# Performance Tuning with the Roofline Model on GPUs and CPUs

| | | |
|---|---|---|
| 2:30pm | Introduction to Roofline | Samuel Williams |
| 3:10pm | Roofline on NVIDIA GPUs | Samuel Williams |
| 3:35pm | NERSC Roofline Use Cases | Jonathan Madsen |
| 4:00pm | break | - |
| 4:30pm | Roofline on Intel GPUs | JaeHyuk Kwack |
| 4:55pm | ALCF Roofline Use Cases | Christopher Knight |
| 5:20pm | Advanced Roofline Topics | Khaled Ibrahim |
| 5:40pm | Profiling with TiMemory | Jonathan Madsen |
| 6:00pm | closing remarks / Q&A | all |

**BERKELEY LAB**
LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

# Presenters

**Samuel Williams**

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

**JaeHyuk Kwack**

ALCF
Argonne National Lab
jkwack@anl.gov

**Jonathan Madsen**

NERSC
Lawrence Berkeley National Lab
JRMadsen@lbl.gov

**Khaled Ibrahim**

Computational Research Division
Lawrence Berkeley National Lab
KZIbrahim@lbl.gov

**Christopher Knight**

ALCF
Argonne National Lab
knightc@anl.gov

# Acknowledgements

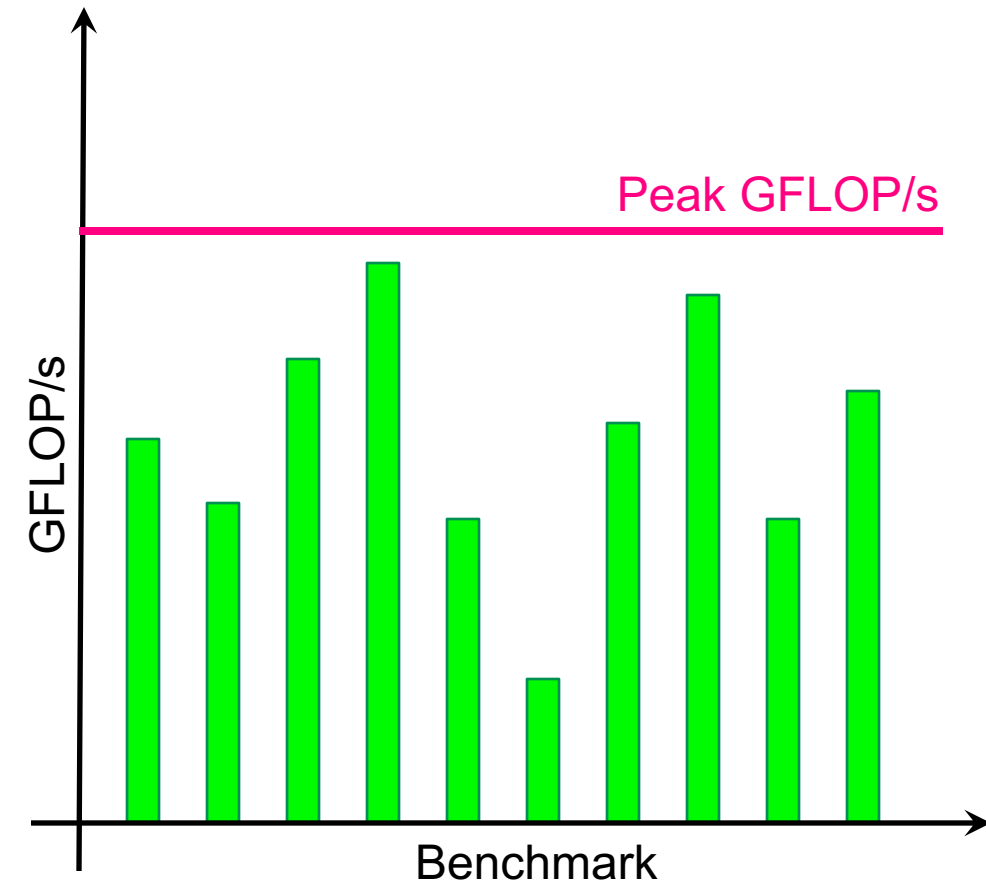**BERKELEY LAB**
LAWRENCE BERKELEY NATIONAL LABORATORY

**U.S. DEPARTMENT OF ENERGY**

We spend millions of dollars porting applications to CPUs and GPUs…

How do we know if we are getting our money's worth?

BERKELEY LAB
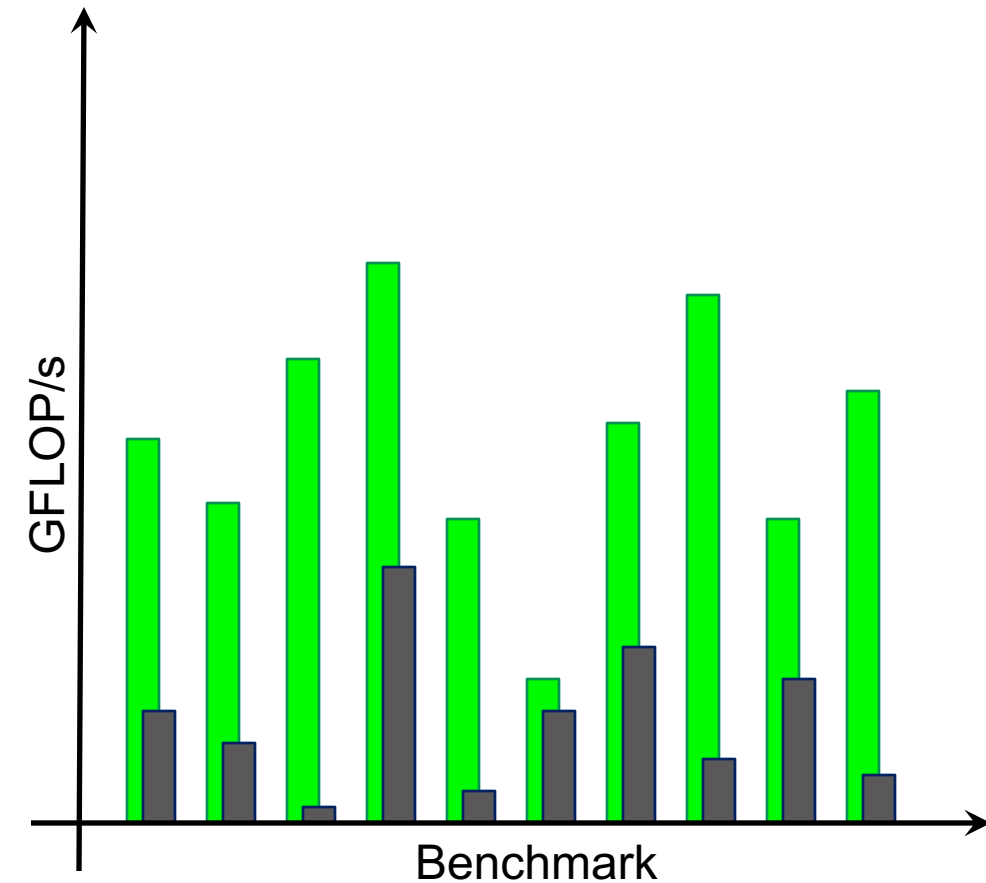LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

# Getting our money's worth?

- Do we get good performance on application benchmarks?

- Imagine profiling a mix of GPU-accelerated benchmarks …
- GFLOP/s alone may not be particularly insightful

# Are we getting good performance?

- ## We could compare performance to a CPU…

  - Speedup may seem random
  - Aren't GPUs always 10x faster than a CPU?
  - If not, what does that tell us about architecture, algorithm or implementation?

  - ➢ **'Speedup' provides no insights into architecture, algorithm, or implementation.**

BERKELEY LAB

# Are we getting good performance?

- We could take a CS approach and look at performance counters…
  - Record microarchitectural events on CPUs/GPUs
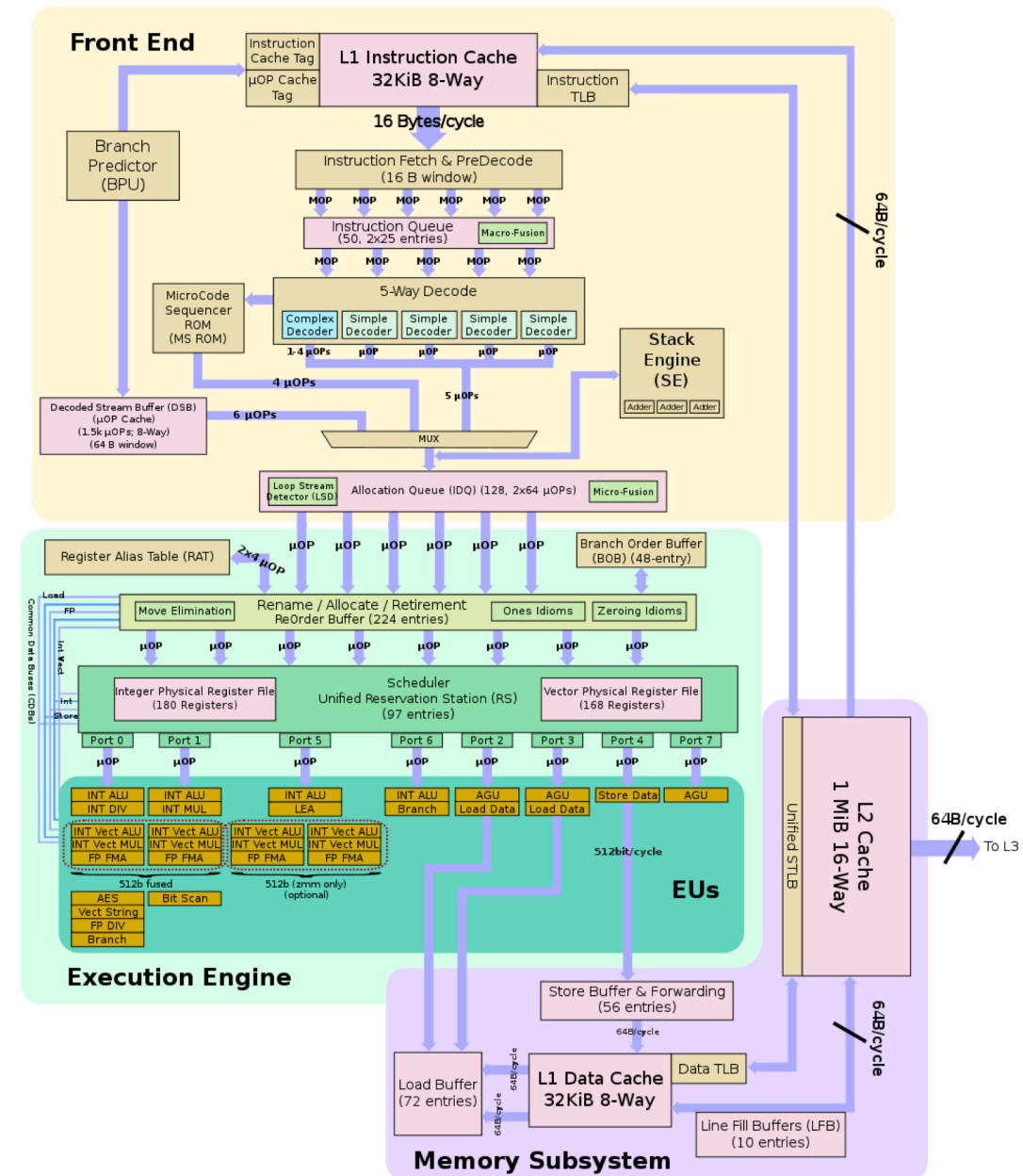  - Use arcane, architecture-specific terminology
  - May be broken

  - We may be able to show correlation between events, but…

  - ➤ **…providing actionable guidance to CS, AM, applications, or procurement can prove elusive.**

```
:
.
FRONTEND_RETIRED.LATENCY_GE_8_PS
FRONTEND_RETIRED.LATENCY_GE_16_PS
FRONTEND_RETIRED.LATENCY_GE_32_PS
RS_EVENTS.EMPTY_END
FRONTEND_RETIRED.L2_MISS_PS
FRONTEND_RETIRED.L1I_MISS_PS
FRONTEND_RETIRED.STLB_MISS_PS
FRONTEND_RETIRED.ITLB_MISS_PS
ITLB_MISSES.WALK_COMPLETED
BR_MISP_RETIRED.ALL_BRANCHES_PS
IDQ.MS_SWITCHES
FRONTEND_RETIRED.LATENCY_GE_2_BUBBLES_GE_1_PS
BR_MISP_RETIRED.ALL_BRANCHES_PS
MACHINE_CLEARS.COUNT
MEM_LOAD_RETIRED.L1_HIT_PS
MEM_LOAD_RETIRED.FB_HIT_PS
MEM_LOAD_UOPS_RETIRED.L1_HIT_PS
MEM_LOAD_UOPS_RETIRED.HIT_LFB_PS
MEM_INST_RETIRED.STLB_MISS_LOADS_PS
MEM_UOPS_RETIRED.STLB_MISS_LOADS_PS
MEM_LOAD_RETIRED.L2_HIT_PSMEM_LOAD_UOPS_RETIRED.L2_HIT_PS
MEM_LOAD_RETIRED.L3_HIT_PS
MEM_LOAD_UOPS_RETIRED.LLC_HIT_PS
MEM_LOAD_UOPS_RETIRED.L3_HIT_PS
MEM_LOAD_RETIRED.L3_MISS_PS
MEM_LOAD_UOPS_RETIRED.LLC_MISS_PS
MEM_LOAD_UOPS_MISC_RETIRED.LLC_MISS_PS
MEM_LOAD_UOPS_RETIRED.L3_MISS_PS
MEM_INST_RETIRED.ALL_STORES_PS
MEM_UOPS_RETIRED.ALL_STORES_PS
ARITH.DIVIDER_ACTIVE
ARITH.DIVIDER_UOPS
ARITH.FPU_DIV_ACTIVE
INST_RETIRED.PREC_DIST
IDQ.MS_UOPS
INST_RETIRED.PREC_DIST
.
.
.
```

# Are we getting good performance?

- We could take the computer architect's approach and build a simulator to understand performance nuances…

  - Modern architectures are incredibly complex
  - Simulators may perfectly reproduce performance
  - Deluge of information interpretable only by computer architects
  - worse, might incur $10^6$x slowdowns

  ➢ **Provide no insights into quality or limits of algorithm or implementation.**

  ➢ **Provide no guidance to CS, AM, applications, or procurement.**

# What's missing…

- Each community speaks their own language and develops specialized tools/methodologies

- Need common mental model of application execution on target system

- Sacrifice accuracy to gain…
  - Architecture independence / extensibility
  - Readily understandable by broad community
  - Intuition, insights, and guidance to CS, AM, apps, procurement, and vendors

➢ **Roofline is just such a model**



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

BERKELEY LAB

# Data Movement or Compute?

- **Which takes longer?**
  - Data Movement
  - Compute



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

BERKELEY LAB

# Data Movement or Compute?

- **Which takes longer?**
  - Data Movement
  - Compute

- **Is performance limited by compute or data movement?**



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

BERKELEY LAB

# Data Movement or Compute?

- **Which takes longer?**
  - Data Movement
  - Compute

- **Is performance limited by compute or data movement?**

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

# Data Movement or Compute?

- **Which takes longer?**
  - Data Movement
  - Compute

- **Is performance limited by compute or data movement?**



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

*Arithmetic Intensity (AI) = measure of data locality*

# Arithmetic Intensity

- Measure of data locality (data reuse)

- Ratio of **Total Flops** performed to **Total Bytes** moved

- For the DRAM Roofline…

  o   Total Bytes to/from DRAM

  o   Includes all cache and prefetcher effects

  o   Can be very different from total loads/stores (bytes requested)

  o   Equal to ratio of sustained GFLOP/s to sustained GB/s (time cancels)

# (DRAM) Roofline Model

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

AI (Arithmetic Intensity) = FLOPs / Bytes (moved to/from DRAM )

- Plot bound on **Log-log scale** as a function of AI (data locality)



*Transition @ AI ==*
*Peak GFLOP/s / Peak GB/s ==*
*'Machine Balance'*

BERKELEY LAB

# (DRAM) Roofline Model

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

AI (Arithmetic Intensity) = FLOPs / Bytes (moved to/from DRAM )

- Plot bound on **Log-log scale** as a function of AI (data locality)
- Roofline tessellates the locality-performance plane into five regions…



*unattainable performance (greater than peak GFLOP/s)*

Peak GFLOP/s

*Compute bound*

Attainable FLOP/s

*unattainable performance (insufficient bandwidth)*

HBM GB/s

**Bandwidth-Bound**

*poor performance*

Arithmetic Intensity (FLOP:Byte)

*Transition @ AI ==
Peak GFLOP/s / Peak GB/s ==
'Machine Balance'*

BERKELEY LAB

# (DRAM) Roofline Model

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

AI (Arithmetic Intensity) = FLOPs / Bytes (moved to/from DRAM )

