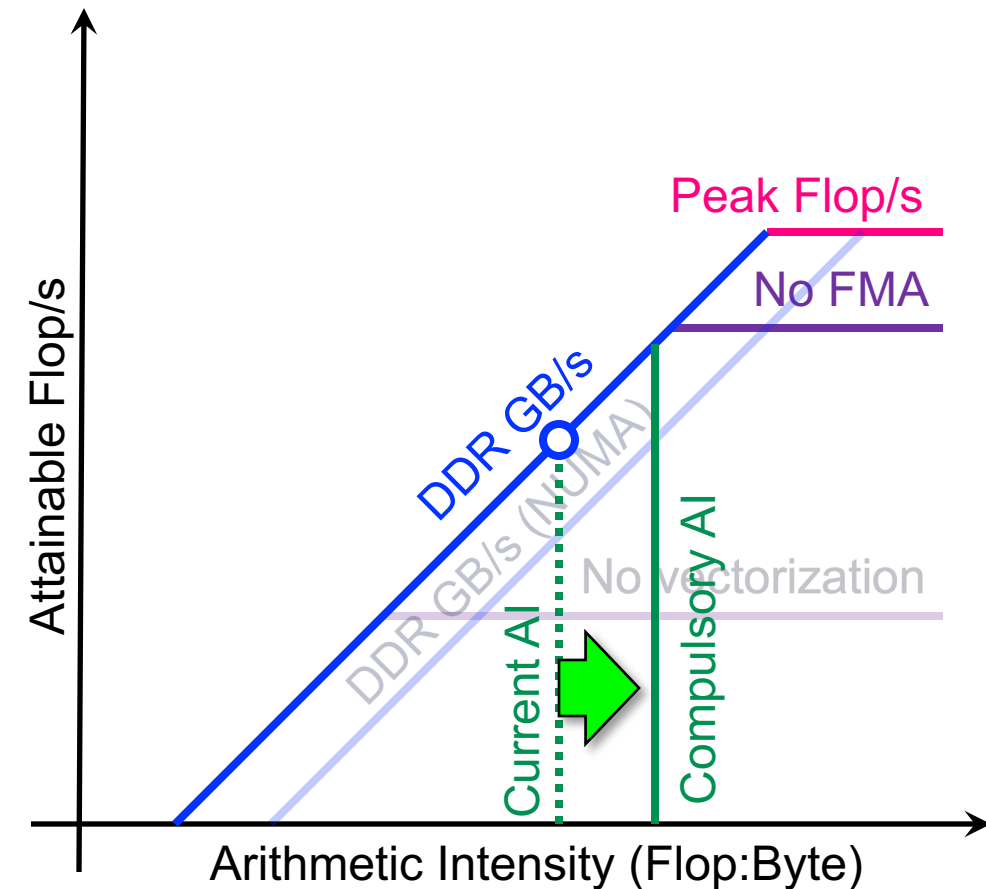- **Maximize memory bandwidth (e.g. NUMA-aware, unit-stride)**

# Roofline-Driven Performance Optimization

- 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, unit stride)

- **Minimize data movement (e.g. cache blocking)**

# Step 1:
## Machine Characterization

# Machine Characterization

- **"Theoretical Performance"** numbers can be highly optimistic…

  - Pin BW vs. sustained 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/

45

# Step 2:
# Application Characterization

# Measuring AI

- To characterize execution with Roofline we need…
  - **Time**
  - **Flops** (=> flop's / time)
  - **Data movement** between each level of memory (=> Flop's / GB's)

- We can look at the full application…
  - Coarse grained, 30-min average
  - Misses many details and bottlenecks

- or we can look at individual loop nests…
  - Requires auto-instrumentation on a loop by loop basis
  - Moreover, we should probably differentiate data movement or flops on a core-by-core basis.

BERKELEY LAB

# How Do We Count Flop's?

## Manual Counting

- Go thru each loop nest and count the number of FP operations
- ✓ Works best for deterministic loop bounds
- ✓ or parameterize by the number of iterations (recorded at run time)
- ✗ Not scalable

## Perf. Counters

- Read counter before/after
- ✓ More Accurate
- ✓ Low overhead (<%) == can run full MPI applications
- ✓ Can detect load imbalance
- ✗ Requires privileged access
- ✗ Requires manual instrumentation (+overhead) or full-app characterization
- ✗ Broken counters = garbage
- ✗ May not differentiate FMADD from FADD
- ✗ No insight into special pipelines

## Binary Instrumentation

- Automated inspection of assembly at run time
- ✓ Most Accurate
- ✓ FMA-, VL-, and mask-aware
- ✓ Can count instructions by class/type
- ✓ Can detect load imbalance
- ✓ Can include effects from non-FP instructions
- ✓ Automated application to multiple loop nests
- ✗ >10x overhead (short runs / reduced concurrency)

BERKELEY LAB

# How Do We Measure Data Movement?

## Manual Counting

- Go thru each loop nest and estimate how many bytes will be moved
- Use a mental model of caches
- ✓ Works best for simple loops that stream from DRAM (stencils, FFTs, spare, …)
- ✗ N/A for complex caches
- ✗ Not scalable

## Perf. Counters

- Read counter before/after
- ✓ Applies to full hierarchy (L2, DRAM,
- ✓ Much more Accurate
- ✓ Low overhead (<%) == can run full MPI applications
- ✓ Can detect load imbalance
- ✗ Requires privileged access
- ✗ Requires manual instrumentation (+overhead) or full-app characterization

## Cache Simulation

- Build a full cache simulator driven by memory addresses
- ✓ Applies to full hierarchy and multicore
- ✓ Can detect load imbalance
- ✓ Automated application to multiple loop nests
- ✗ Ignores prefetchers
- ✗ >10x overhead (short runs / reduced concurrency)

BERKELEY LAB

# Previously Cobbled 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/

# Intel Advisor

- **Includes Roofline Automation…**
  - ✓ Automatically instruments applications (one dot per loop nest/function)
  - ✓ Computes FLOPS and AI for each function (**CARM**)
  - ✓ 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



Memory-bound, invest into cache blocking etc.