- Plot bound on **Log-log scale** as a function of AI (data locality)
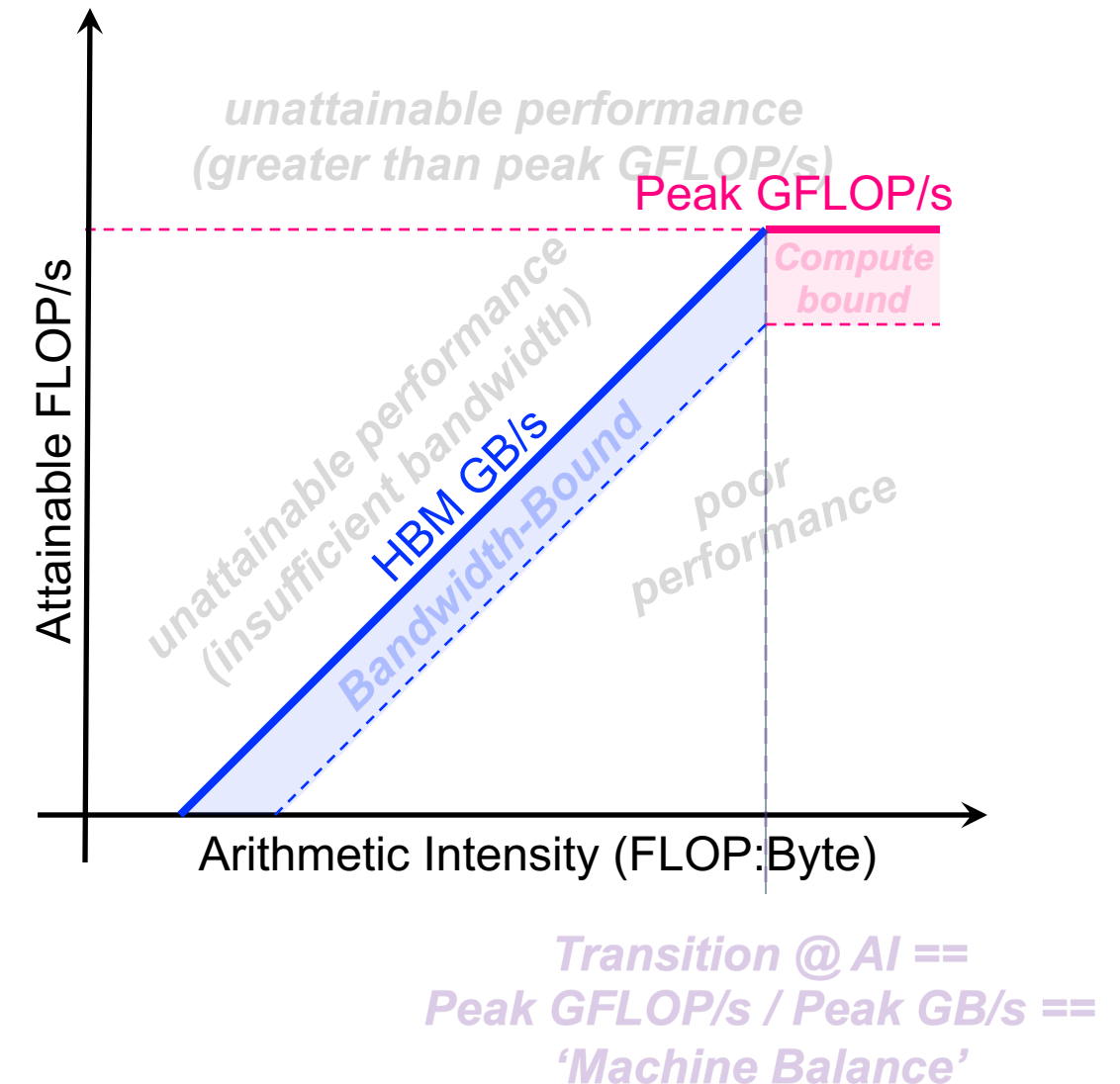
- Roofline tessellates the locality-performance plane into five regions…

- Measure application (AI,GF/s) and plot in the 2D locality-performance plane.



*unattainable performance (greater than peak GFLOP/s)*

Peak GFLOP/s

*Compute bound*

Attainable FLOP/s

*unattainable performance (insufficient bandwidth)*

HBM GB/s

**Bandwidth-Bound**

*poor performance*

Arithmetic Intensity (FLOP:Byte)

*Transition @ AI == Peak GFLOP/s / Peak GB/s == 'Machine Balance'*

BERKELEY LAB

# Roofline Examples

# 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 GFLOP/s

Attainable FLOP/s

HBM GB/s

GFLOP/s ≤ AI * HBM GB/s

TRIAD

0.083     5.0

Arithmetic Intensity (FLOP:Byte)

BERKELEY LAB

# Roofline Example #2

- Conversely, 7-point constant coefficient stencil…



```
#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  ];
}}}
```
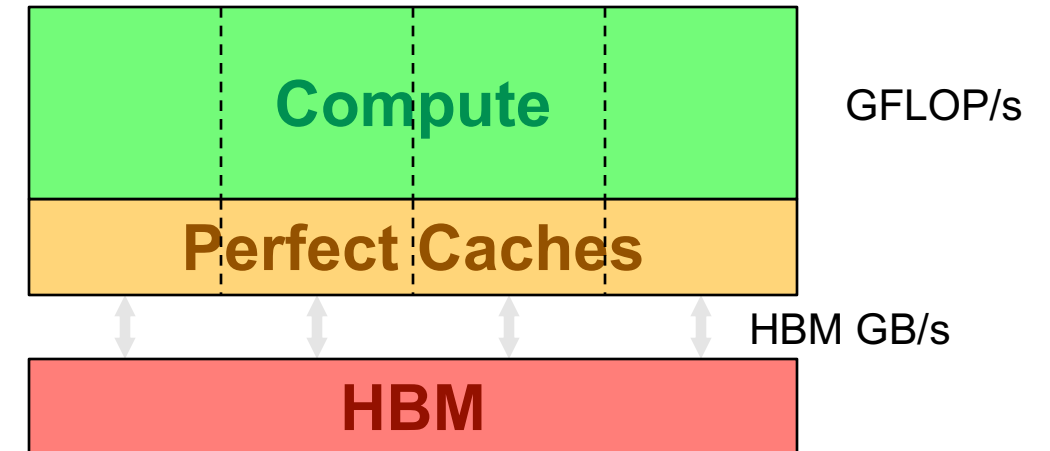
# Roofline Example #2

- Conversely, 7-point constant coefficient stencil…
  - 7 FLOPs
  - 8 memory references (7 reads, 1 store) per point
  - **AI = 7 / (8*8) = 0.11 FLOPs per byte**
  
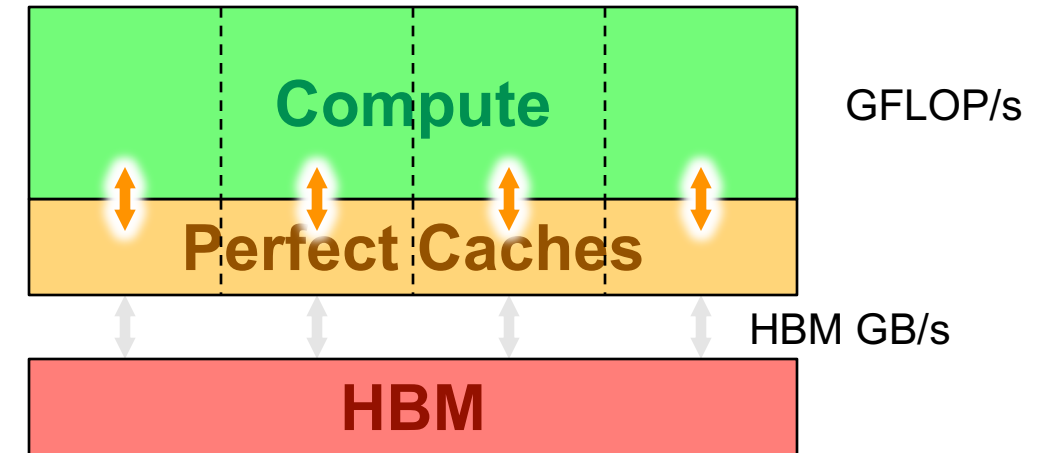    **(measured at the L1)**



```
#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

# Roofline Example #2

- Conversely, 7-point constant coefficient stencil…
  - 7 FLOPs
  - 8 memory references (7 reads, 1 store) per point
  - **Ideally, cache will filter all but 1 read and 1 write per point**
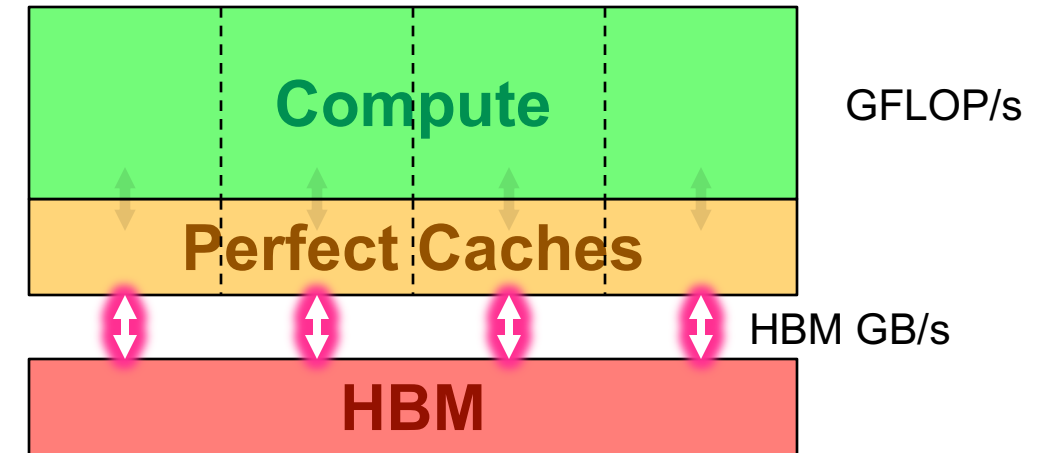
```
#pragma omp parallel for
for(k=1;k<dim+1;k++){
for(j=1;j<dim+1;j++){
for(i=1;i<dim+1;i++){
   new[k][j][i] = -6.0*old[k  ][j  ][i  ]
                  + old[k  ][j  ][i-1]
                  + old[k  ][j  ][i+1]
                  + old[k  ][j-1][i  ]
                  + old[k  ][j+1][i  ]
                  + old[k-1][j  ][i  ]
                  + old[k+1][j  ][i  ]
}}}
```

# Roofline Example #2

- **Conversely, 7-point constant coefficient stencil…**
  - ○ 7 FLOPs
  - ○ 8 memory references (7 reads, 1 store) per point
  - ○ Ideally, cache will filter all but 1 read and 1 write per point
  - ➢ **7 / (8+8) = 0.44 FLOPs per byte (DRAM)**
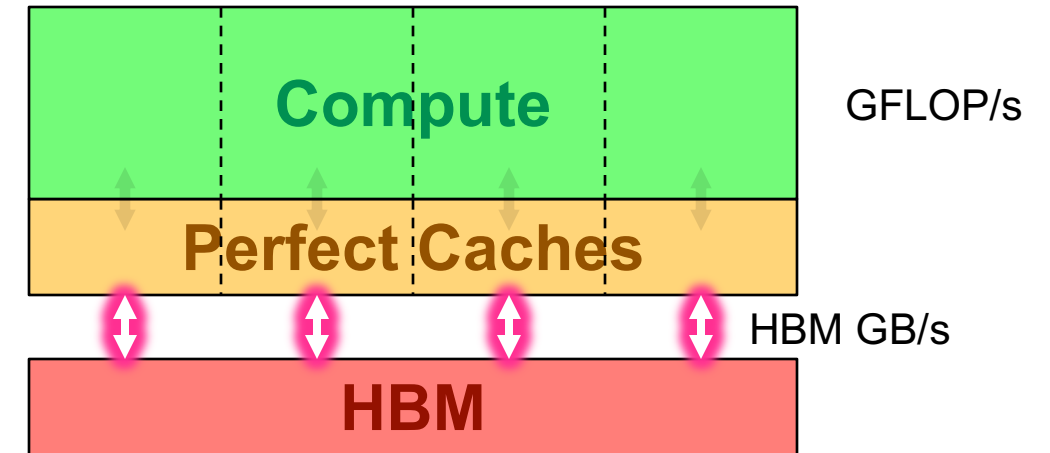


```
#pragma omp parallel for
for(k=1;k<dim+1;k++){
for(j=1;j<dim+1;j++){
for(i=1;i<dim+1;i++){
  new[k][j][i] = -6.0*old[k  ][j  ][i  ]
                    + old[k  ][j  ][i-1]
                    + old[k  ][j  ][i+1]
                    + old[k  ][j-1][i  ]
                    + old[k  ][j+1][i  ]
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

}}}
```

# Roofline Example #2

- ## Conversely, 7-point constant coefficient stencil…

  - 7 FLOPs
  - 8 memory references (7 reads, 1 store) per point
  - Ideally, cache will filter all but 1 read and 1 write per point
  - **7 / (8+8) = 0.44 FLOPs per byte (DRAM)**

    **== memory bound, but 5x the FLOP rate as TRIAD**

```
#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  ];
}}}
```
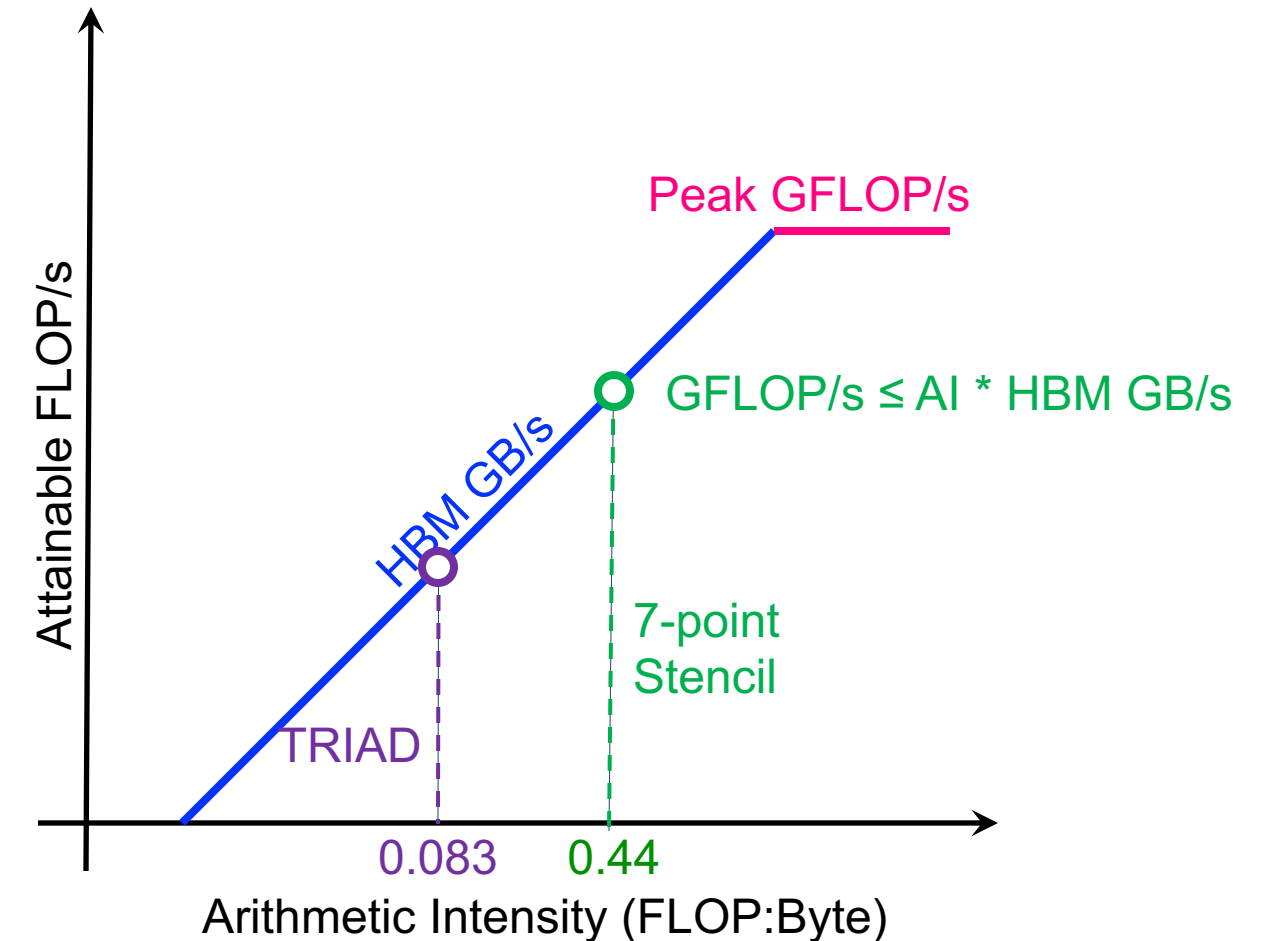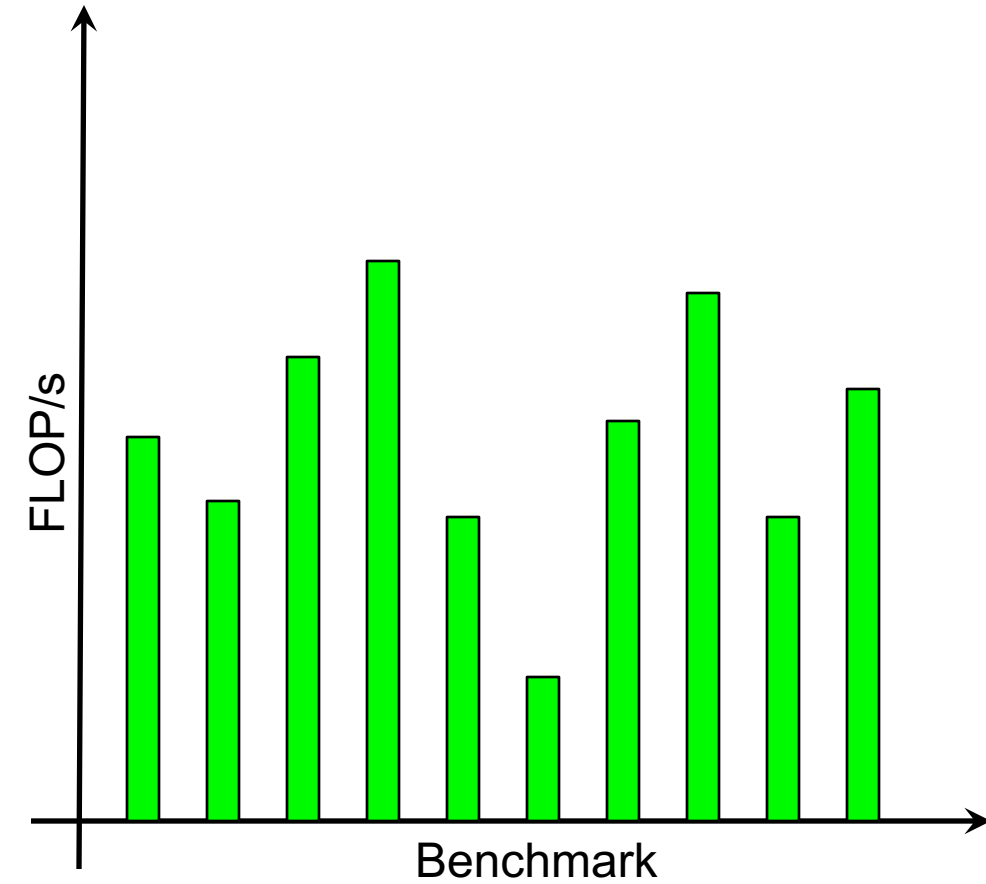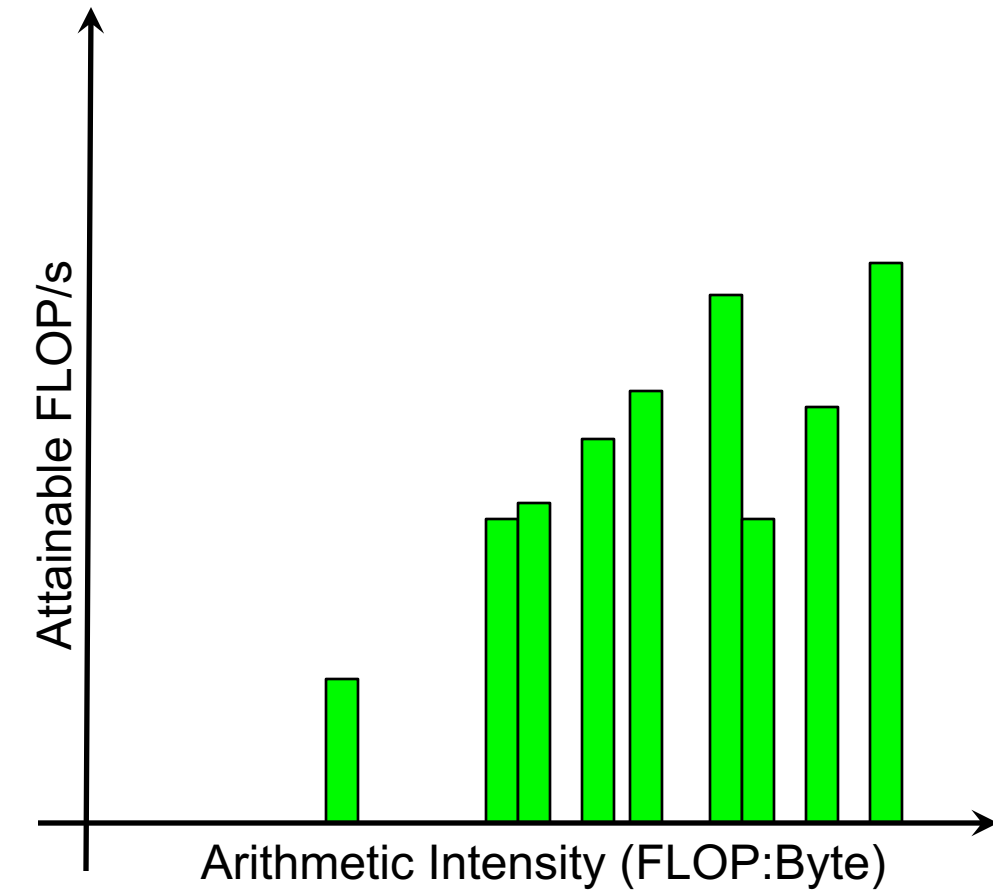


GFLOP/s ≤ AI * HBM GB/s

# Are we getting good performance?

- Think back to our mix of benchmarks…
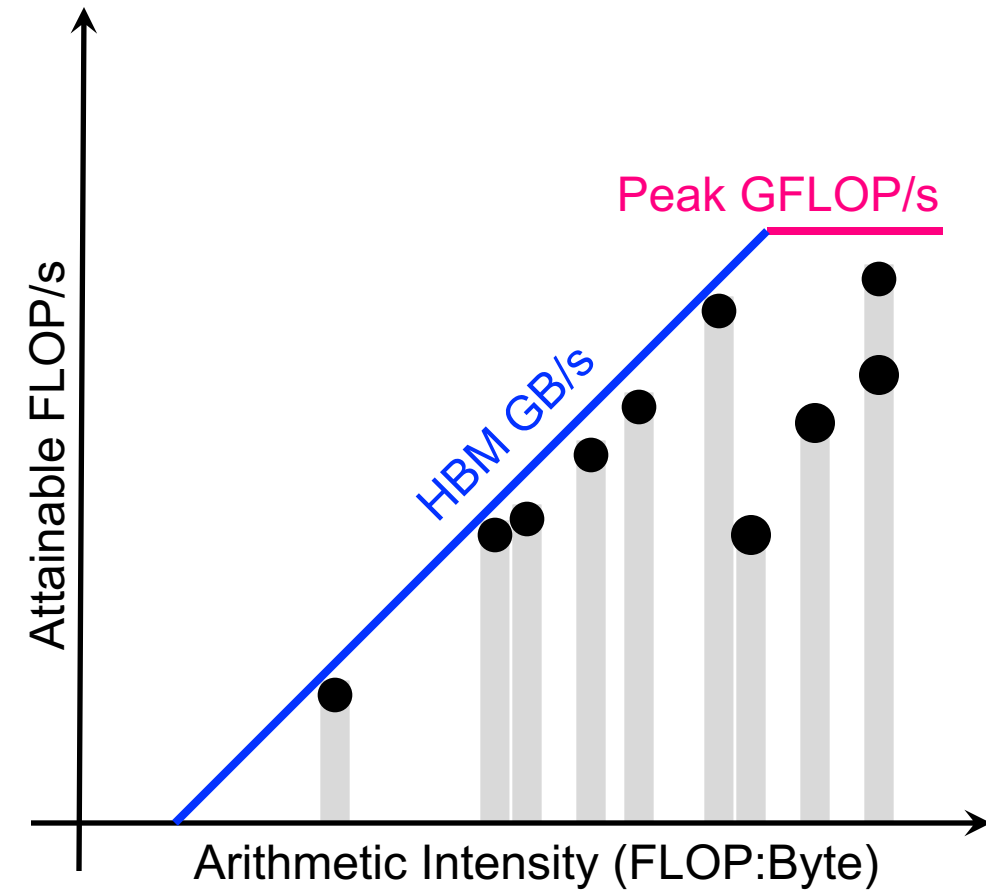
# Are we getting good performance?

- We can sort benchmarks by arithmetic intensity…

BERKELEY LAB

# Are we getting good performance?

- We can sort benchmarks by arithmetic intensity…

- … and compare performance relative to machine capabilities

BERKELEY LAB

# Are we getting good performance?

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

BERKELEY LAB

# Are we getting good performance?

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

    ➢ benchmarks can have **low performance** (GFLOP/s), but make **good use** (%STREAM) of a machine

BERKELEY LAB

# Are we getting good performance?

- Benchmarks near the roofline are making **good use** of computational resources
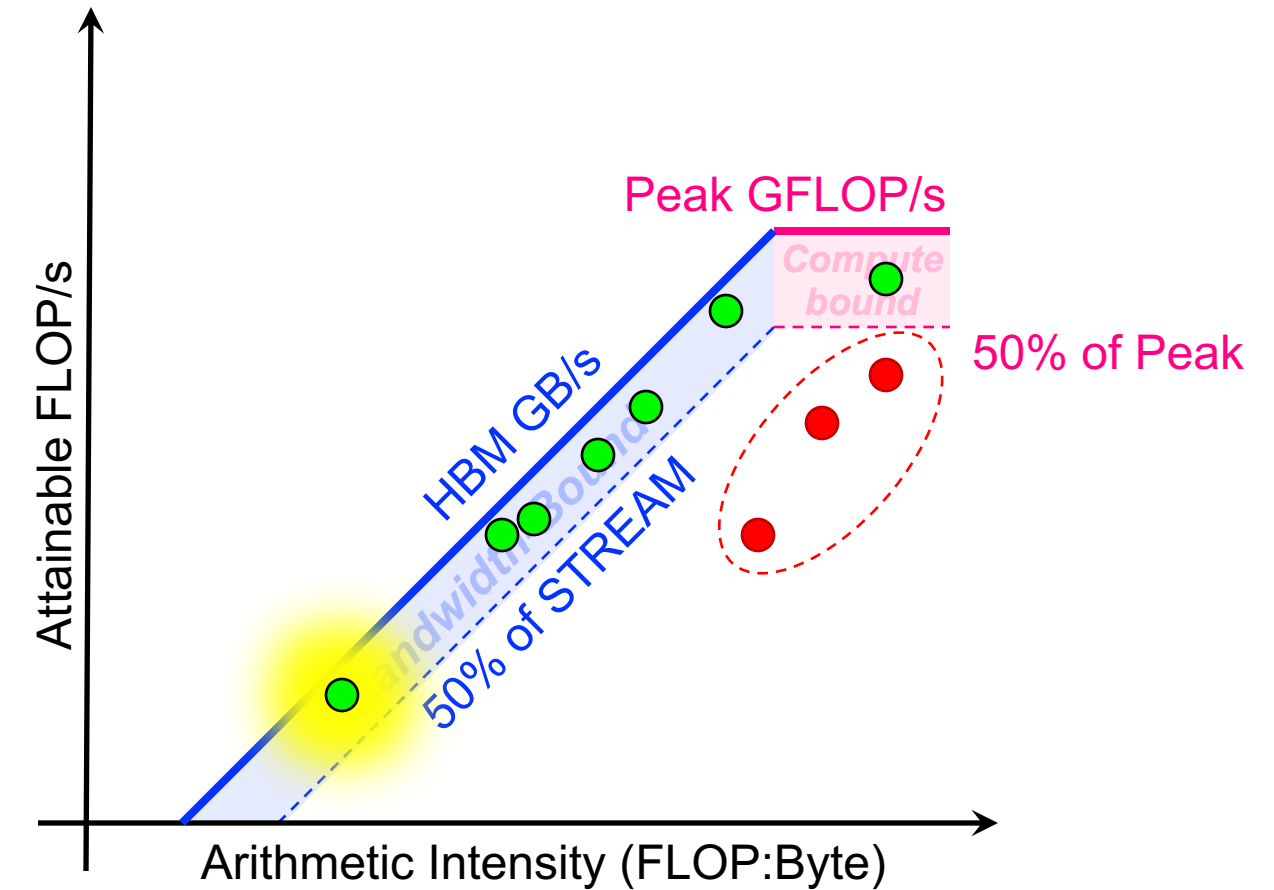
  - benchmarks can have **low performance** (GFLOP/s), but make **good use** (%STREAM) of a machine

  - benchmarks can have **high performance** (GFLOP/s), but still make **poor use** of a machine (%peak)

BERKELEY LAB

# Recap: Roofline is made of two components

- **Machine Model**

  o Lines defined by peak GB/s and GF/s
  (**Benchmarking**)

  o Unique to each architecture

  o Common to all apps on that architecture

BERKELEY LAB

# Recap: Roofline is made of two components

- **Machine Model**
    - Lines defined by peak GB/s and GF/s (**Benchmarking**)
    - Unique to each architecture
    - Common to all apps on that architecture
- **Application Characteristics**
    - Dots defined by application GFLOPs, GBs, and run time

      (**Application Instrumentation**)
    - Unique to each application
    - Unique to each architecture

# Recap: Optimization Strategy

1. Get to the Roofline

# Recap: Optimization Strategy

1. Get to the Roofline
2. Reduce Data movement when bandwidth-limited
   - Bandwidth-bound implies run time is tied to data movement and peak GB/s.
   - *Optimizations that reduce data movement will improve performance*

BERKELEY LAB

# Recap: Optimization Strategy

1. Get to the Roofline
2. Reduce Data movement when bandwidth-limited
   - Bandwidth-bound implies run time is tied to data movement and peak GB/s.
   - Optimizations that reduce data movement will improve performance
3. Reduce the number of #FLOPs when compute-bound
   - Compute bound implies run time is tied to #FLOPs and peak GFLOP/s
   - *Optimizations that eliminate FLOPs will improve time-to-solution (but may reduce GFLOP/s)*
   - Subtlety, this will reduce AI, but increase performance

Peak GFLOP/s

50% of Peak

Attainable FLOP/s

HBM GB/s

50% of STREAM

Arithmetic Intensity (FLOP:Byte)

BERKELEY LAB

# How can performance ever be below the Roofline?

# Theoretical vs. Empirical

- **Theoretical Roofline:**
  - Pin bandwidth == bits * GHz
  - Peak FLOPs == FPUs * GHz
  - 1 C++ FLOP = 1 ISA FLOP
  - Data movement = Compulsory Misses

# Theoretical vs. Empirical / Benchmarking

- **Theoretical Roofline:**
  - Pin bandwidth == bits * GHz
  - Peak FLOPs == FPUs * GHz
  - 1 C++ FLOP = 1 ISA FLOP
  - Data movement = Compulsory Misses
- **Empirical Roofline:**
  - Empirical bandwidth (STREAM) <= theoretical