Compute bound: invest into SIMD…

[1]Experimental Feature, the look and feel and exact behavior is subject for change

51

# 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

[1]Experimental Feature, the look and feel and exact behavior is subject for change

53

# 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*

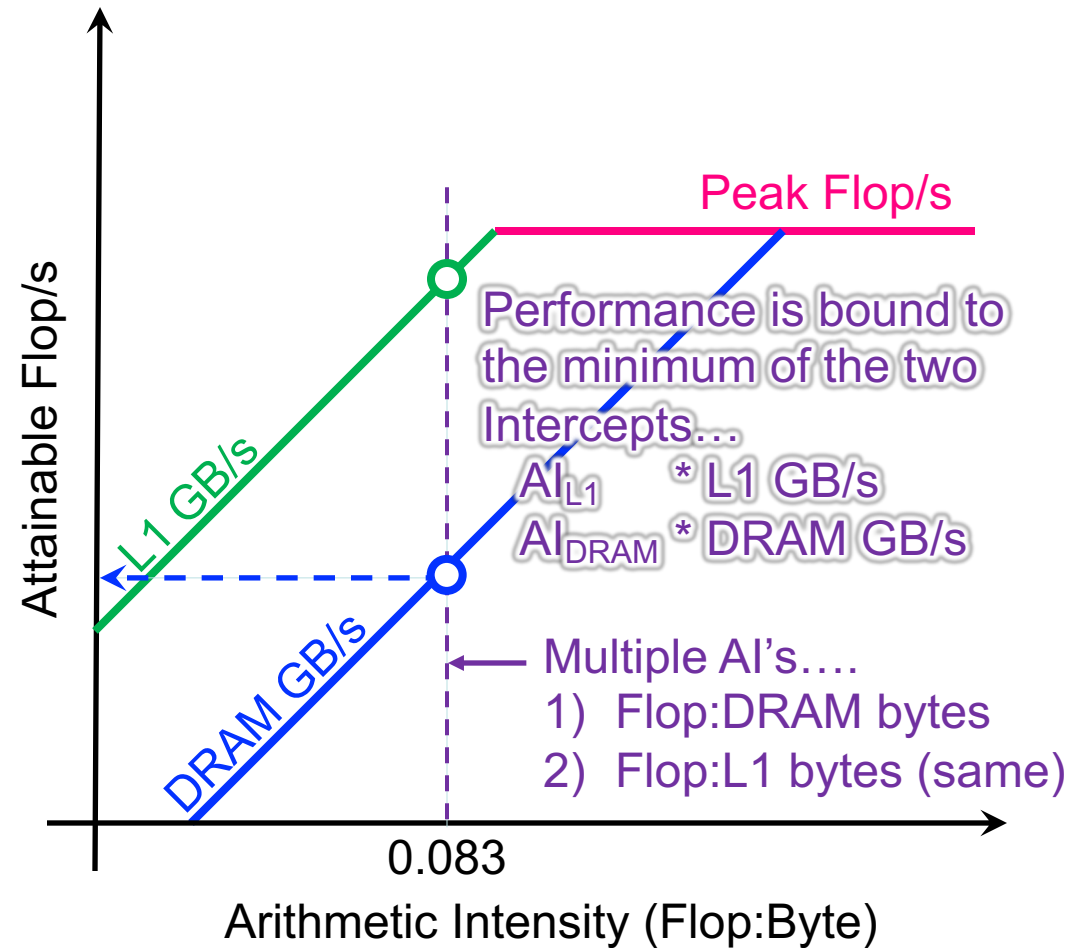- Requires static analysis or *binary instrumentation* to measure AI

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];
}
```
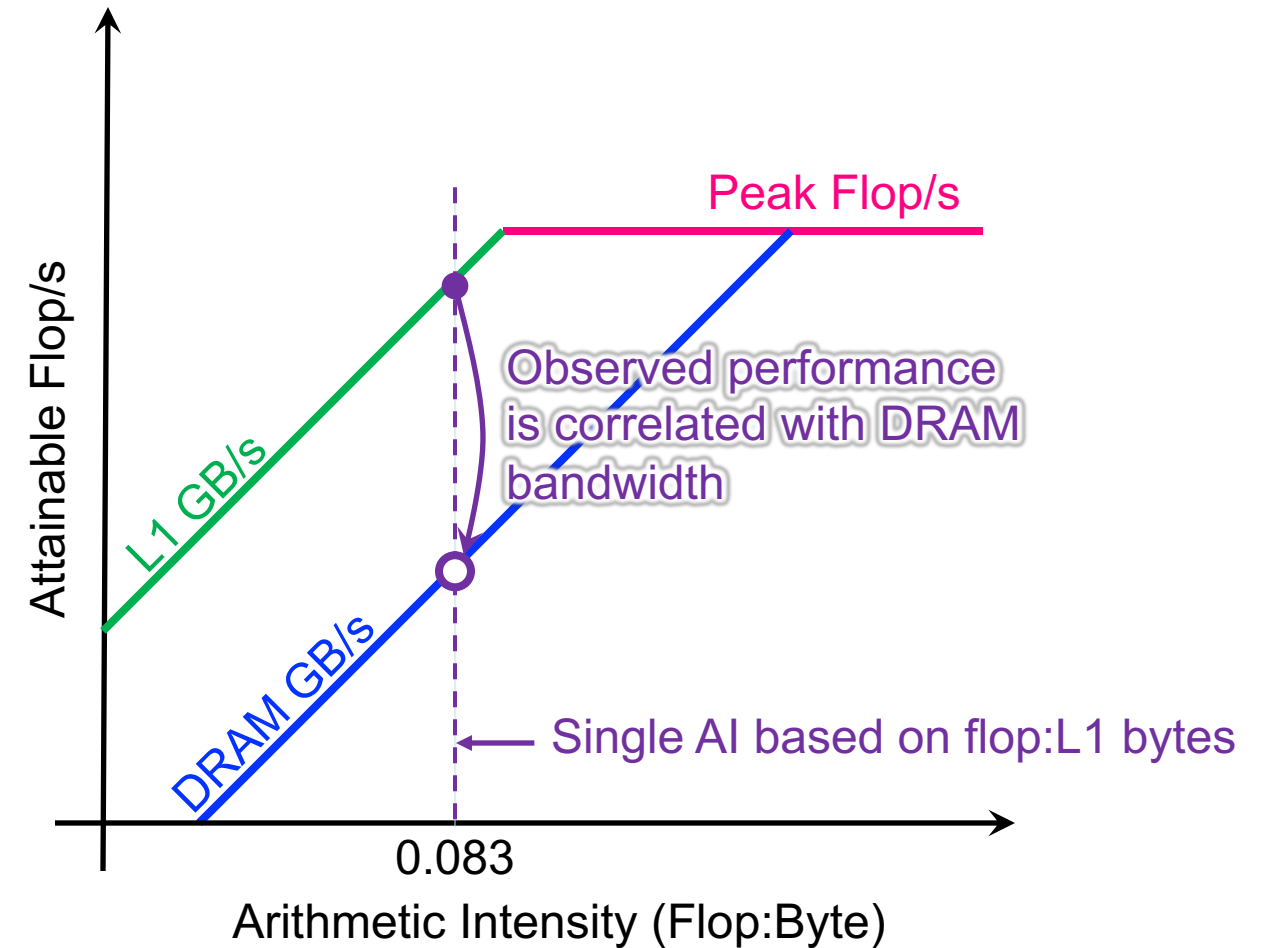
# Example: STREAM

## 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)

56

BERKELEY LAB
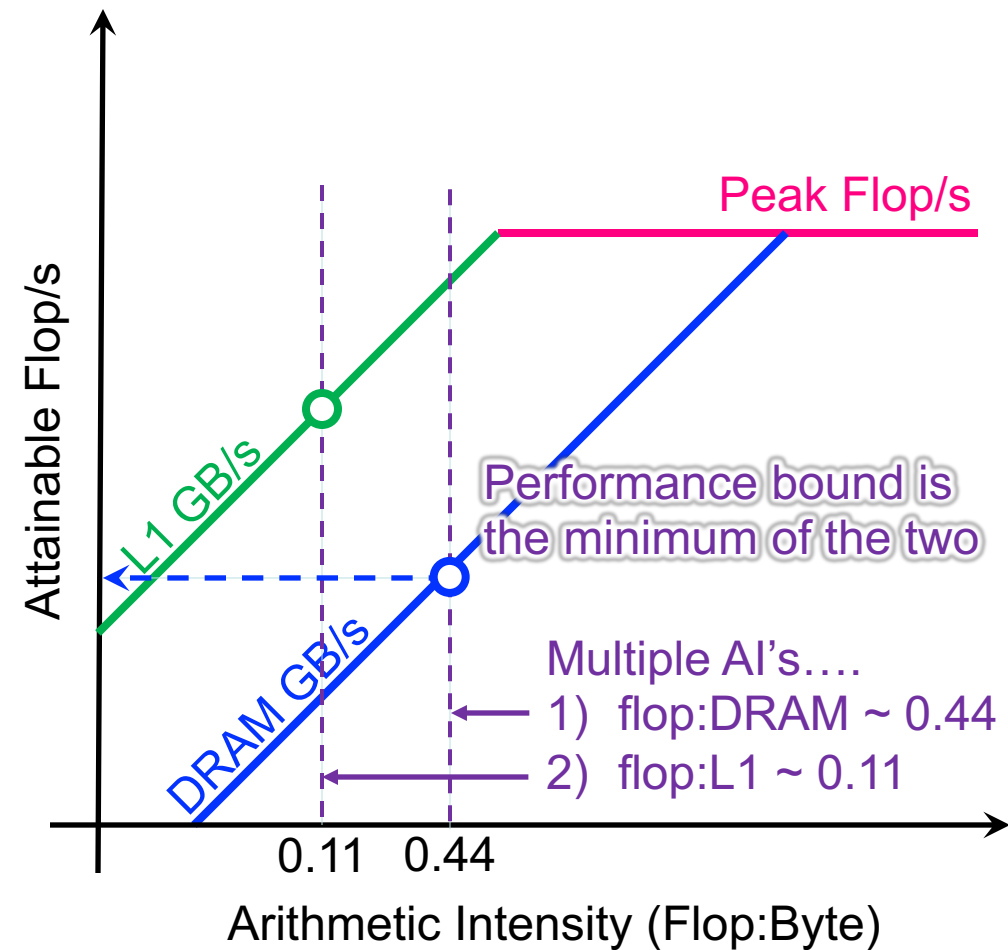
# Example: 7-point Stencil (Small Problem)

- ## 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[k][j][i+1]` is reused on next iteration of i.
  - `old[k][j+1][i]` is reused on next iteration of j.
  - `old[k+1][j][i]` is reused on next 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++){
   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  ];

}}}
```
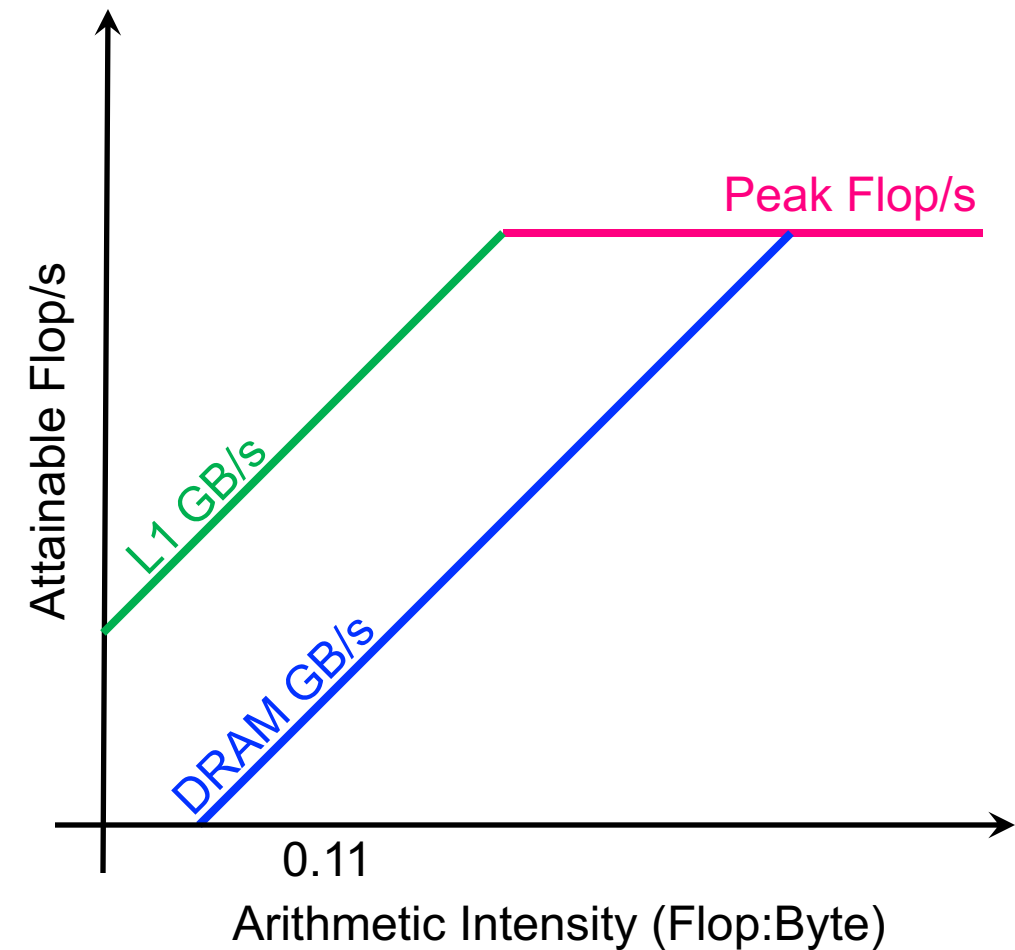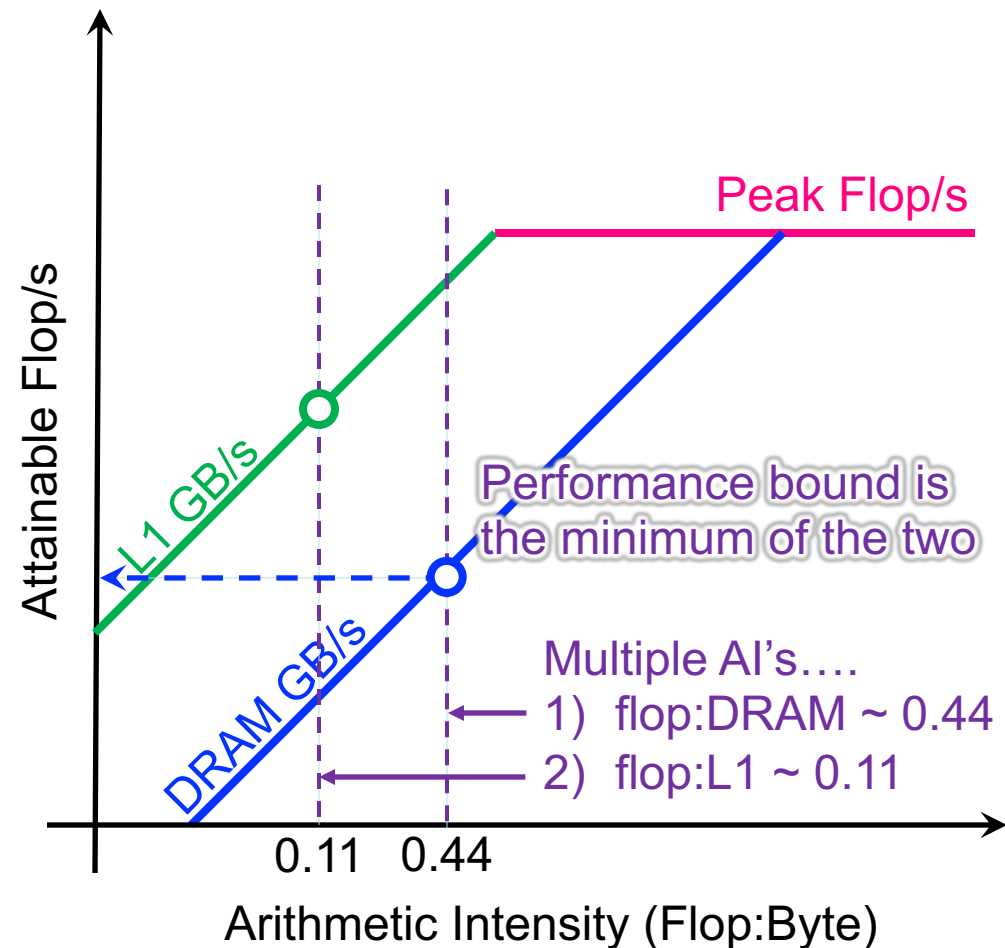
BERKELEY LAB