  - Empirical peak FLOP/s <= theoretical

# Theoretical vs. Empirical / FLOPs

- **Theoretical Roofline:**
  - Pin bandwidth == bits * GHz
  - Peak FLOPs == FPUs * GHz
  - 1 C++ FLOP = 1 ISA FLOP
  - Data movement = Compulsory Misses

- **Empirical Roofline:**
  - Empirical bandwidth (STREAM) <= theoretical
  - Empirical peak FLOP/s <= theoretical
  - 1 C++ FLOP >= 1 ISA FLOP (e.g. divide)

Theoretical GFLOP/s

Empirical GFLOP/s

Attainable FLOP/s

Theoretical GB/s

Empirical GB/s

AI using empirical FLOPs

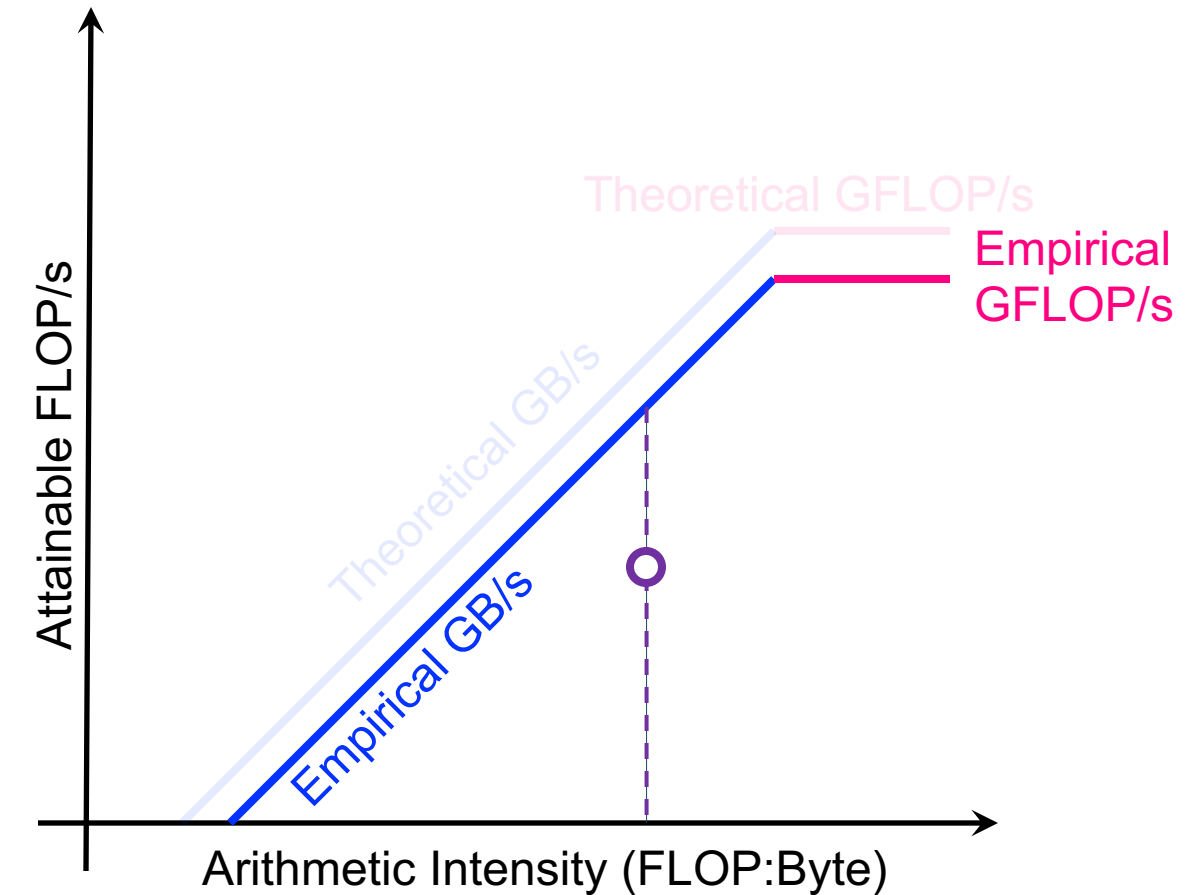Arithmetic Intensity (FLOP:Byte)

BERKELEY LAB

# Theoretical vs. Empirical / Bytes

- **Theoretical Roofline:**
    - Pin bandwidth == bits * GHz
    - Peak FLOPs == FPUs * GHz
    - 1 C++ FLOP = 1 ISA FLOP
    - Data movement = Compulsory Misses
- **Empirical Roofline:**
    - Empirical bandwidth (STREAM) <= theoretical
    - Empirical peak FLOP/s <= theoretical
    - 1 C++ FLOP >= 1 ISA FLOP (e.g. divide)
    - Data movement >> Compulsory Misses

# Theoretical vs. Empirical / Bytes

- Theoretical Roofline:
  o Pin bandwidth == bits * GHz
  o Peak FLOPs == FPUs * GHz
  o 1 C++ FLOP = 1 ISA FLOP
  o Data movement = Compulsory Misses
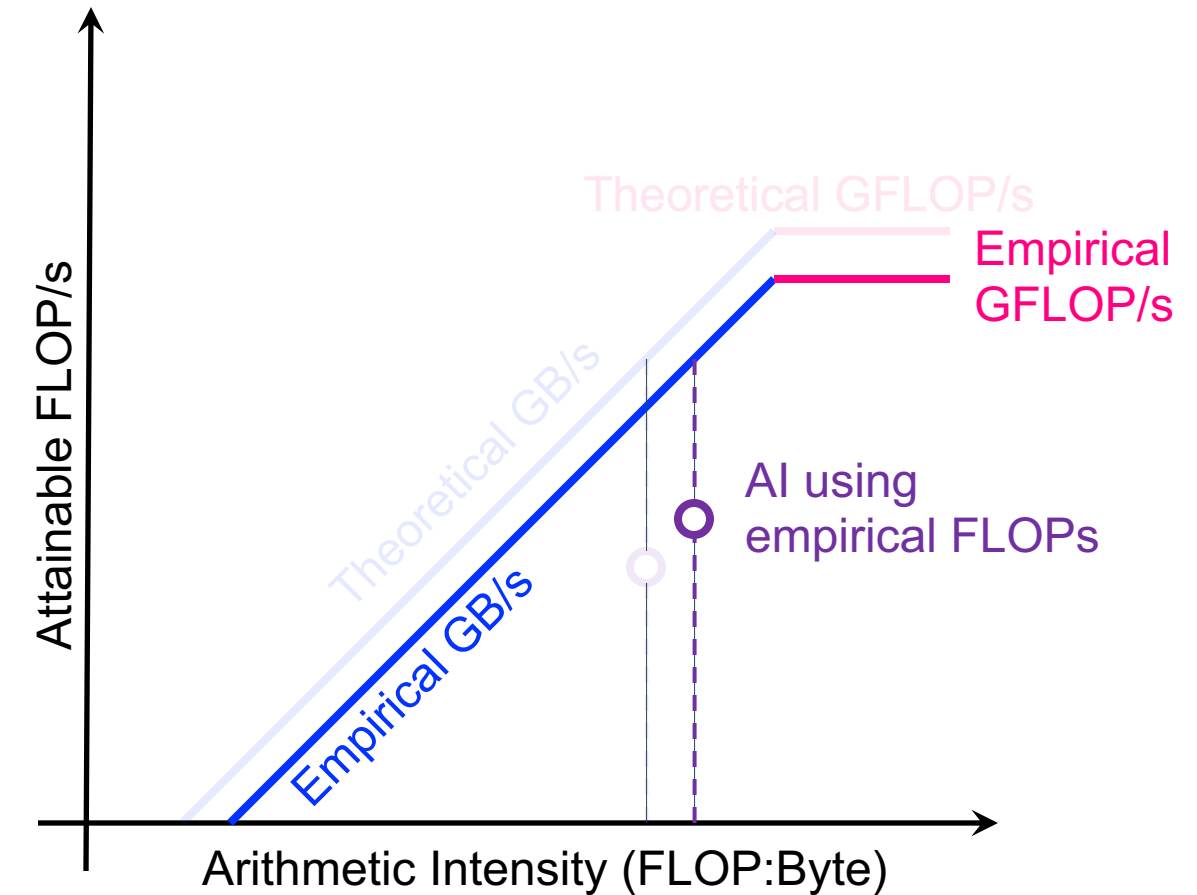
- Empirical Roofline:
  o Empirical bandwidth (STREAM) <= theoretical
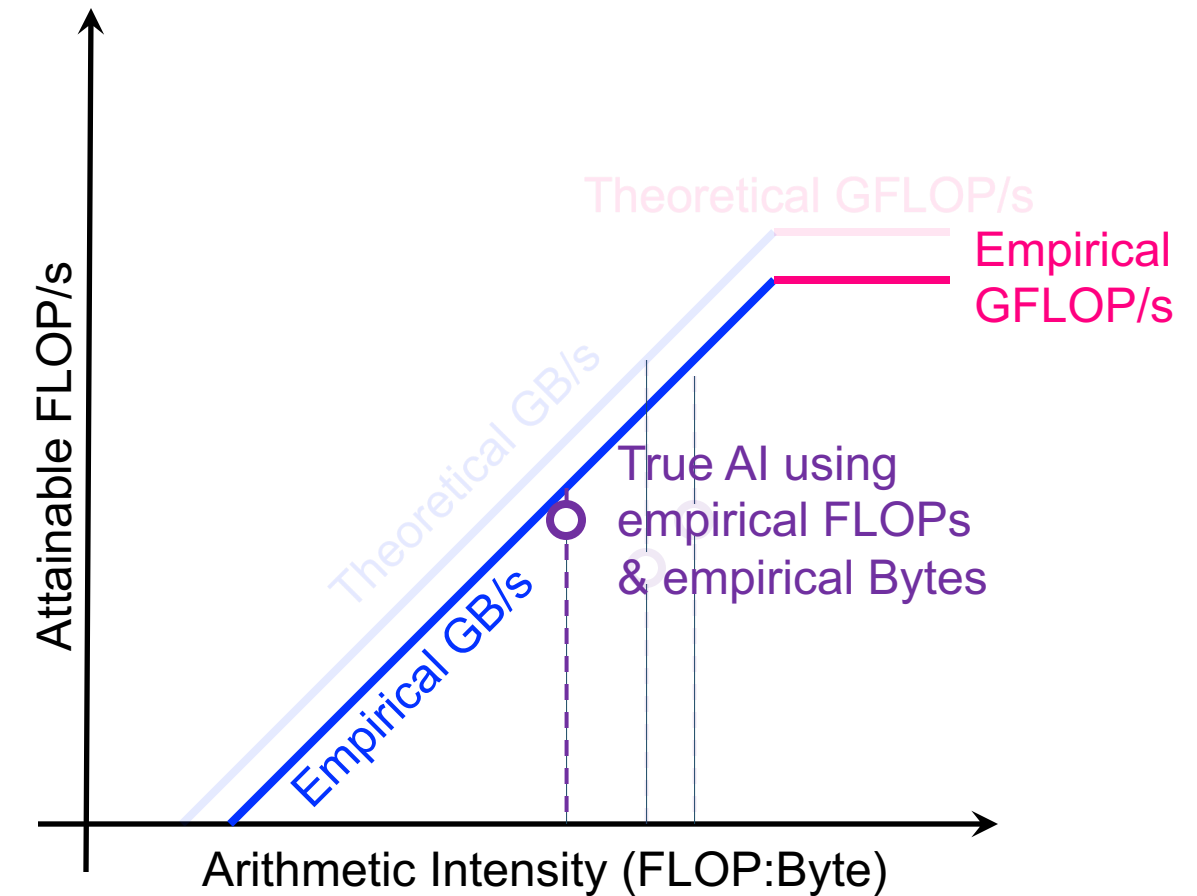  o Empirical peak FLOP/s <= theoretical
  o 1 C++ FLOP >= 1 ISA FLOP (e.g. divide)
  o Data movement >> Compulsory Misses

- **Use benchmarking tools to construct the Roofline model (ceilings)**
- **Use Profiling tools to populate the Roofline model (dots)**

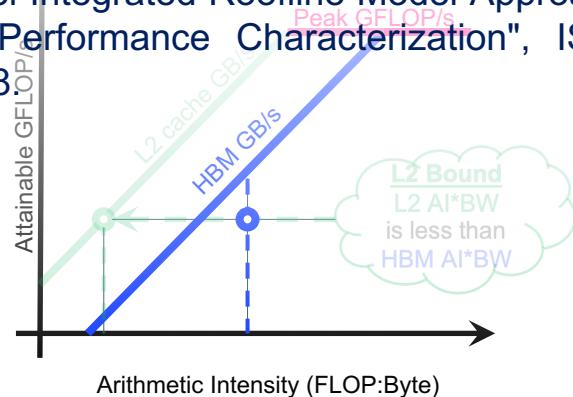# How else can performance be below the Roofline?

*Simple DRAM model can be insufficient for a variety of reasons…*

**DRAM's not the bottleneck…**
- Cache bandwidth and cache locality
- PCIe bandwidth

*…The Hierarchical Roofline Model*

T. Koskela, Z. Matveev, C. Yang, A. Adedoyin, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, S. Williams, "A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization", ISC, 2018.

**Not enough of Vector/Tensor instr.**
- No FMA
- Mixed Precision
- No Tensor Core OPs

*… Additional Ceilings*

C. Yang, T. Kurth, S. Williams, "Hierarchical Roofline analysis for GPUs: Accelerating performance optimization for the NERSC-9 Perlmutter system", CCPE, 2019.

**Lack of Parallelism…**
- Idle Cores/SMs
- Insufficient ILP/TLP
- Divergence and Predication

*… Roofline Scaling Trajectories*

K. Ibrahim, S. Williams, L. Oliker, "Performance Analysis of GPU Programming Models using the Roofline Scaling Trajectories", BEST PAPER, Bench, 2019.

**Integer-heavy Codes…**
- Non-FP inst. impede FLOPs
- No FP instructions

*…The Instruction Roofline Model*

N. Ding, S. Williams, "An Instruction Roofline Model for GPUs", BEST PAPER, PMBS, 2019.

BERKELEY LAB

# Below the Roofline?
## Memory Hierarchy and Cache Bottlenecks

# Memory Hierarchy

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

Core

L1 D$

L2 D$

L3 D$

DRAM

# Memory Hierarchy

- CPUs/GPUs have different bandwidths for each level

**Bandwidth**

Core

L1 GB/s

L1 D$

L2 GB/s

L2 D$

L3 GB/s

L3 D$

DRAM GB/s

DRAM

BERKELEY LAB

# Memory Hierarchy

- CPUs/GPUs have different bandwidths for each level
  - different **machine balances** for each level

**Machine Balance**

GFLOP/s
L1 GB/s

GFLOP/s
L2 GB/s

GFLOP/s
L3 GB/s

GFLOP/s
DRAM GB/s