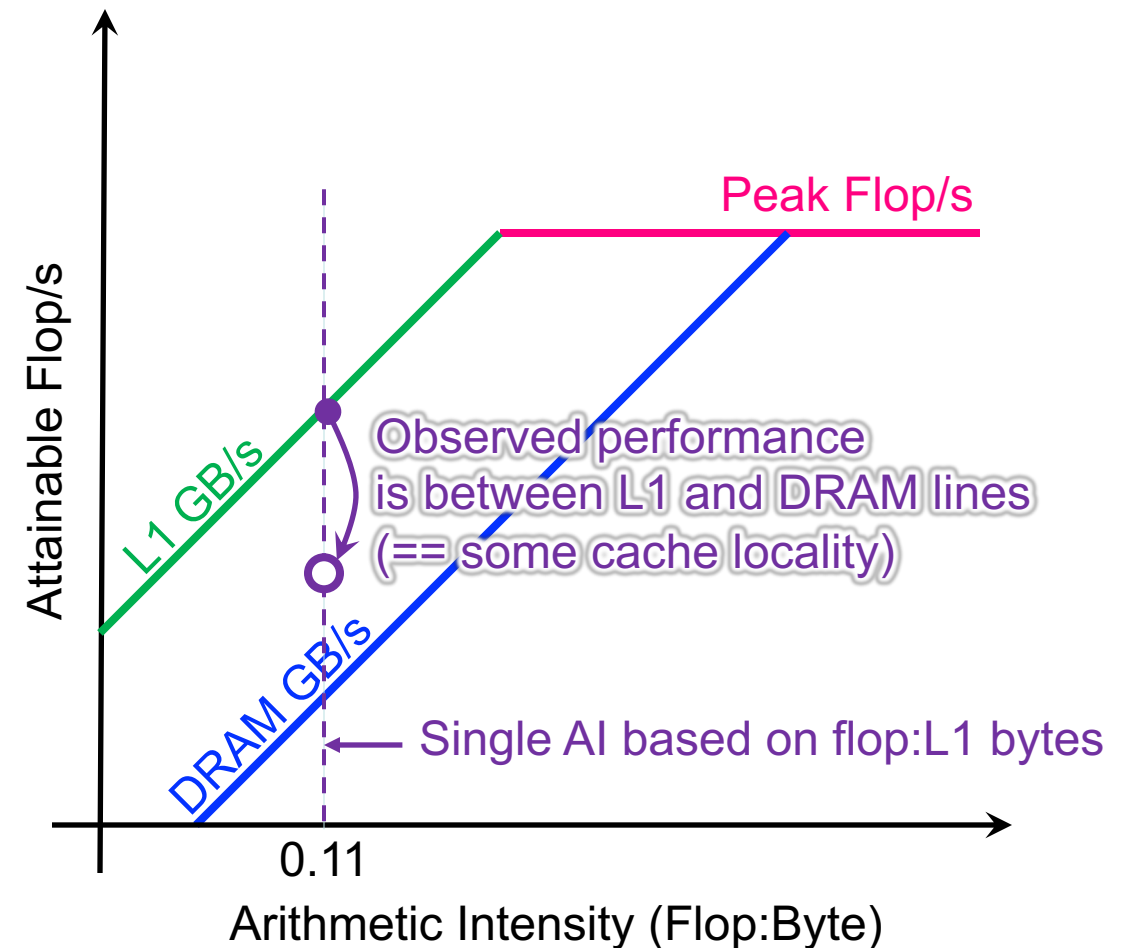# Example: 7-point Stencil (Small Problem)