```
Core

L1 D$

L2 D$

L3 D$

DRAM
```

BERKELEY LAB

# Memory Hierarchy

- **CPUs/GPUs have different bandwidths for each level**
  - different machine balances for each level

- **Applications have locality in each level**
  - different **data movements** for each level

**Machine Balance**    **Data Movement**

```
        Core
GFLOP/s                    L1 GB
L1 GB/s
        L1 D$
GFLOP/s                    L2 GB
L2 GB/s
        L2 D$
GFLOP/s                    L3 GB
L3 GB/s
        L3 D$
GFLOP/s                    DRAM GB
DRAM GB/s
        DRAM
```

BERKELEY LAB

# Memory Hierarchy

- **CPUs/GPUs have different bandwidths for each level**
  - different machine balances for each level

- **Applications have locality in each level**
  - different data movements for each level
  - different **arithmetic intensity** for each level

**Machine Balance**          **Arithmetic Intensity**

Core

$\dfrac{GFLOP/s}{L1\ GB/s}$          $\dfrac{GFLOPs}{L1\ GB}$

**L1 D$**

$\dfrac{GFLOP/s}{L2\ GB/s}$          $\dfrac{GFLOPs}{L2\ GB}$

**L2 D$**

$\dfrac{GFLOP/s}{L3\ GB/s}$          $\dfrac{GFLOPs}{L3\ GB}$

**L3 D$**

$\dfrac{GFLOP/s}{DRAM\ GB/s}$          $\dfrac{GFLOPs}{DRAM\ GB}$

**DRAM**

BERKELEY LAB

# Cache Bottlenecks

- For each additional level of the memory hierarchy, we can add another term to our model…

$$\text{GFLOP/s} = \min \begin{cases} \text{Peak GFLOP/s} \\ \text{AI}_{\text{DRAM}} * \text{DRAM GB/s} \end{cases}$$

AI$_x$ (Arithmetic Intensity at level "x") = FLOPs / Bytes (moved to/from level "x" )

BERKELEY LAB

# Cache Bottlenecks

- For each additional level of the memory hierarchy, we can add another term to our model…

$$GFLOP/s = \min \begin{cases} \textbf{Peak GFLOP/s} \\ \textbf{AI}_{DRAM} \textbf{ * DRAM GB/s} \\ \textbf{AI}_{L2} \textbf{ * L2 GB/s} \end{cases}$$

$AI_x$ (Arithmetic Intensity at level "x") = FLOPs / Bytes (moved to/from level "x" )

BERKELEY LAB

# Cache Bottlenecks

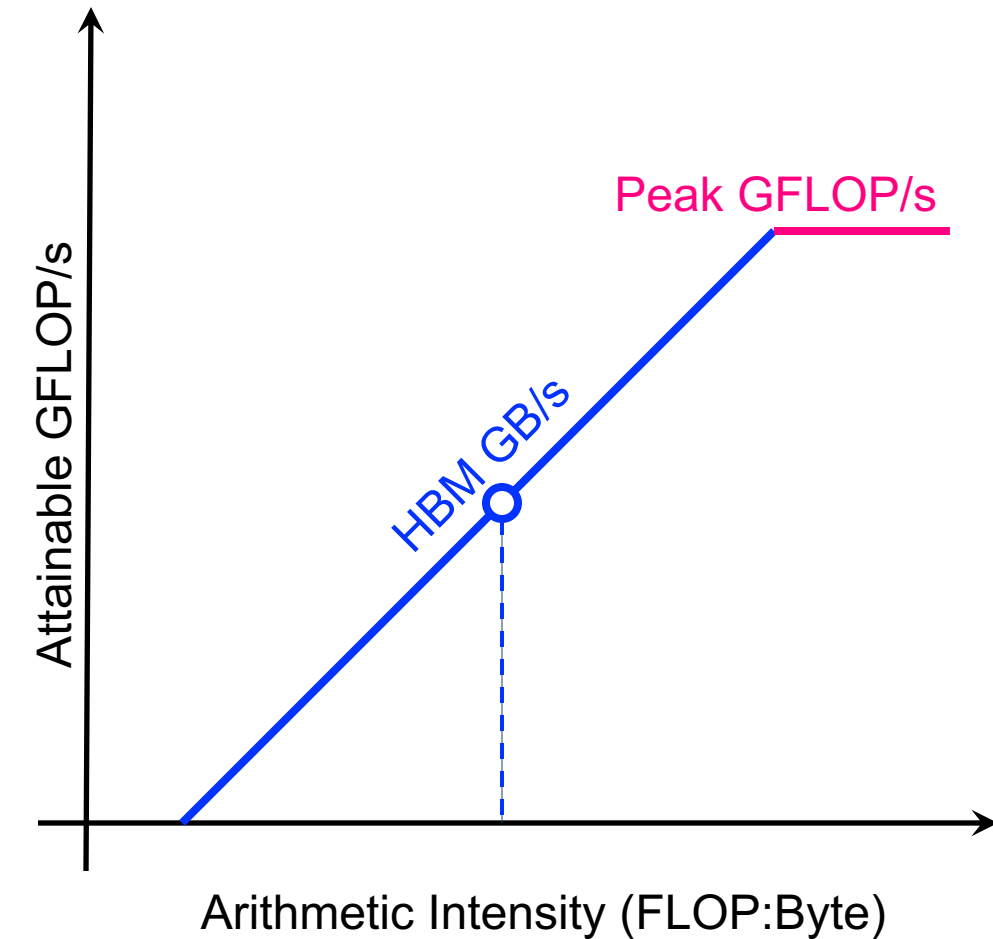- For each additional level of the memory hierarchy, we can add another term to our model…

$$\text{GFLOP/s} = \min \begin{cases} \text{Peak GFLOP/s} \\ \text{AI}_{\text{DRAM}} * \text{DRAM GB/s} \\ \text{AI}_{\text{L2}} * \text{L2 GB/s} \\ \text{AI}_{\text{L1}} * \text{L1 GB/s} \end{cases}$$

$\text{AI}_x$ (Arithmetic Intensity at level "x") = FLOPs / Bytes (moved to/from level "x" )

BERKELEY LAB

# Cache Bottlenecks

- Plot equation in a single figure…
  - "**Hierarchical Roofline**" Model

T. Koskela, Z. Matveev, C. Yang, A. Adedoyin, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, S. Williams, "A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization", ISC, 2018.

# Cache Bottlenecks

- Plot equation in a single figure…