## Hierarchical Roofline



Peak 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

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

0.11

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

# Example: 7-point Stencil (Small Problem)

## Hierarchical Roofline



Performance bound is the minimum of the two

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

Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

0.11  0.44

Arithmetic Intensity (Flop:Byte)

## Cache-Aware Roofline



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

Single AI based on flop:L1 bytes

Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

0.11

Arithmetic Intensity (Flop:Byte)

BERKELEY LAB

# 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)

BERKELEY LAB

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

## Hierarchical Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/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

DRAM GB/s

Attainable Flop/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



Attainable Flop/s (y-axis)
Arithmetic Intensity (Flop:Byte) (x-axis)

Peak Flop/s
L1 GB/s
DRAM GB/s

0.11   0.20

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

## Cache-Aware Roofline



Attainable Flop/s (y-axis)
Arithmetic Intensity (Flop:Byte) (x-axis)

Peak Flop/s
L1 GB/s
DRAM GB/s

0.11

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

Single AI based on flop:L1 bytes

# Don't forget to take the Survey...

http://bit.ly/sc18-eval

# Backup

# Refining Roofline:
## NUMA

# NUMA Effects

- **Cori's Haswell nodes are built from 2 Xeon processors (sockets)**

  - Memory attached to each socket (fast)

  - Interconnect that allows remote memory access (slow == NUMA)

  - Improper memory allocation can result in more than a 2x performance penalty
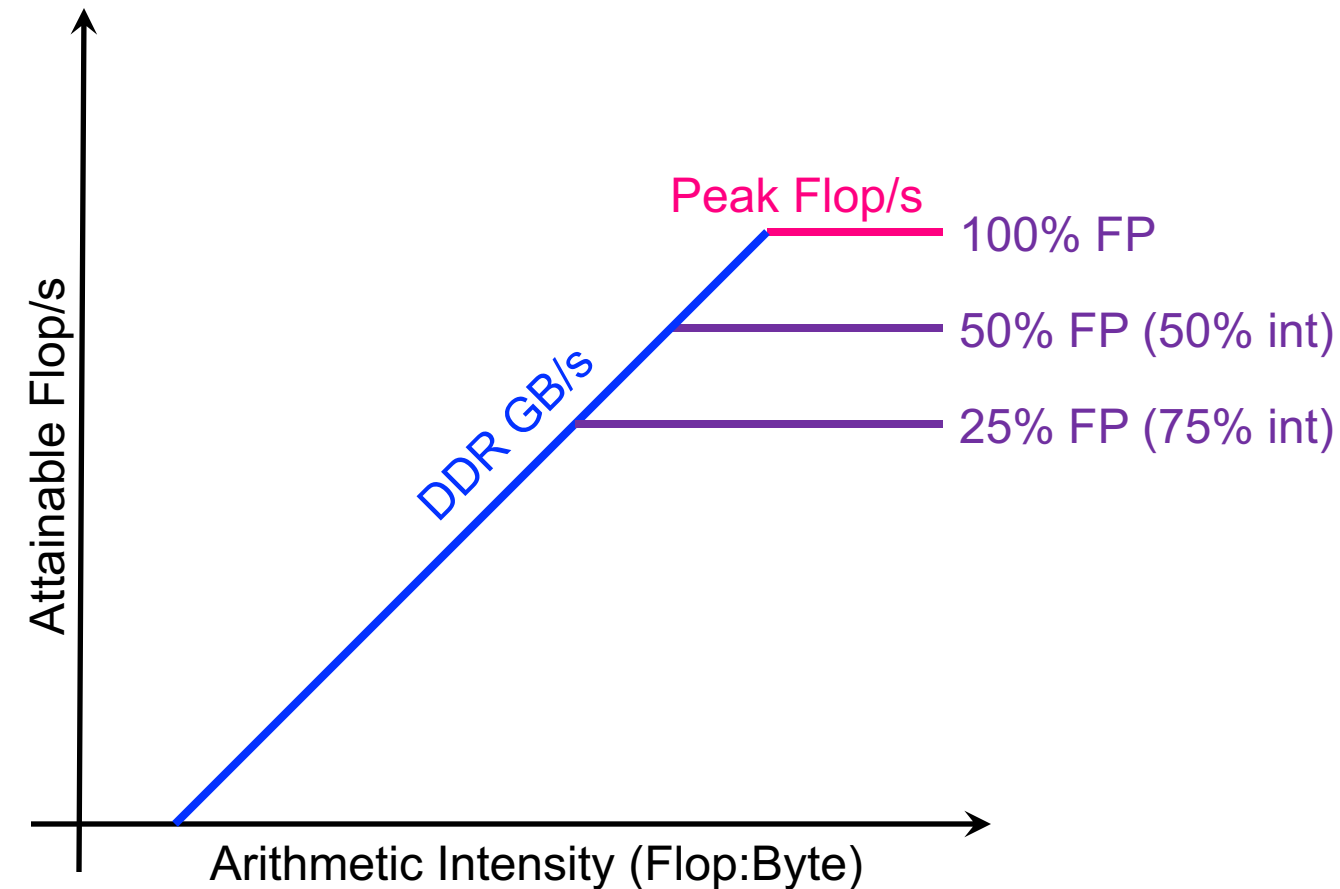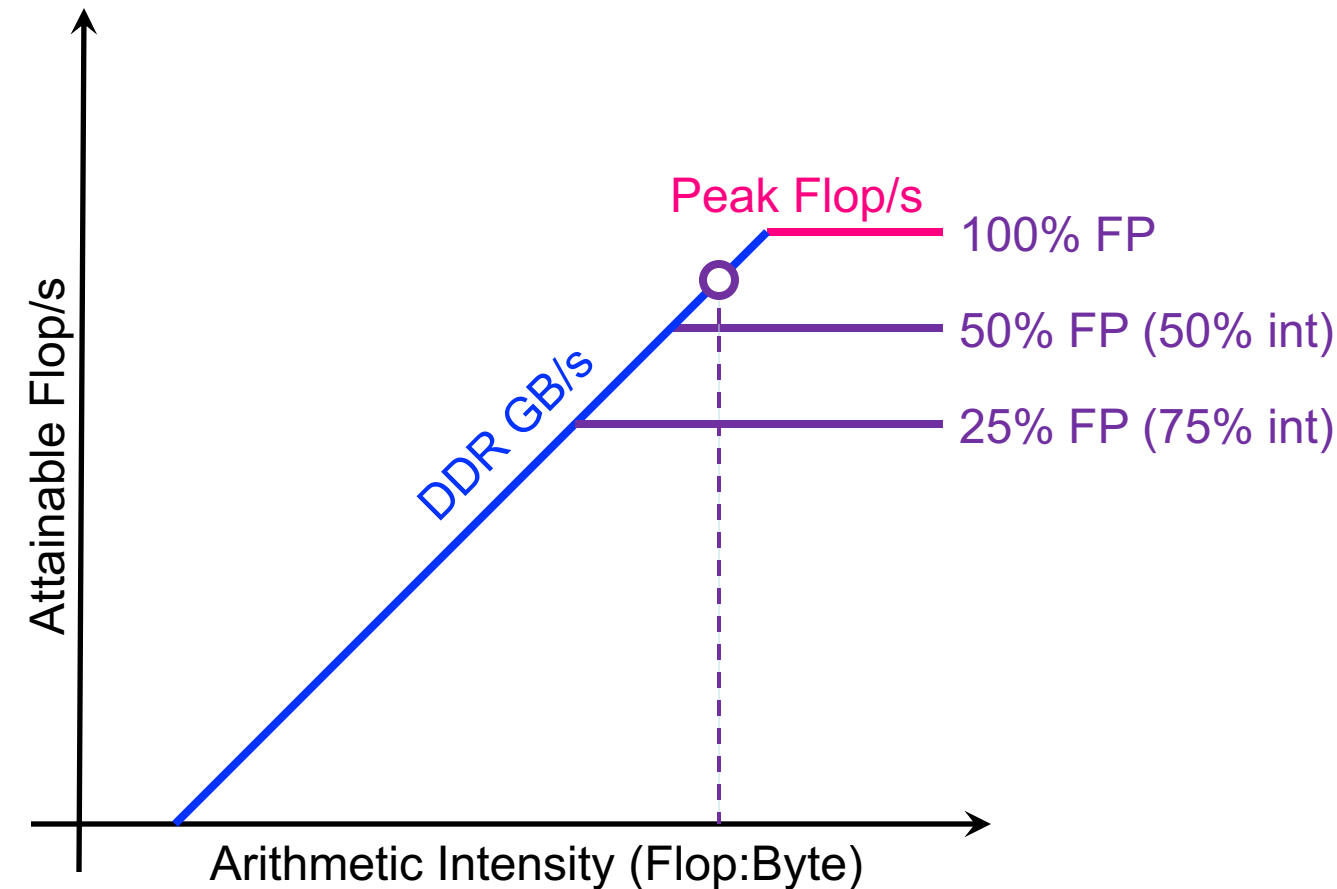
# Superscalar vs. Instruction mix

- 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
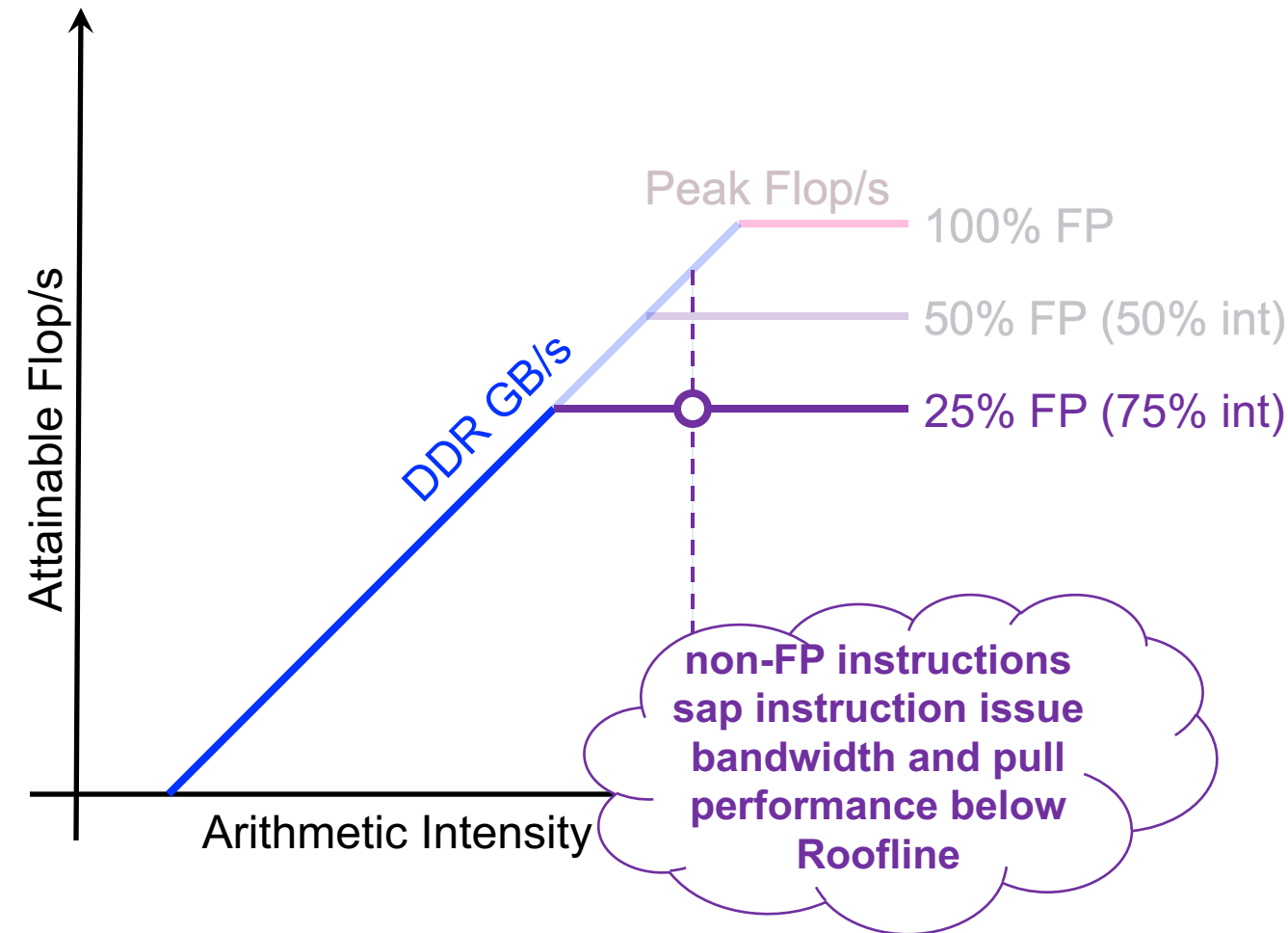  - Requires 100% of the instructions to be FP to get peak performance