  - "**Hierarchical Roofline**" Model
  - Bandwidth ceiling (diagonal line) for each level of memory

T. Koskela, Z. Matveev, C. Yang, A. Adedoyin, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, S. Williams, "A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization", ISC, 2018.
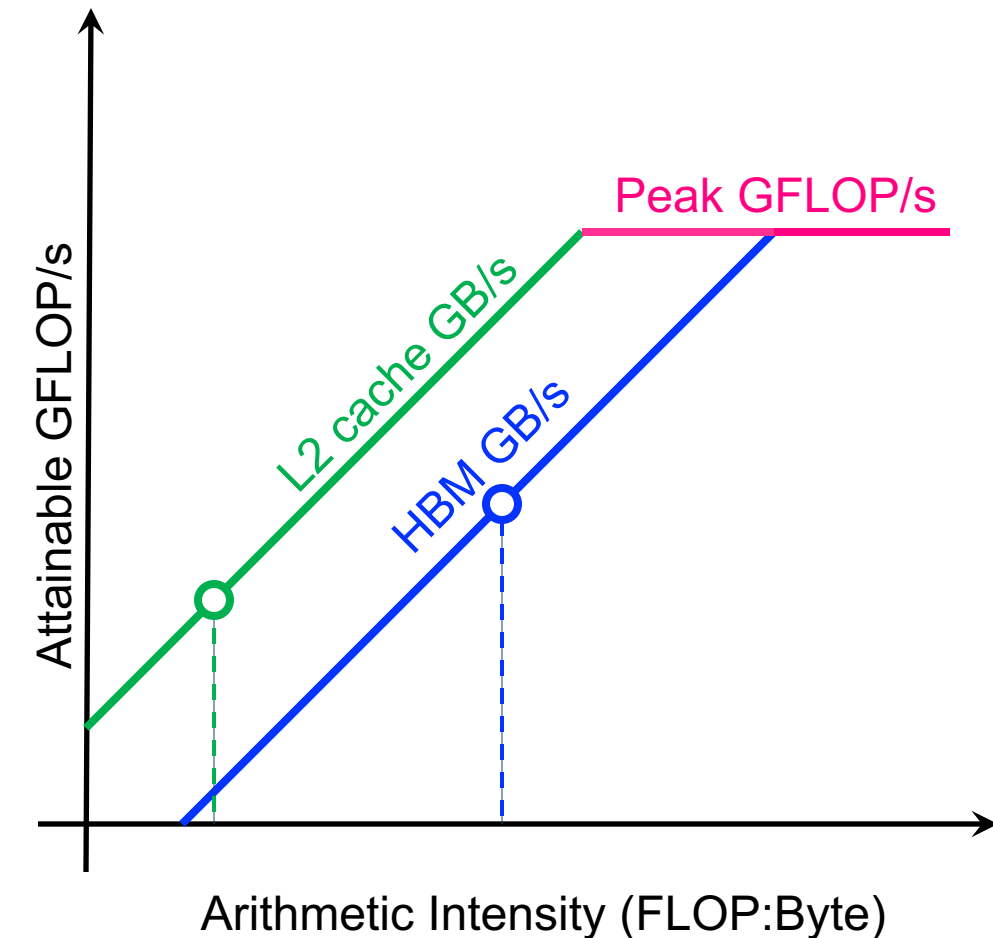
BERKELEY LAB

# Cache Bottlenecks

- Plot equation in a single figure…
  - "**Hierarchical Roofline**" Model
  - Bandwidth ceiling (diagonal line) for each level of memory
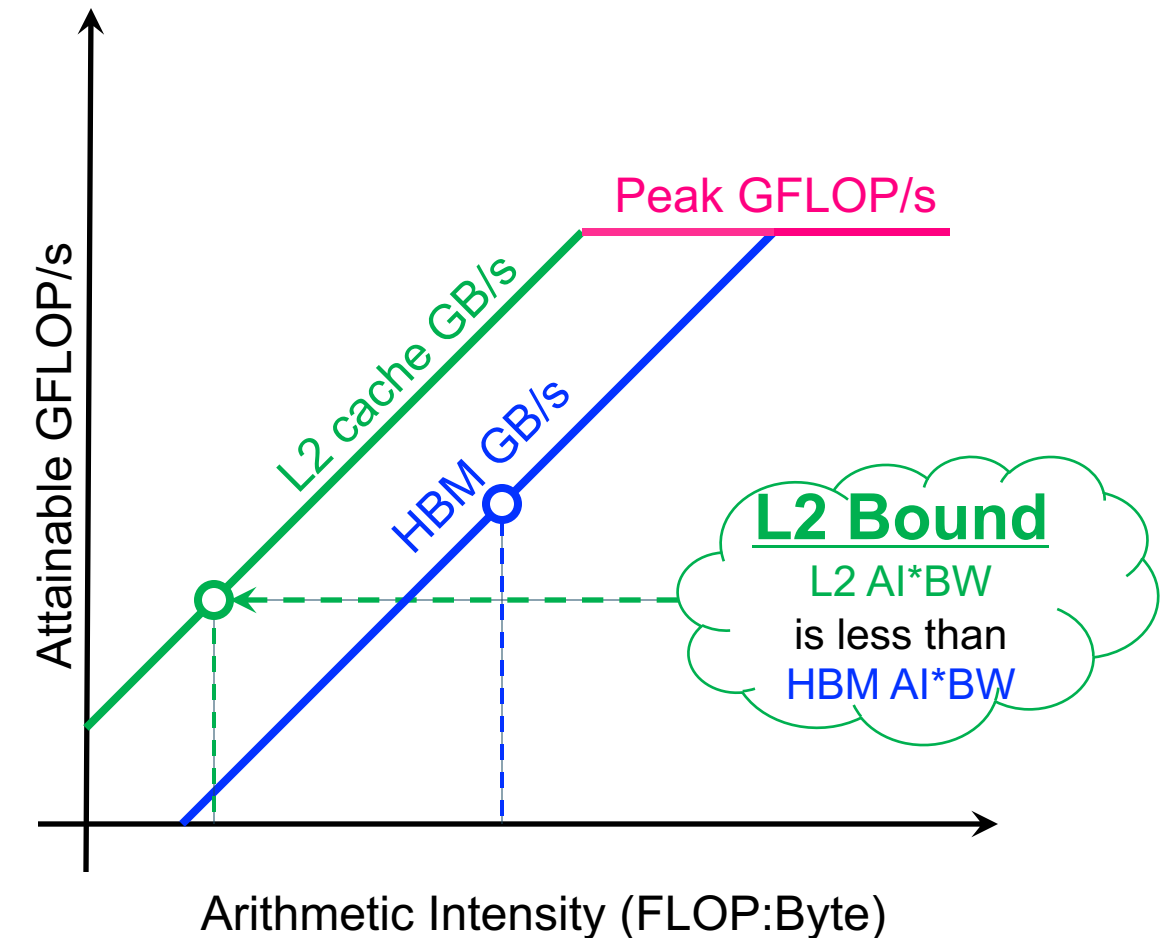  - Arithmetic Intensity (dot) for each level of memory

BERKELEY LAB

# Cache Bottlenecks

- Plot equation in a single figure…
  - "**Hierarchical Roofline**" Model
  - Bandwidth ceiling (diagonal line) for each level of memory
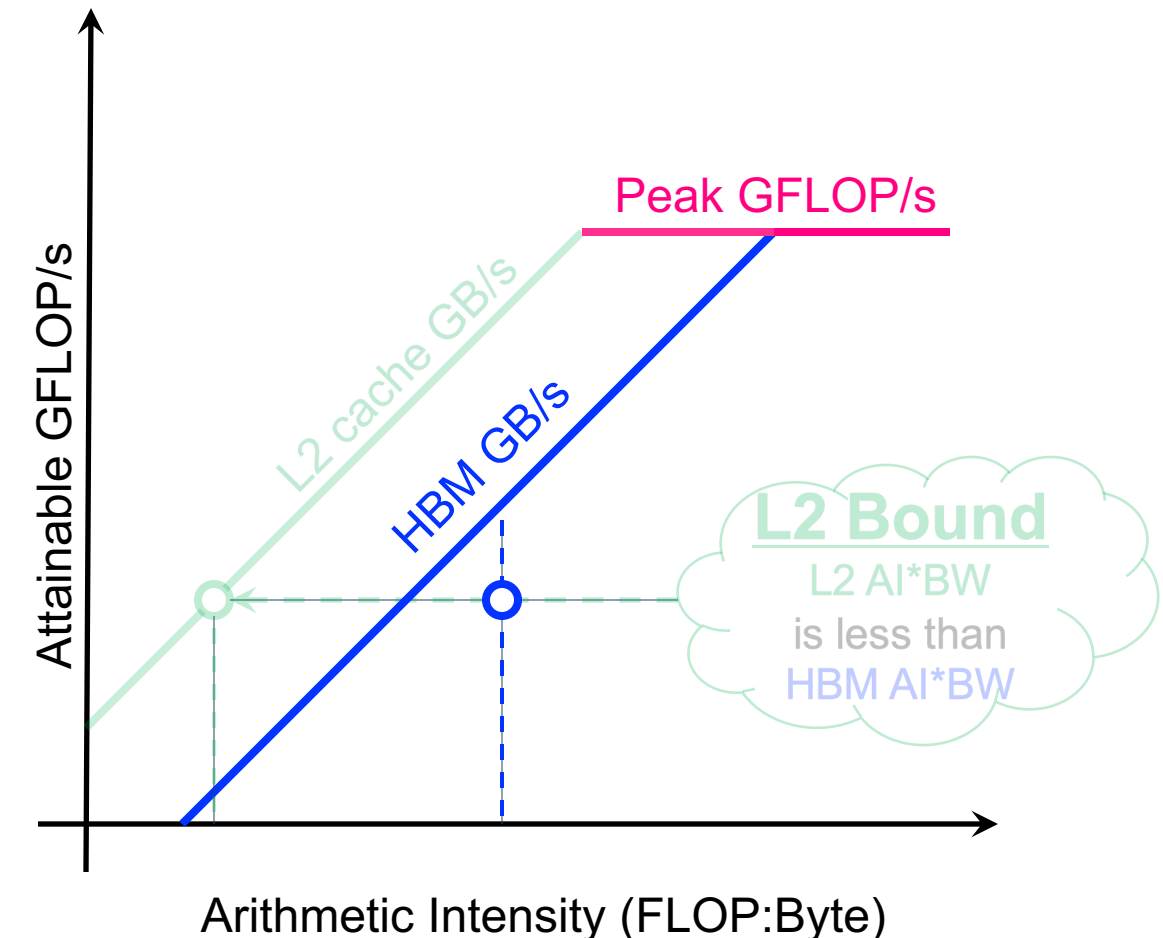  - Arithmetic Intensity (dot) for each level of memory
  - **performance is ultimately the minimum of these bounds**



*Peak GFLOP/s*

*L2 cache GB/s*

*HBM GB/s*

*Attainable GFLOP/s*

*Arithmetic Intensity (FLOP:Byte)*

**L2 Bound**
L2 AI*BW
is less than
HBM AI*BW

T. Koskela, Z. Matveev, C. Yang, A. Adedoyin, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, S. Williams, "A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization", ISC, 2018.

BERKELEY LAB

# Cache Bottlenecks

- Plot equation in a single figure…
  - "**Hierarchical Roofline**" Model
  - Bandwidth ceiling (diagonal line) for each level of memory
  - Arithmetic Intensity (dot) for each level of memory
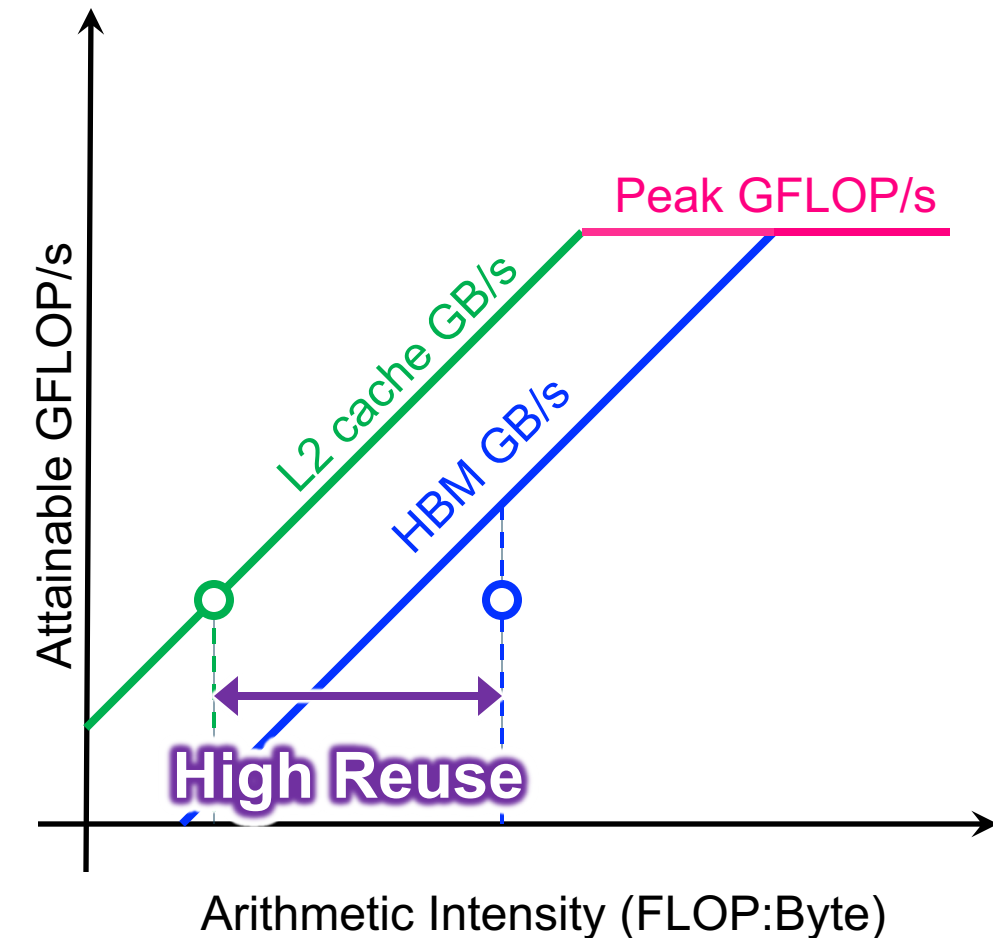  - ➢ **performance is ultimately the minimum of these bounds**

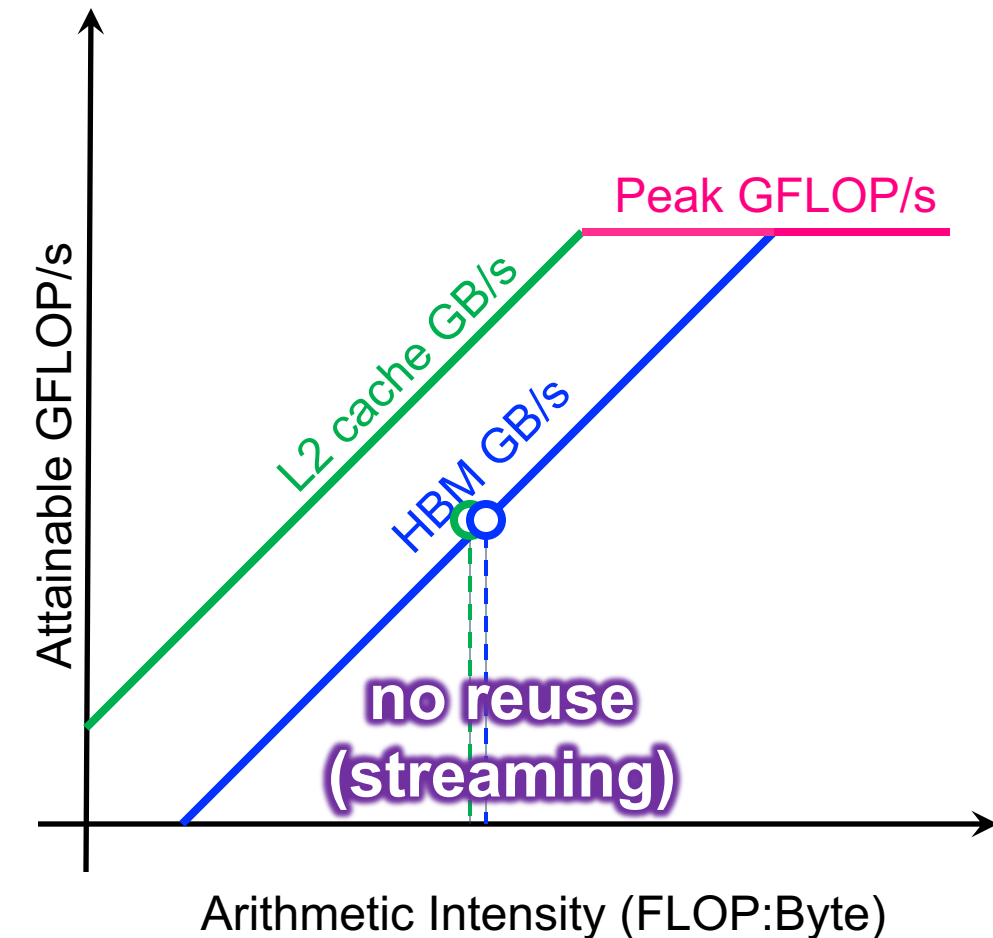- **If L2 bound, we see DRAM dot well below DRAM ceiling**

# Cache Hit Rates

- Widely separated Arithmetic Intensities indicate high reuse in the (L2) cache

T. Koskela, Z. Matveev, C. Yang, A. Adedoyin, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, S. Williams, "A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization", ISC, 2018.

# Cache Hit Rates

- Widely separated Arithmetic Intensities indicate high reuse in the (L2) cache

- Similar Arithmetic Intensities indicate effectively no (L2) cache reuse (**== streaming**)

BERKELEY LAB
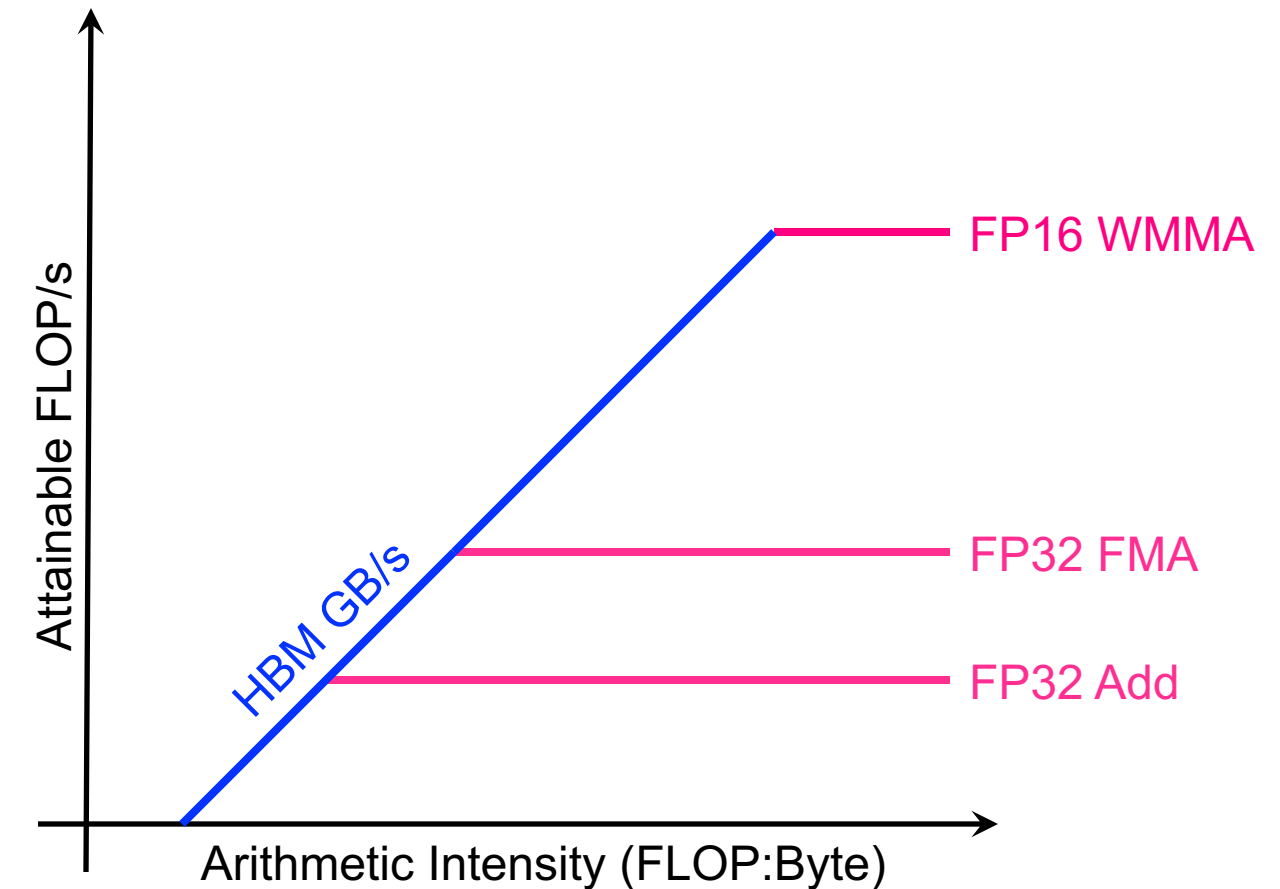
# Below the Roofline?
## Return of CISC

# Return of CISC

- Vectors have their limits (finite DLP, register file energy scales with VL, etc…)
- Death of Moore's Law is reinvigorating Complex Instruction Set Computing (CISC)

- Modern CPUs and GPUs are increasingly reliant on special (fused) instructions that perform multiple operations (fuse common instruction sequences)…
  - FMA (Fused Multiply Add):      $z=a*x+y$      …*z,x,y are vectors or scalars*