Peak Flop/s

100% FP

50% FP (50% int)

25% FP (75% int)

DDR GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

# 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
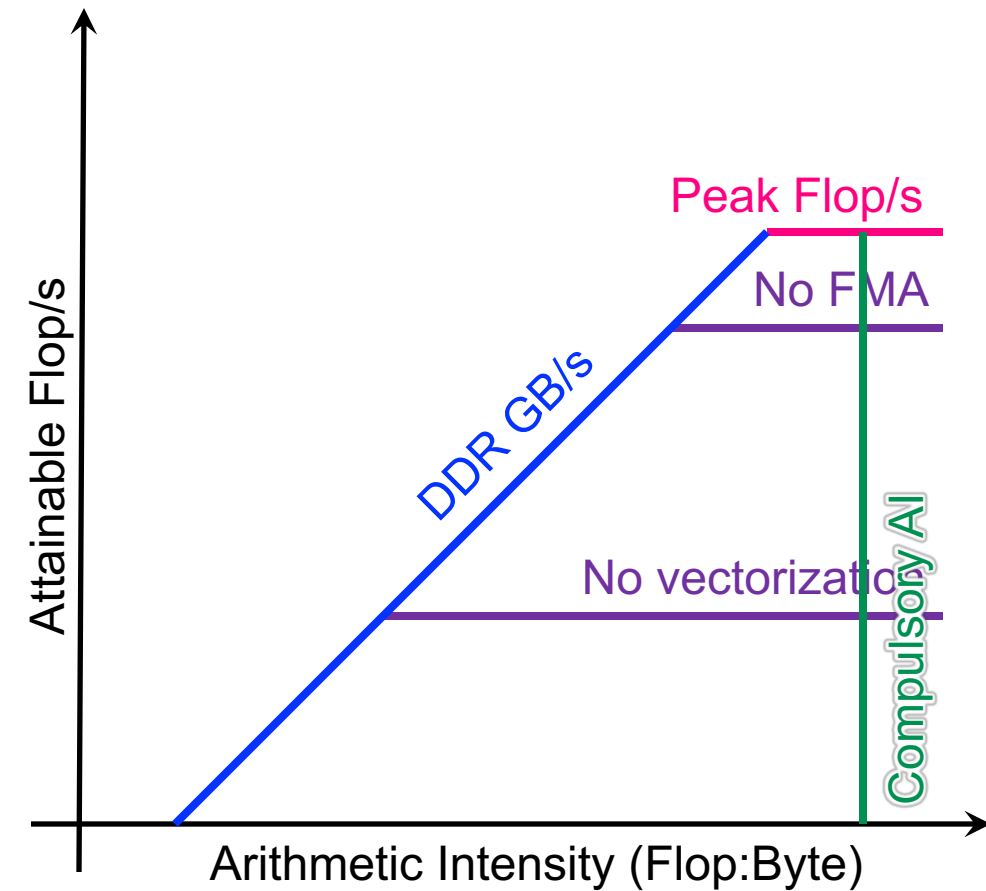  - Requires 100% of the instructions to be FP to get peak performance

# Refining Roofline:
## Compulsory, Capacity, and Conflict misses

# Locality Walls

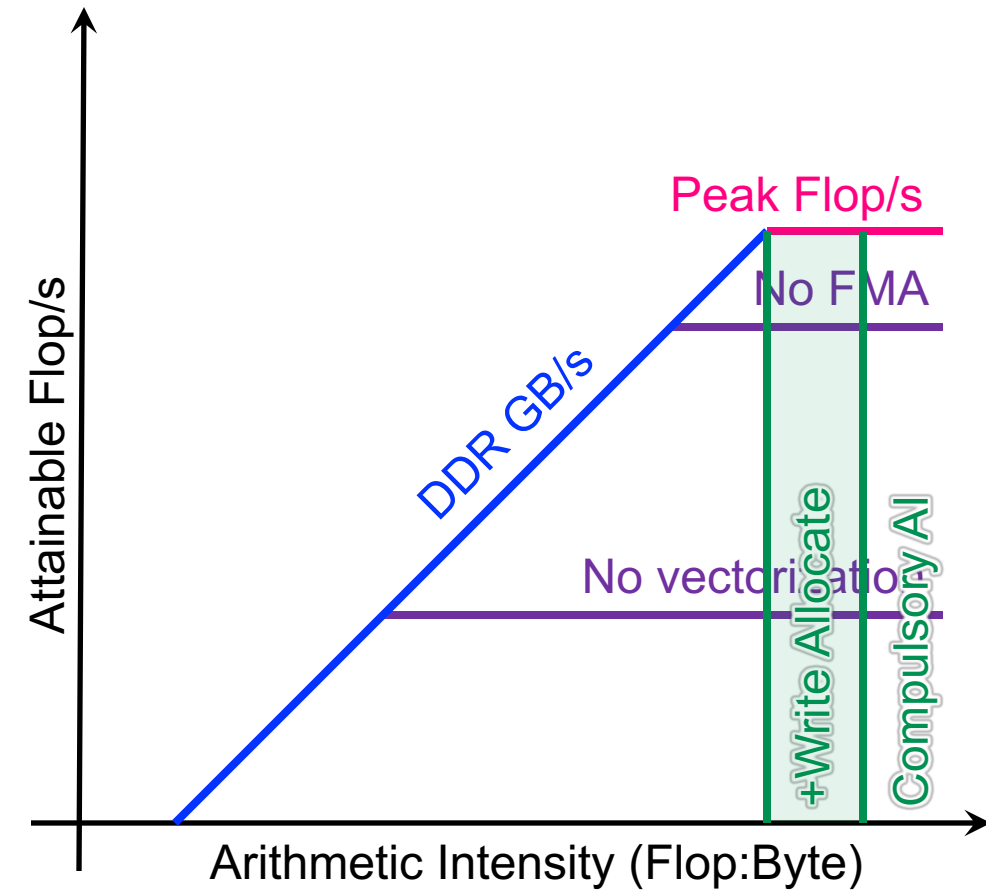- Naively, we can bound AI using only compulsory cache misses



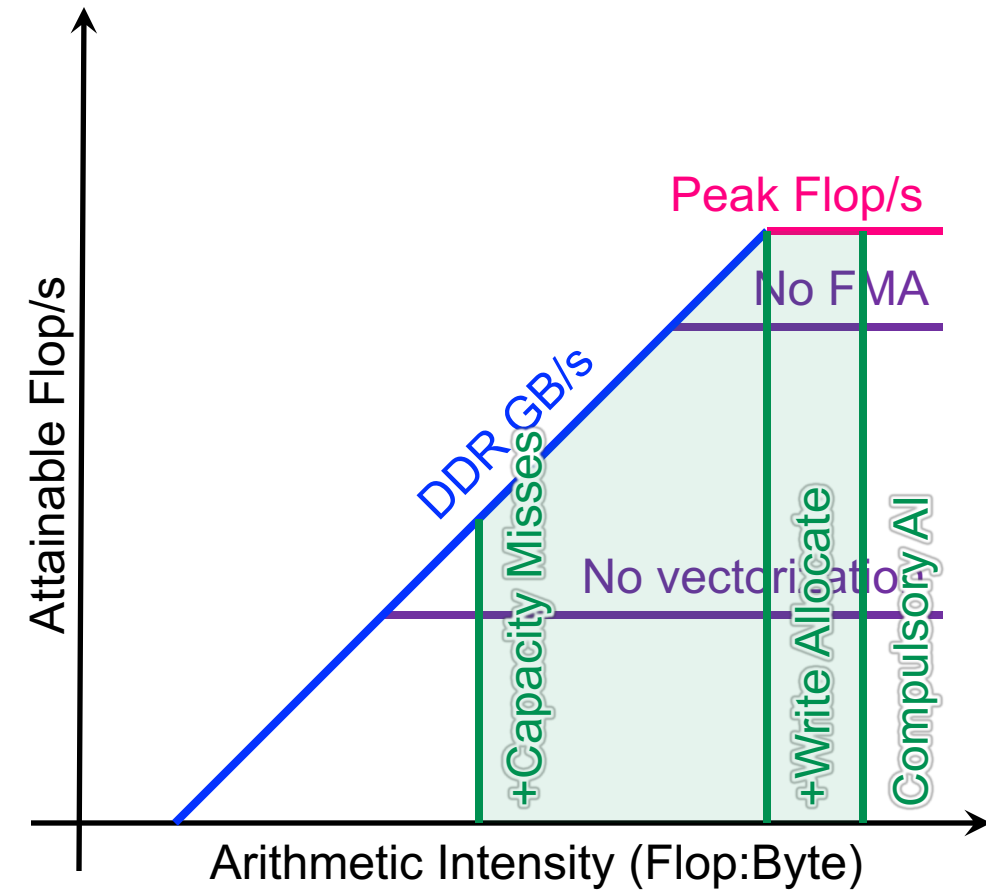$$AI = \frac{\#Flop's}{Compulsory\ Misses}$$

# Locality Walls

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

- However, write allocate caches can lower AI



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

# 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
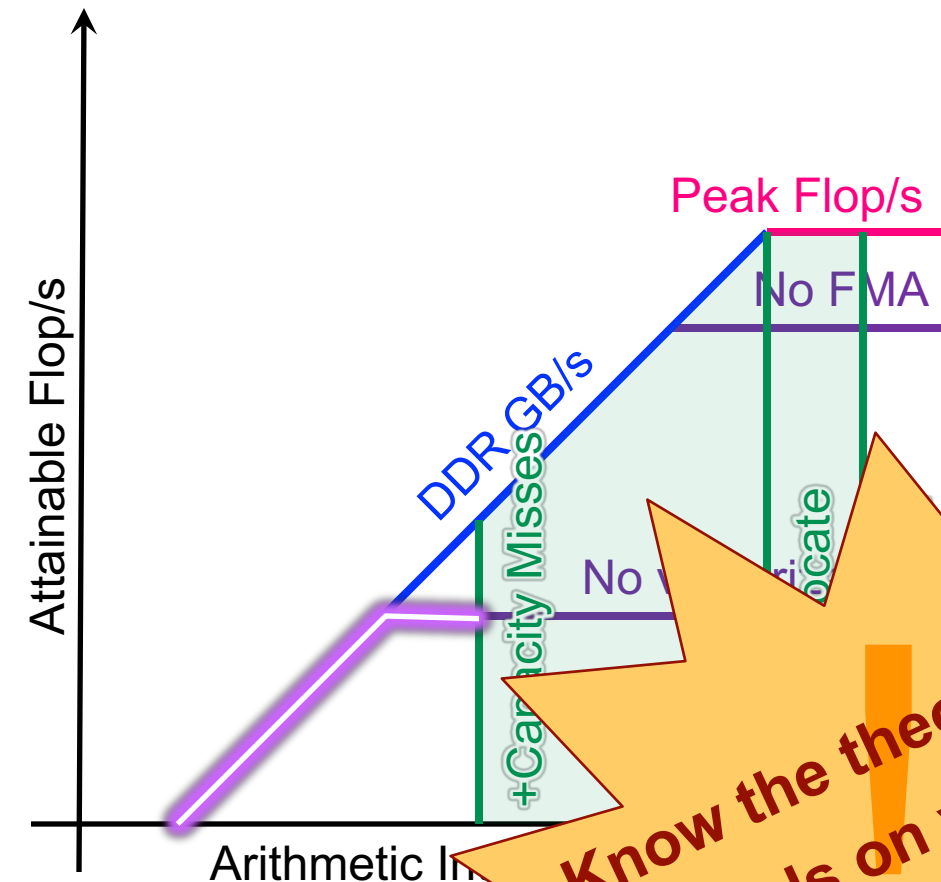
- Cache capacity misses can have a huge penalty

➢ **Compute bound became memory bound**

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

Peak Flop/s

No FMA

DDR GB/s

+Capacity Misses

No write allocate

Attainable Flop/s

Arithmetic In...

**Know the theoretical bounds on your AI.**

BERKELEY LAB

# LIKWID

- LIKWID provides easy to use wrappers for measuring performance counters...

  ✓ **Works on NERSC production systems**

  ✓ Distills counters into user-friendly metrics (e.g. MCDRAM Bandwidth)

  ✓ Minimal overhead (<1%)

  ✓ Scalable in distributed memory (MPI-friendly)

  ✓ Fast, high-level characterization

  ✗ No timing breakdowns

  ✗ Suffers from Garbage-in/Garbage Out

  (i.e. hardware counter must be sufficient and correct)

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

http://www.nersc.gov/users/software/performance-and-debugging-tools/likwid

BERKELEY LAB