  - 4FMA (Quad FMA):      $z=A*x+z$      …*A is a FP32 matrix; x,z are vectors*
  - WMMA (Tensor Core):      $Z=AB+C$      …*A,B are FP16 matrices; Z,C are FP32*

- **Define a set of "ceilings" based on instruction type (all tensor, all FMA, or all FADD)**
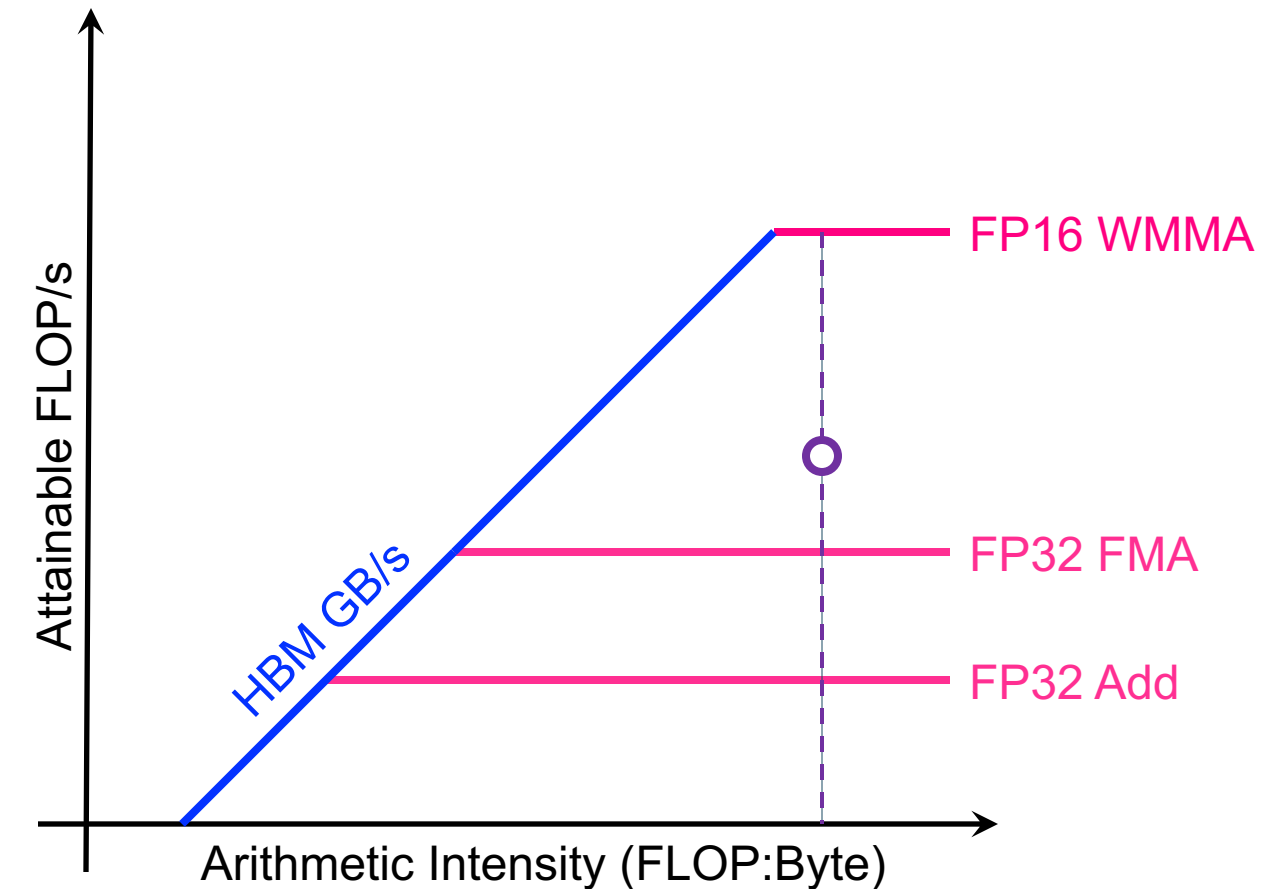
BERKELEY LAB

# Floating-Point and Mixed Precision Ceilings

- Consider NVIDIA Volta GPU
- We may define 3 performance ceilings…
  - 15 TFLOPS for FP32 FMA
  - 7.5 TFLOPs for FP32 Add
  - ~100 TFLOPs for FP16 Tensor
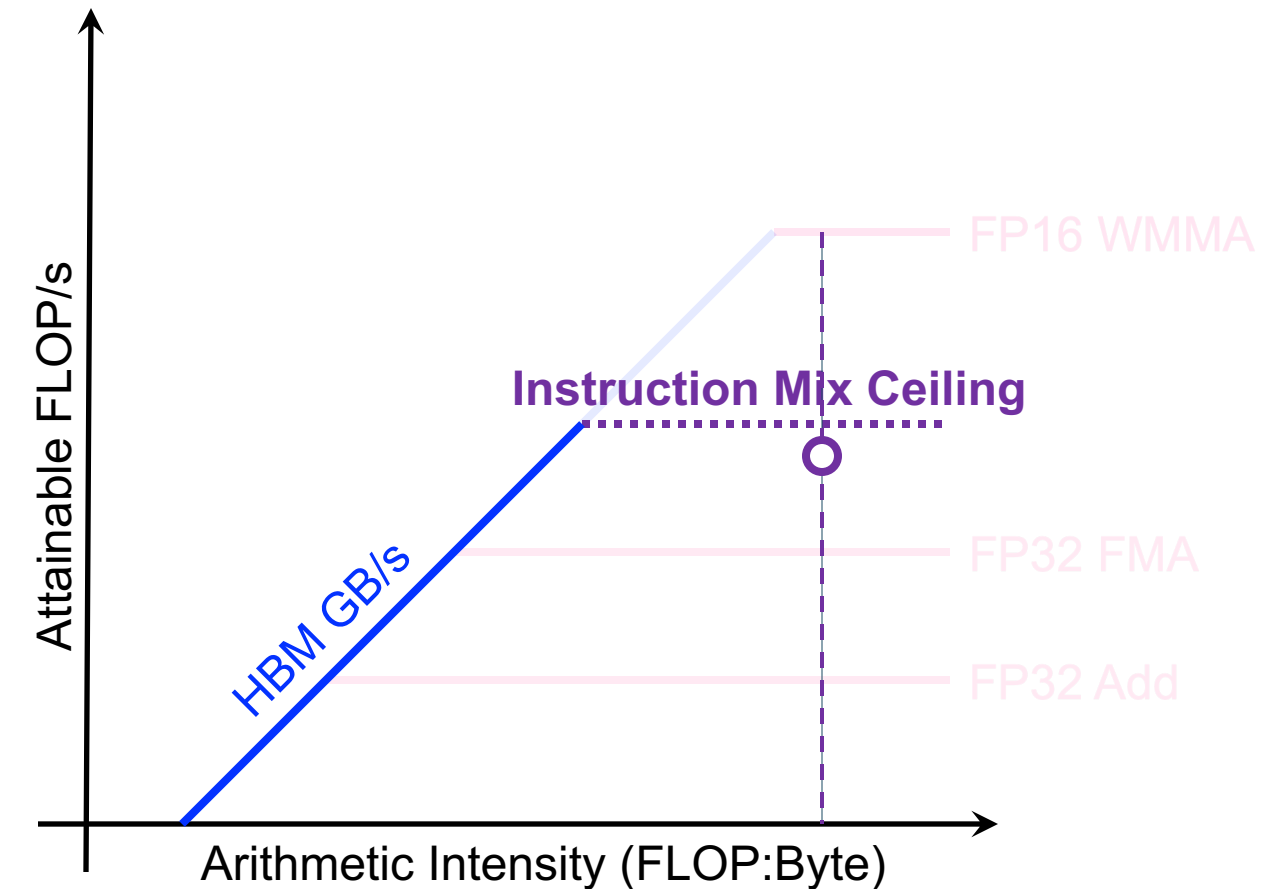
BERKELEY LAB

# Floating-Point and Mixed Precision Ceilings

- When calculating (AI,GFLOP/s), count the <u>total</u> FLOPs from <u>all</u> types of instructions

- DL performance can often be well below nominal Tensor Core peak

# Floating-Point and Mixed Precision Ceilings

- When calculating (AI,GFLOP/s), count the <u>total</u> FLOPs from <u>all</u> types of instructions

- DL performance can often be well below nominal Tensor Core peak

- DL applications are a mix Tensor, FP16, and FP32 instructions

- Thus, there is an <u>**ceiling**</u> on performance defined by the mix of instructions
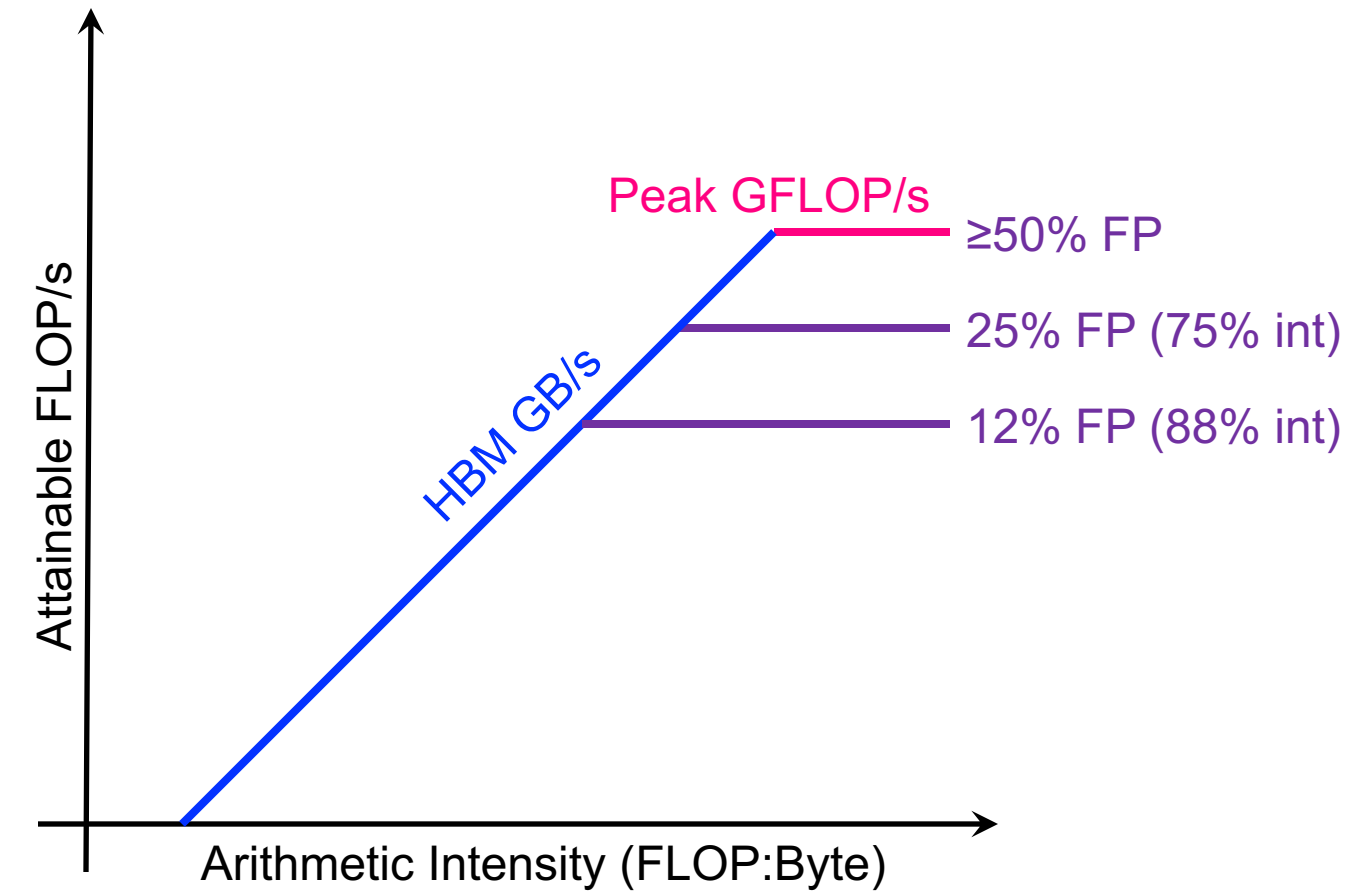
# Below the Roofline?

## FPU Starvation

# FPU Starvation

- CPUs and GPUs have finite instruction fetch/decode/issue bandwidth

- The number of FPUs dictates the FP issue rate required to hit peak

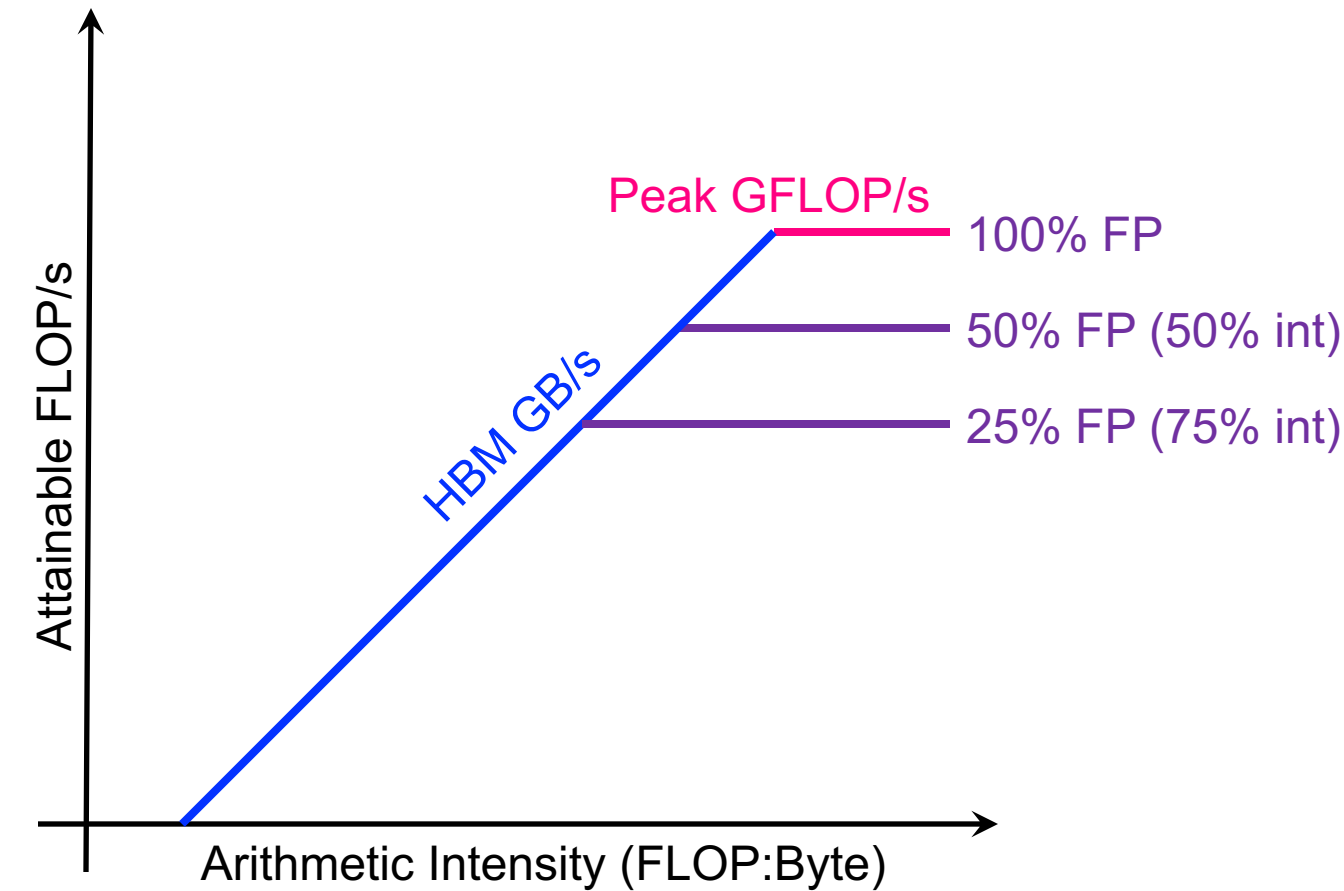➢ **Ratio of these two rates is the minimum FP instruction fraction required to hit peak**

BERKELEY LAB

# FPU Starvation

- ## Consider…
  - 4-issue CPU (or GPU)
  - 2 FP data paths
  - ➢ **>50% of the instructions** must be FP to have any chance at peak performance

# FPU Starvation

- ## Conversely,
  - ○ Keeping 2 FP data paths,
  - ○ but downscaling to 2-issue CPU (or GPU)
  - ➢ **100% of the instructions must be FP to get peak performance**
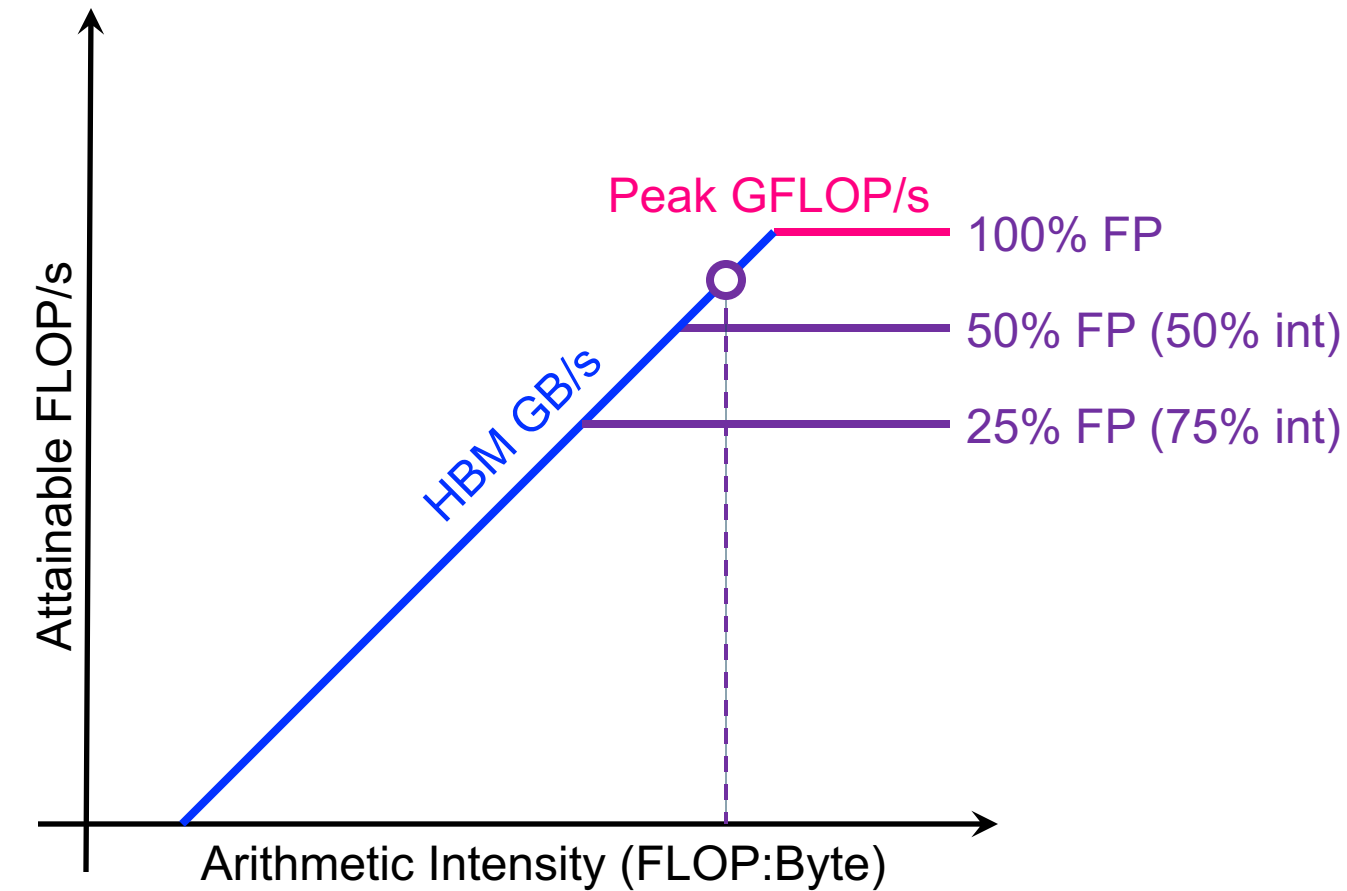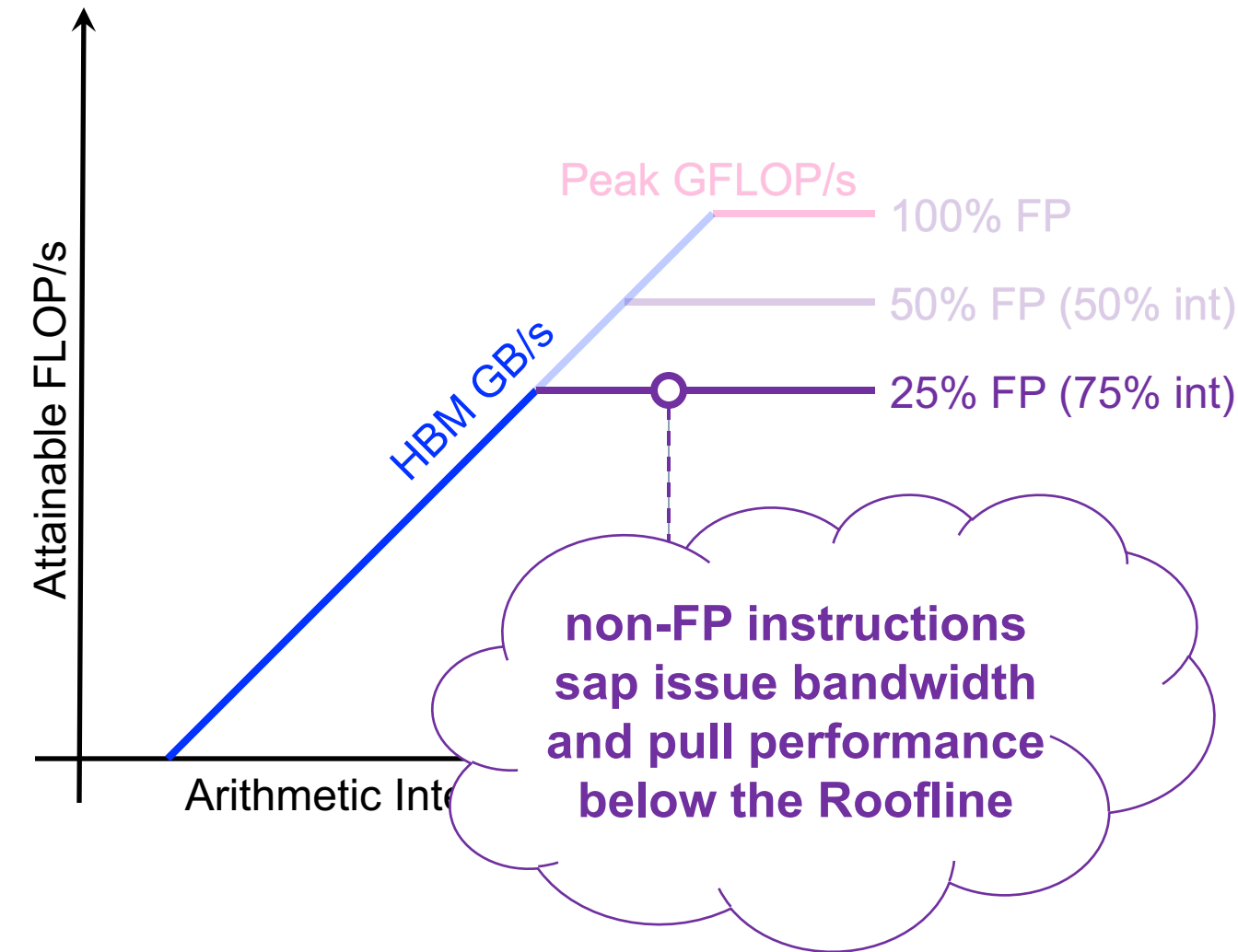
# FPU Starvation

- Conversely,
  - Keeping 2 FP data paths,
  - but downscaling to 2-issue CPU (or GPU)
  - **100% of the instructions must be FP to get peak performance**

# FPU Starvation

- ## Conversely,
  - o Keeping 2 FP data paths,
  - o but downscaling to 2-issue CPU (or GPU)
  - o 100% of the instructions must be FP to get peak performance
  - ➤ **Codes that would have been memory-bound are now decode/issue-bound.**



Peak GFLOP/s

100% FP

50% FP (50% int)

25% FP (75% int)

HBM GB/s

Attainable FLOP/s

Arithmetic Int[

**non-FP instructions sap issue bandwidth and pull performance below the Roofline**
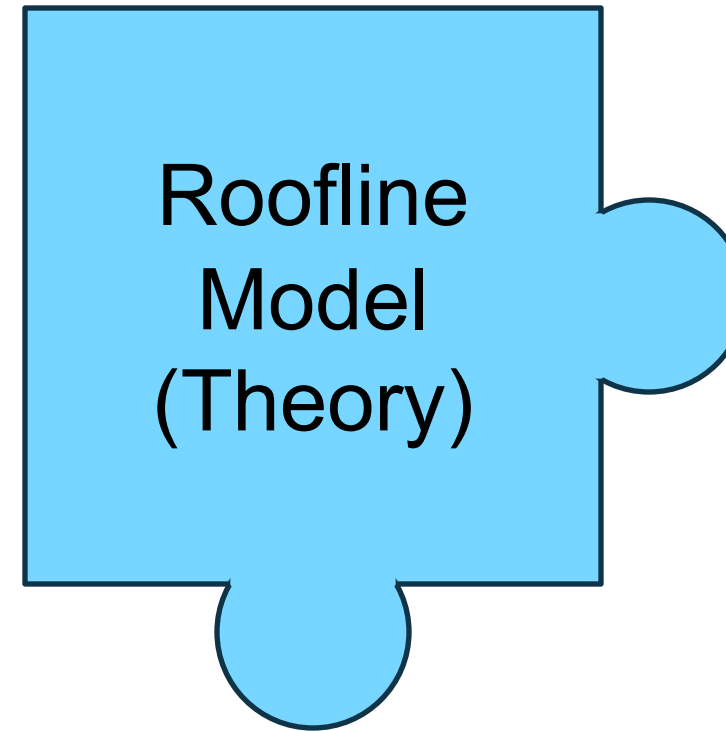
BERKELEY LAB

# Recap

# Recap

- Roofline bounds performance as a function of Arithmetic Intensity
  - Horizontal Lines = Compute Ceilings
  - Diagonal Lines = Bandwidth Ceilings
  - Bandwidth ceilings are always parallel on log-log scale
  - **Collectively, define an upper limit on performance (speed-of-light)**

- Loop Arithmetic Intensity (for each level of memory)
  - **Total FLOPs / Total Data Movement** (for that level of memory)
  - Measure of a loop's temporal locality
  - Includes **all** cache effects

- Plotting loops on the (Hierarchical) Roofline
  - **Each loop has one dot per level of memory**
  - x-coordinate = arithmetic intensity at that level
  - y-coordinate = performance (e.g. GFLOP/s)
  - Proximity to associated ceiling is indicative of a performance bound
  - Proximity of dots to each other is indicative of **streaming** behavior (low cache hit rates)

BERKELEY LAB

# Why would you use Roofline?

- Understand performance differences between Architectures, Programming Models, implementations, etc…
  - Why do some Architectures/Implementations move more data than others?
  - Why do some compilers outperform others?

- Predict performance on future machines / architectures
  - Set realistic performance expectations
  - Drive for HW/SW Co-Design

- Identify performance bottlenecks & motivate software optimizations

- Determine when we're done optimizing code
  - Assess performance relative to machine capabilities
  - Track progress towards optimality
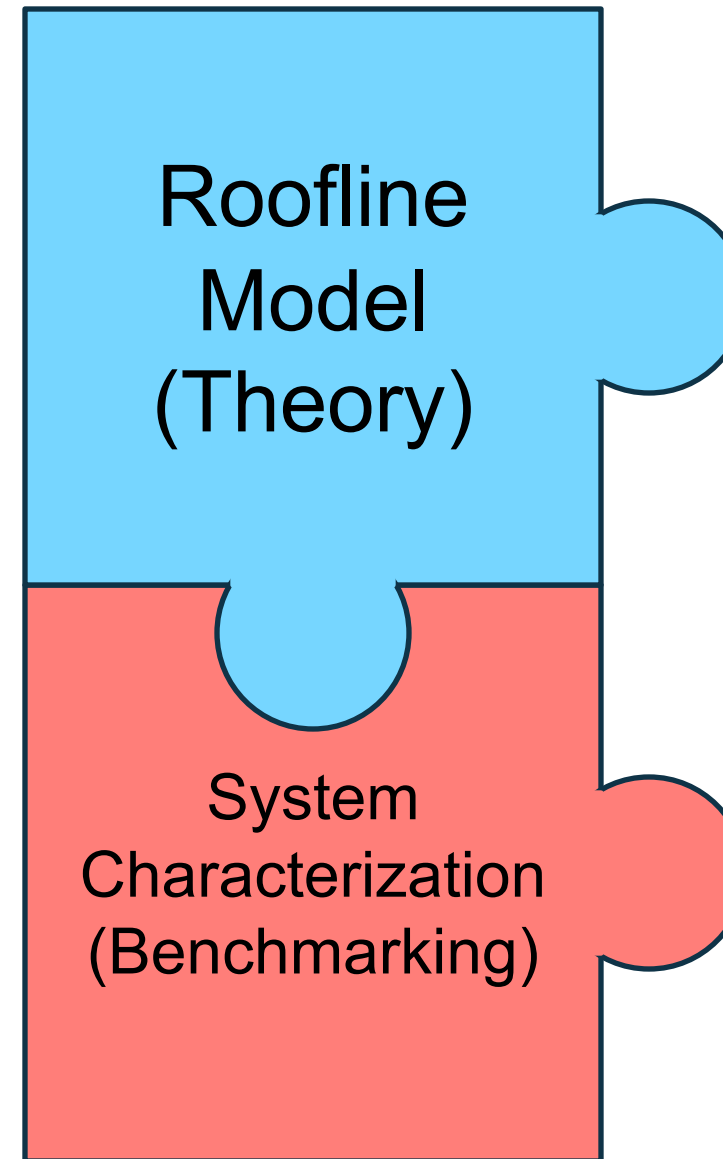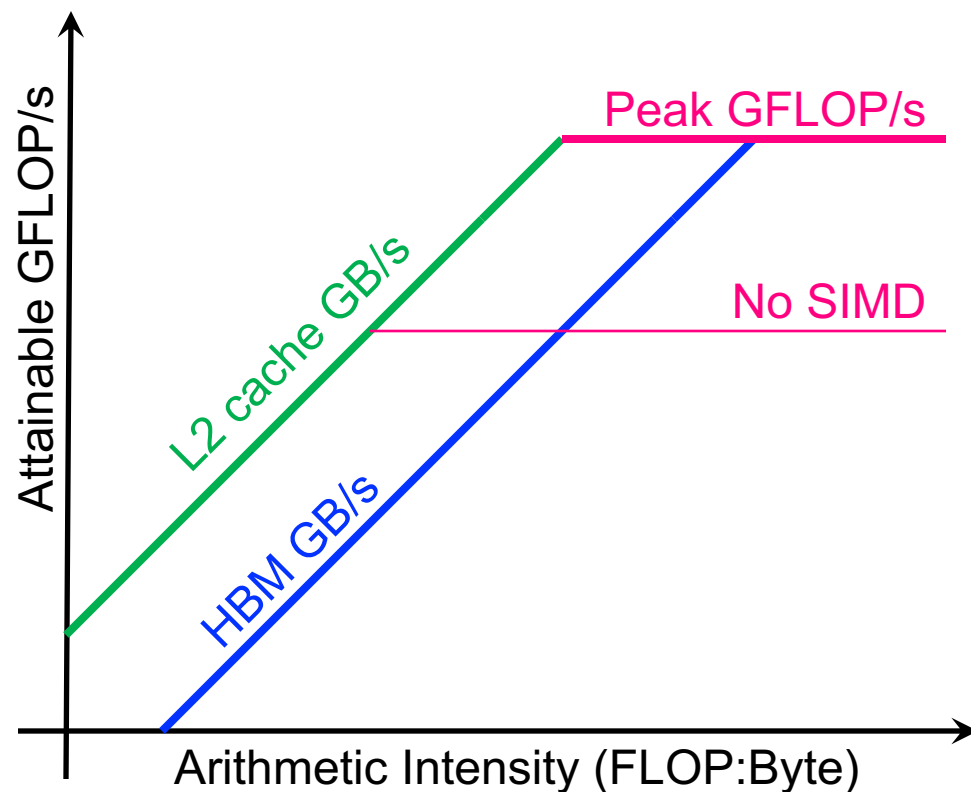  - Motivate need for algorithmic changes

BERKELEY LAB

# Model is just one piece of the puzzle…

- Roofline Model defines the basic concepts and equations.
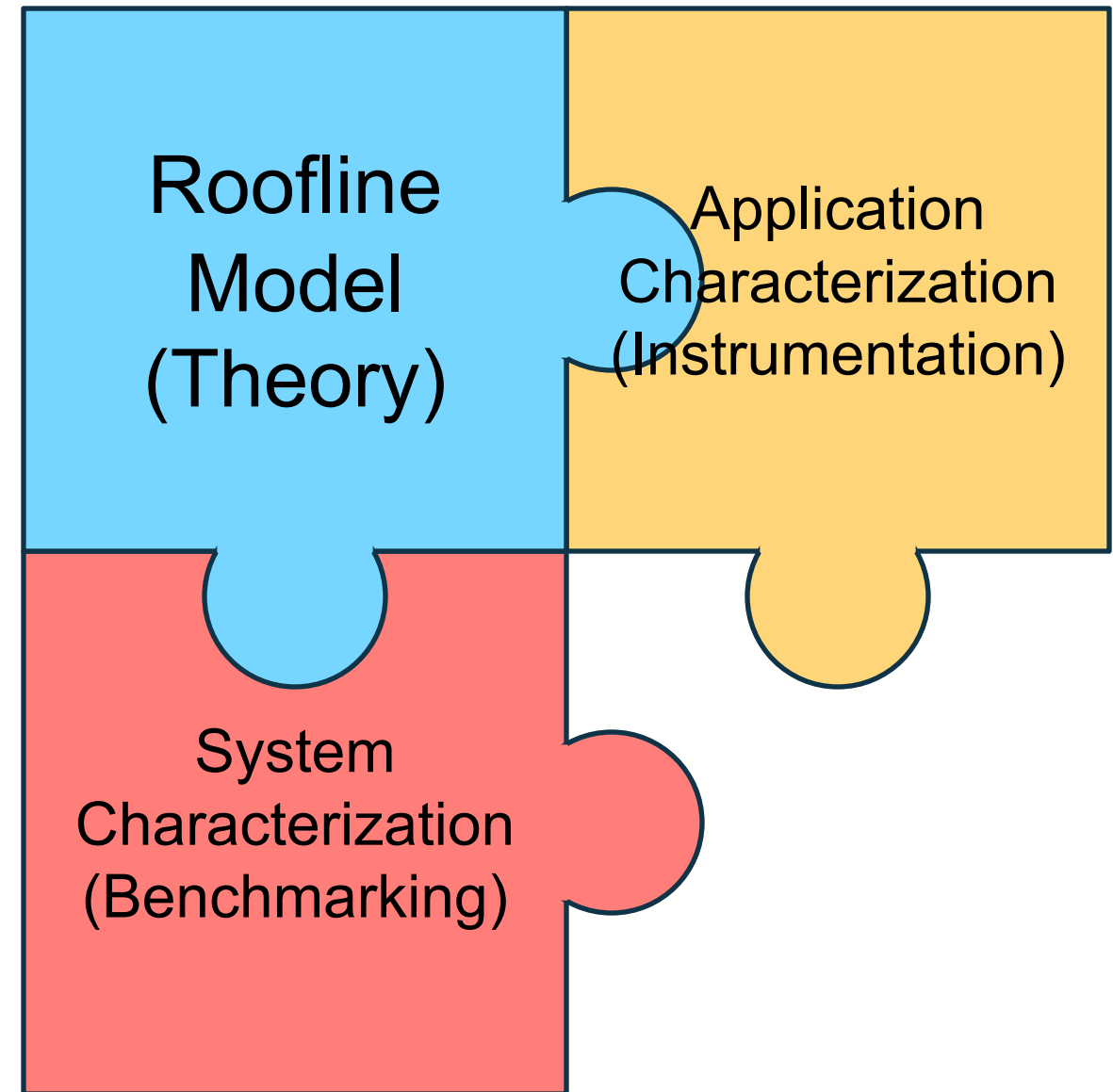
Roofline
Model
(Theory)

BERKELEY LAB

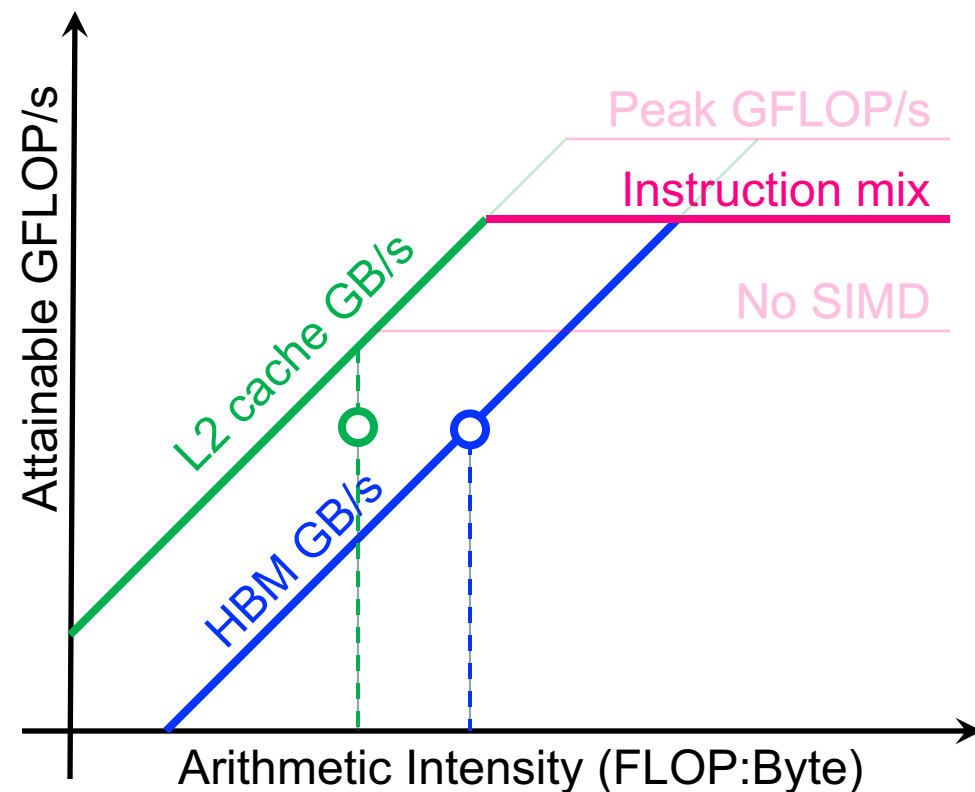# Model is just one piece of the puzzle...

- System Characterization defines the shape of the Roofline (peak bandwidths and FLOP/s)
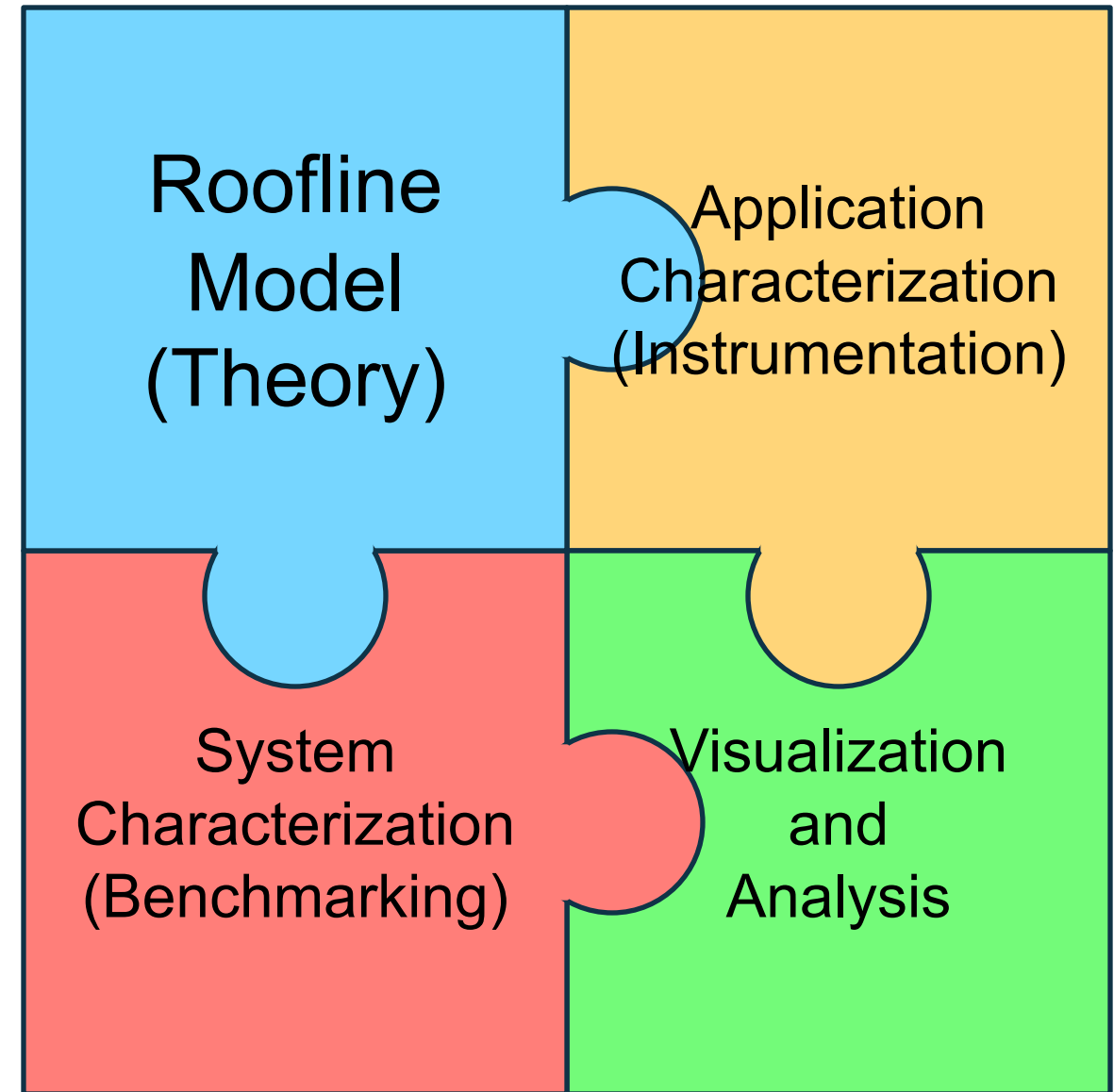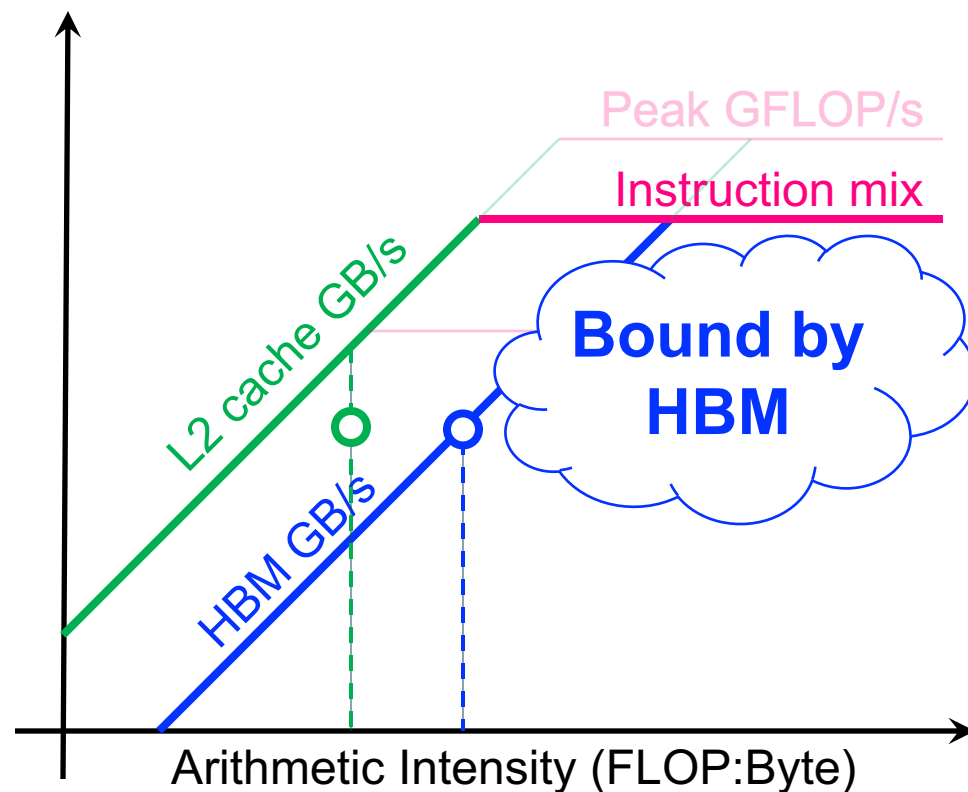
BERKELEY LAB

# Model is just one piece of the puzzle…

- Application Characterization determines…
  - Intensity and Performance of each loop
  - Position of any implicit ceilings

BERKELEY LAB

# Model is just one piece of the puzzle…

- Visualization tools combine all data together and provide analytical capability

BERKELEY LAB

# Rest of Tutorial

- Tools for Roofline analysis on NVIDIA GPUs and use cases at NERSC

- Tools for Roofline analysis on Intel GPUs/CPUs and use cases at ALCF

- Advanced Roofline topics

- Using TiMemory for portable application profiling

BERKELEY LAB

# Questions?

BERKELEY LAB
LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY