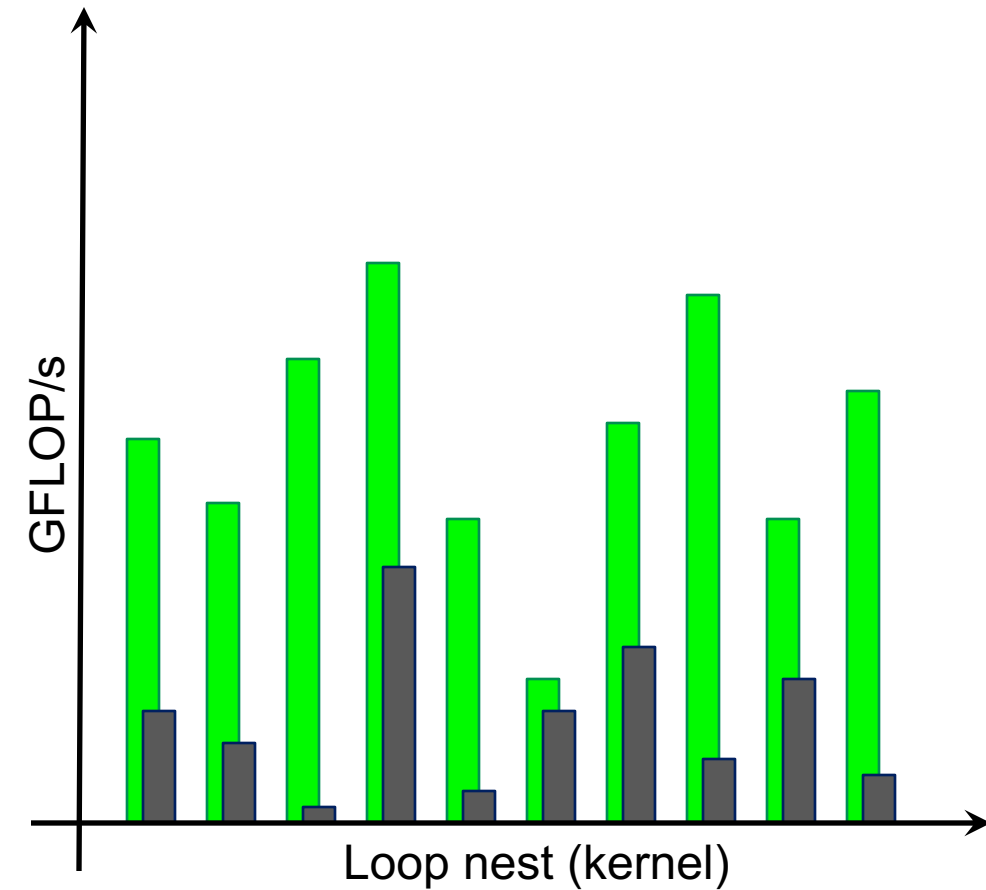You just spent 6 months porting your application to GPUs

Are you done?

# What is "Good" Performance?

- Imagine profiling the mix of loop nests in an application when running on the GPU

  - GFLOP/s alone may not be particularly insightful
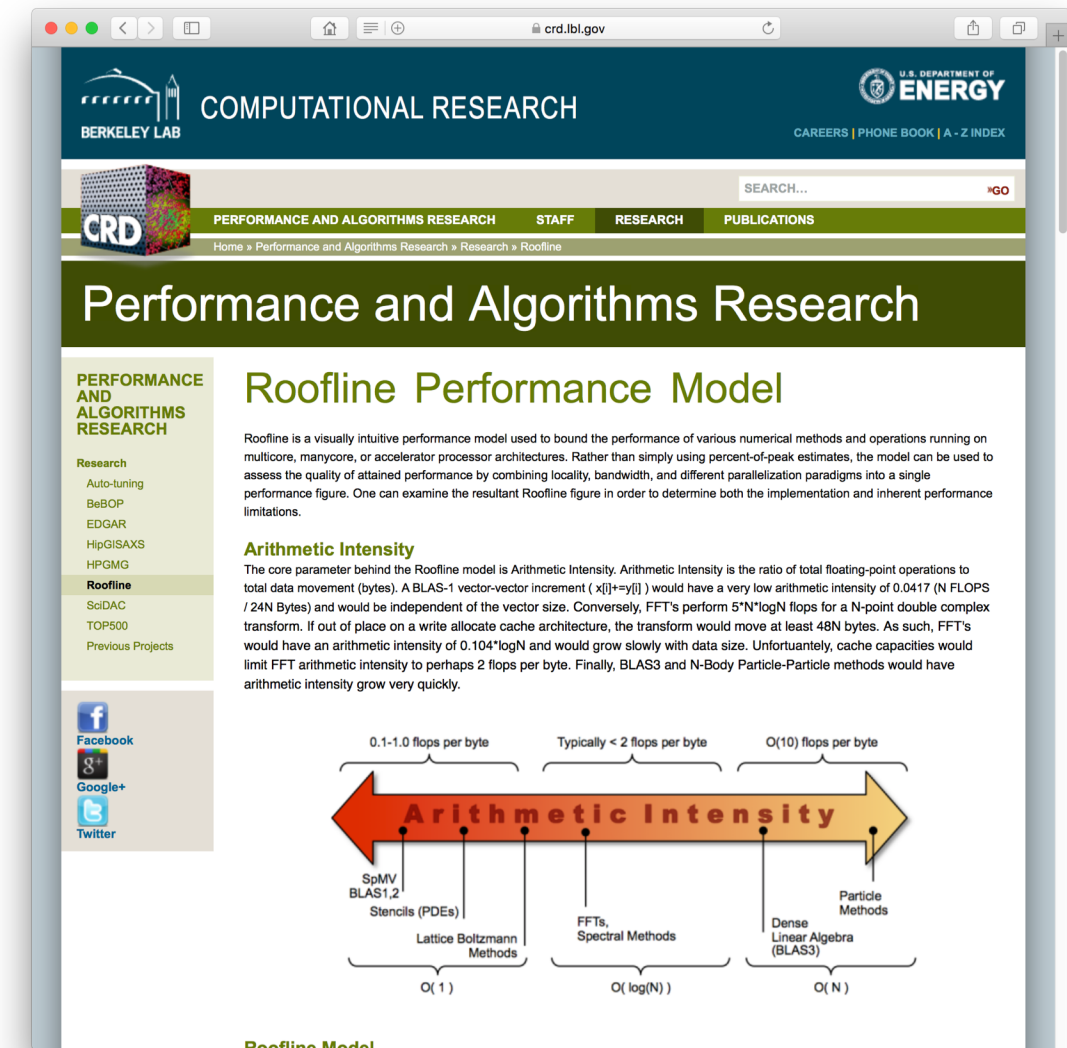
  - speedup relative to a Xeon may seem random

BERKELEY LAB

# What is "Good" Performance?

- Two fundamental aspects to "Good" performance…

1. Operating in the throughput-limited regime
   *not sensitive to Amdahl effects, D2H/H2D transfers, launch overheads, etc…*

2. making good use of the GPU's **compute** and/or **bandwidth** capabilities

➤ **Ultimately, we need a quantitative model rather than qualitative statements like "good"**

# Roofline Model

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

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

- Independent of ISA and architecture

- applies to CPUs, GPUs, Google TPUs[1], FPGAs, etc…

- Helps quantify **Good Performance**



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.

BERKELEY LAB

# Reduced Model

- Superscalar architectures can be complex
- Don't model / simulate full architecture
- Created simplified processor architecture

# Reduced Model

- Superscalar architectures can be complex

- Don't model / simulate full architecture

- Created simplified processor architecture

- Make assumptions on performance and usage…

  o Cores can attain peak GFLOP/s on local data

  o Cores execute load-balanced SPMD code

  o NoC bisection bandwidth is sufficient

  o There is sufficient cache bandwidth and capacity such that they do not affect performance

  ➢ **Basis for DRAM Roofline Model**

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?**



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

BERKELEY LAB

# Data Movement or Compute?

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

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

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

**Compute** — GFLOP/s

**Perfect Caches**

DRAM GB/s

**DRAM**

BERKELEY LAB

# Data Movement or Compute?

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

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



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

AI (Arithmetic Intensity) = FLOPs / Bytes (as presented to DRAM )

BERKELEY LAB

# Arithmetic Intensity

- Measure of data locality (data reuse)

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

- For the DRAM Roofline…
  - Total Bytes to/from DRAM
  - Includes all cache and prefetcher effects
  - Can be very different from total loads/stores (bytes requested)
  - Equal to ratio of sustained GFLOP/s to sustained GB/s (time cancels)

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 Roofline bound using Arithmetic Intensity as the x-axis

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

BERKELEY LAB

# (DRAM) Roofline Model

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

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

- 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…



*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 )

- Roofline tessellates this 2D view of performance into 5 regions…



Attainable FLOP/s

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

Peak GFLOP/s

*Compute bound*

*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 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

$GFLOP/s \leq AI * HBM\ GB/s$

TRIAD

0.083     5.0

Arithmetic Intensity (FLOP:Byte)

# 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  ];

}}}
```

BERKELEY LAB

# 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**
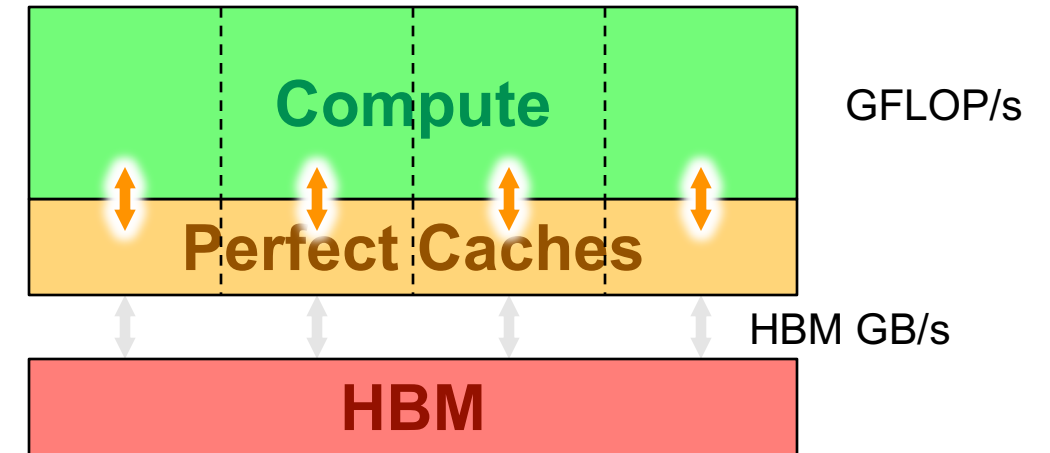  - **(measured at the L1)**



```
#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  ]

}}}
```

# 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**
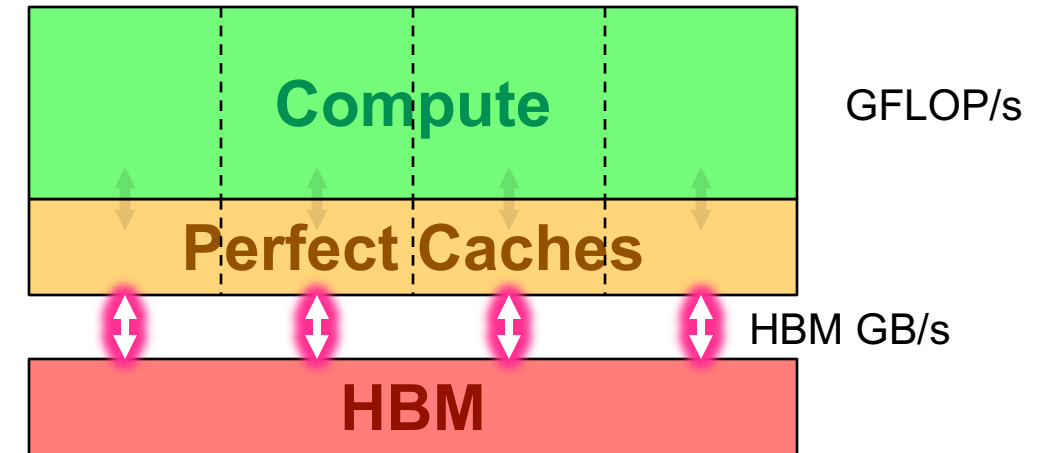


```
#pragma omp parallel for
for(k=1;k<dim+1;k++){
for(j=1;j<dim+1;j++){
for(i=1;i<dim+1;i++){
  new[k][j][i] = -6.0*old[k  ][j  ][i  ]
                     + old[k  ][j  ][i-1]
                     + old[k  ][j  ][i+1]
                     + old[k  ][j-1][i  ]
                     + old[k  ][j+1][i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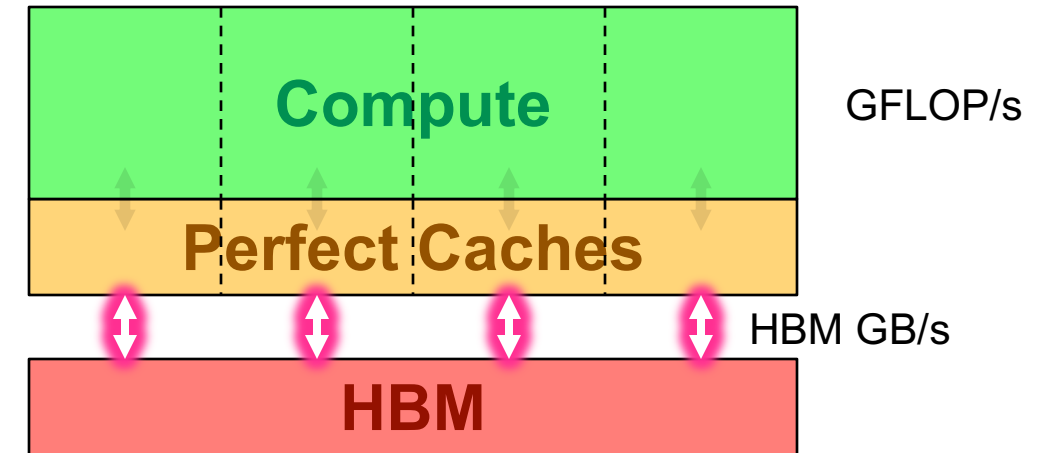


```
#pragma omp parallel for
for(k=1;k<dim+1;k++){
for(j=1;j<dim+1;j++){
for(i=1;i<dim+1;i++){
   new[k][j][i] = -6.0*old[k  ][j  ][i  ]
                      + old[k  ][j  ][i-1]
                      + old[k  ][j  ][i+1]
                      + old[k  ][j-1][i  ]
                      + old[k  ][j+1][i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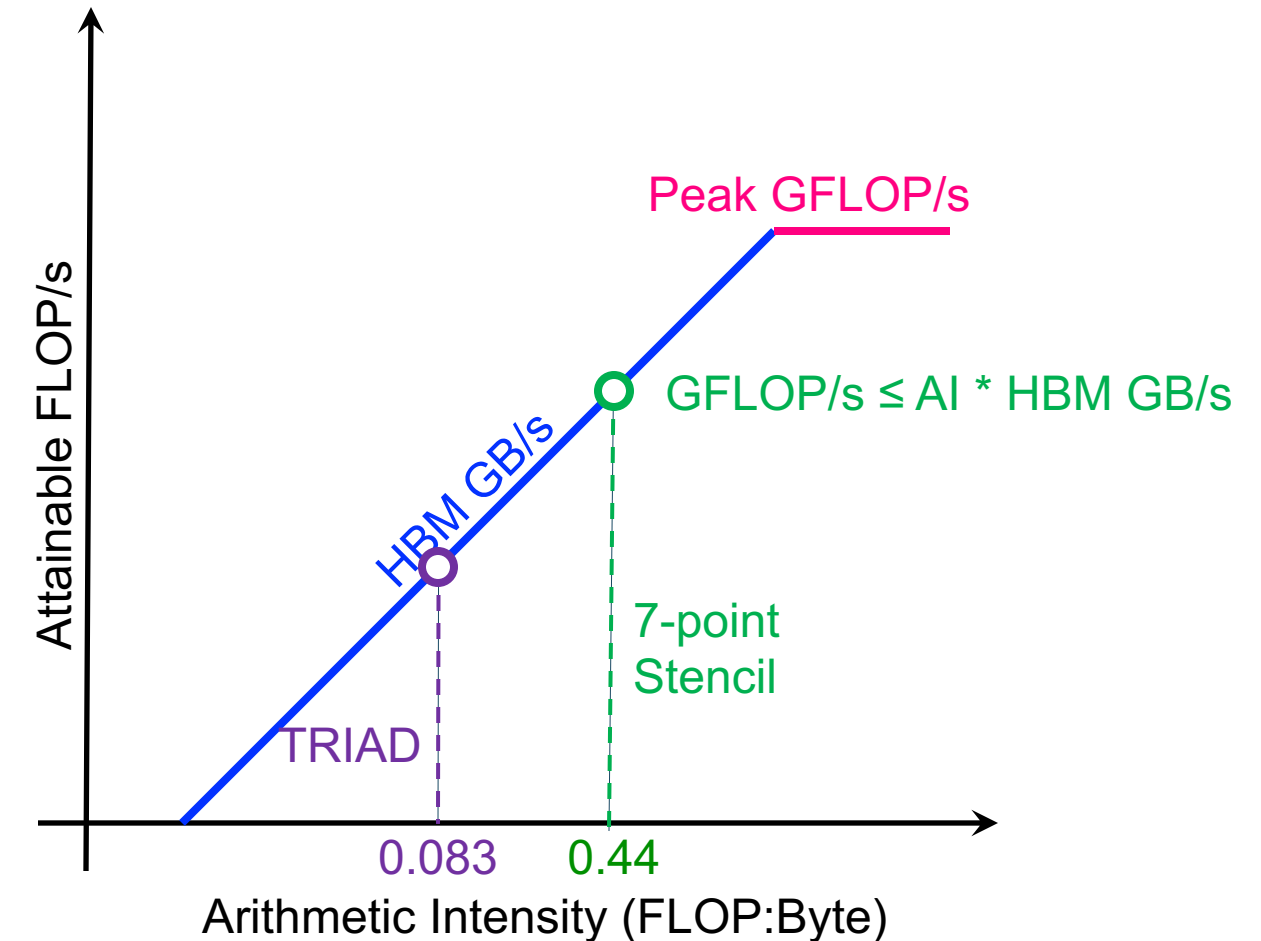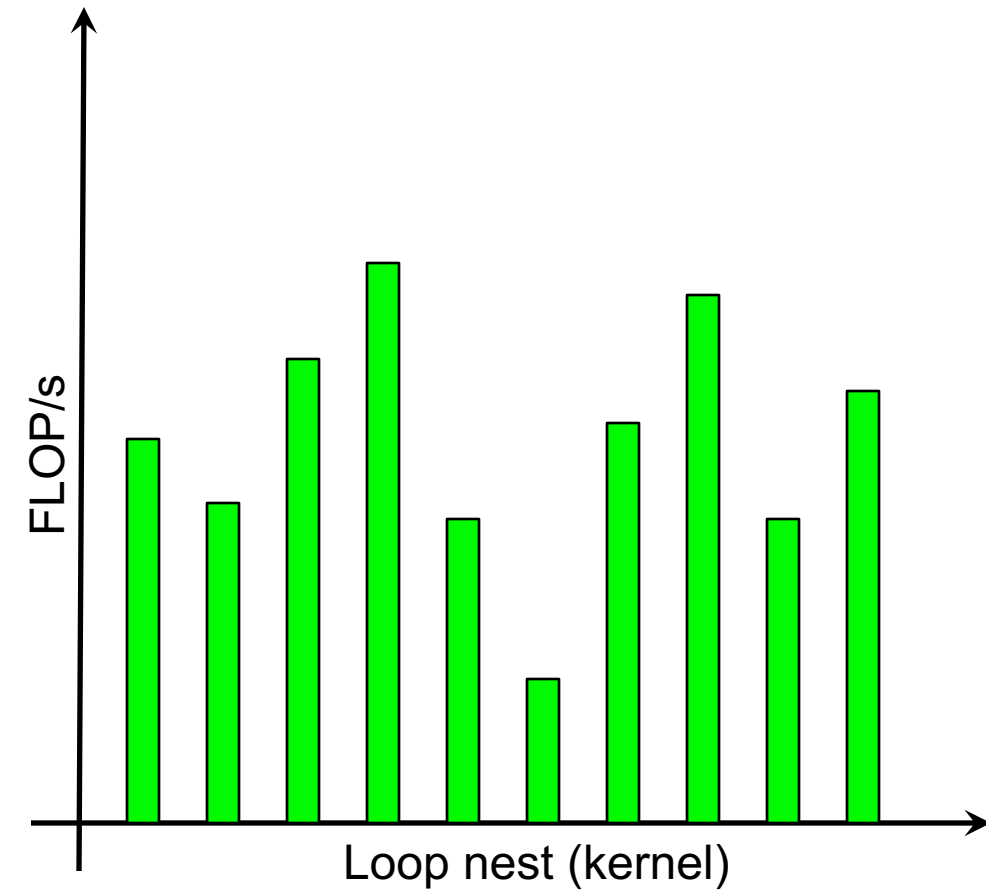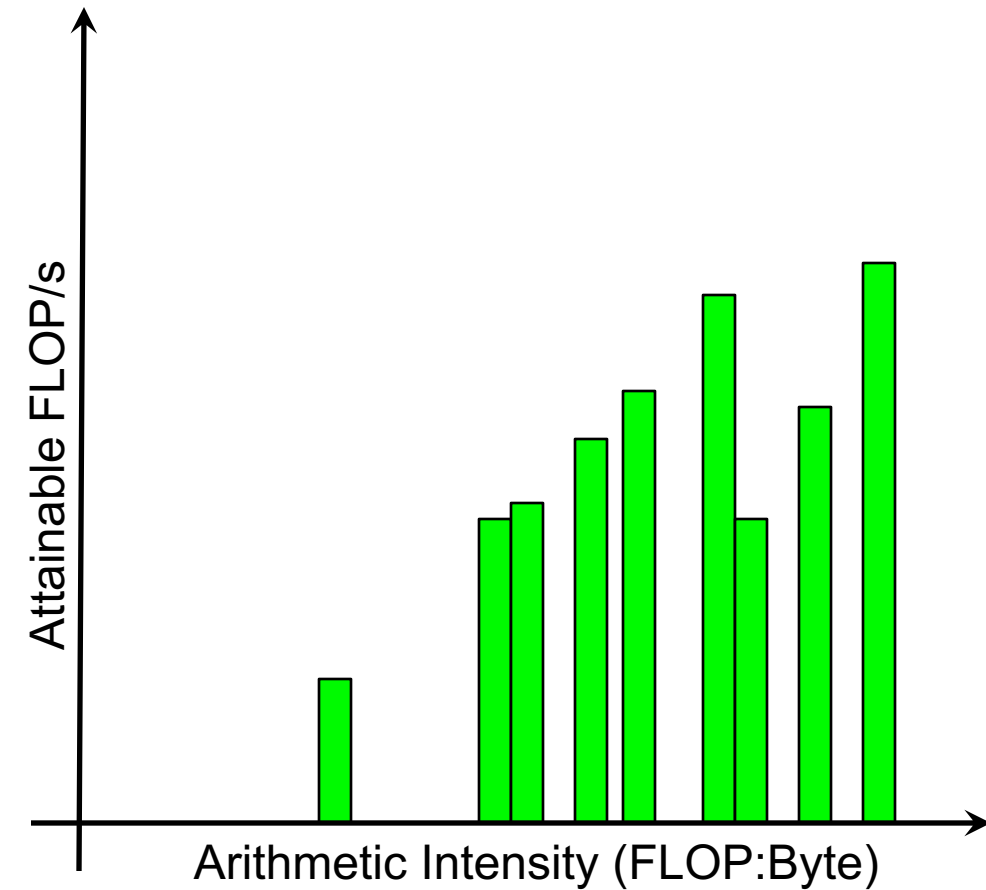
Attainable FLOP/s

Peak GFLOP/s

HBM GB/s

GFLOP/s ≤ AI * HBM GB/s

7-point Stencil

TRIAD

0.083    0.44

Arithmetic Intensity (FLOP:Byte)

BERKELEY LAB

# What is "Good" Performance?
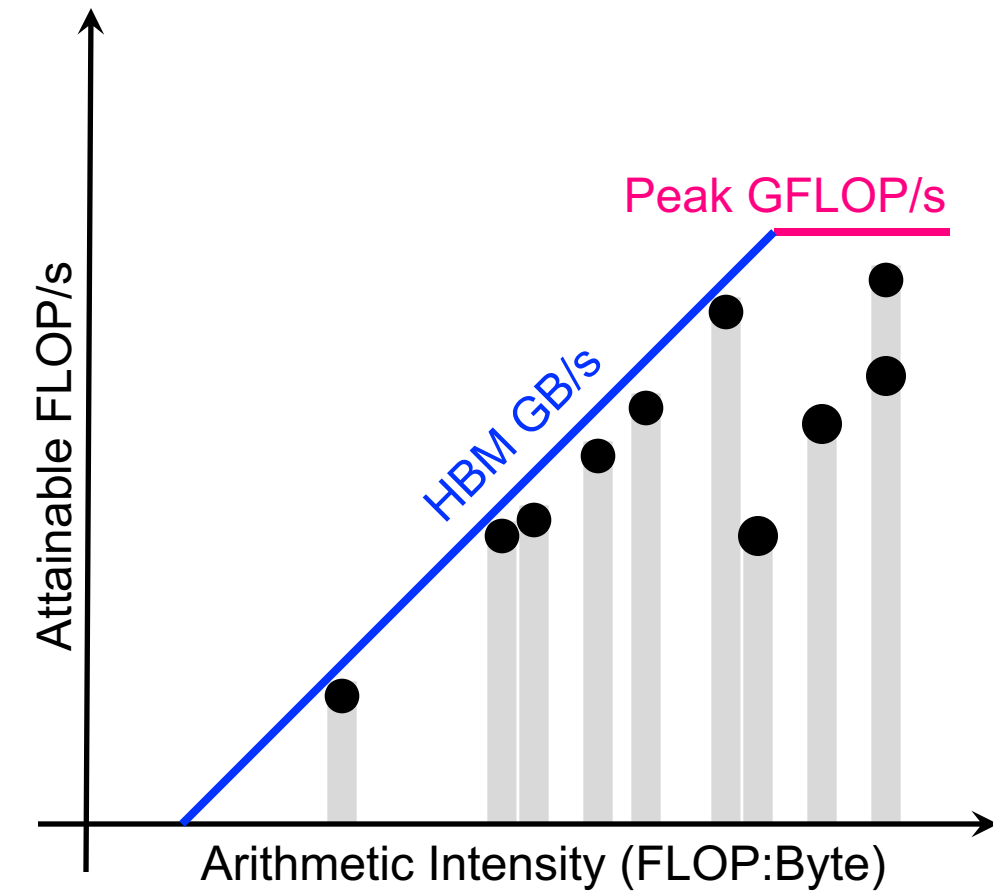
- Think back to our mix of loop nests…

# What is "Good" Performance?
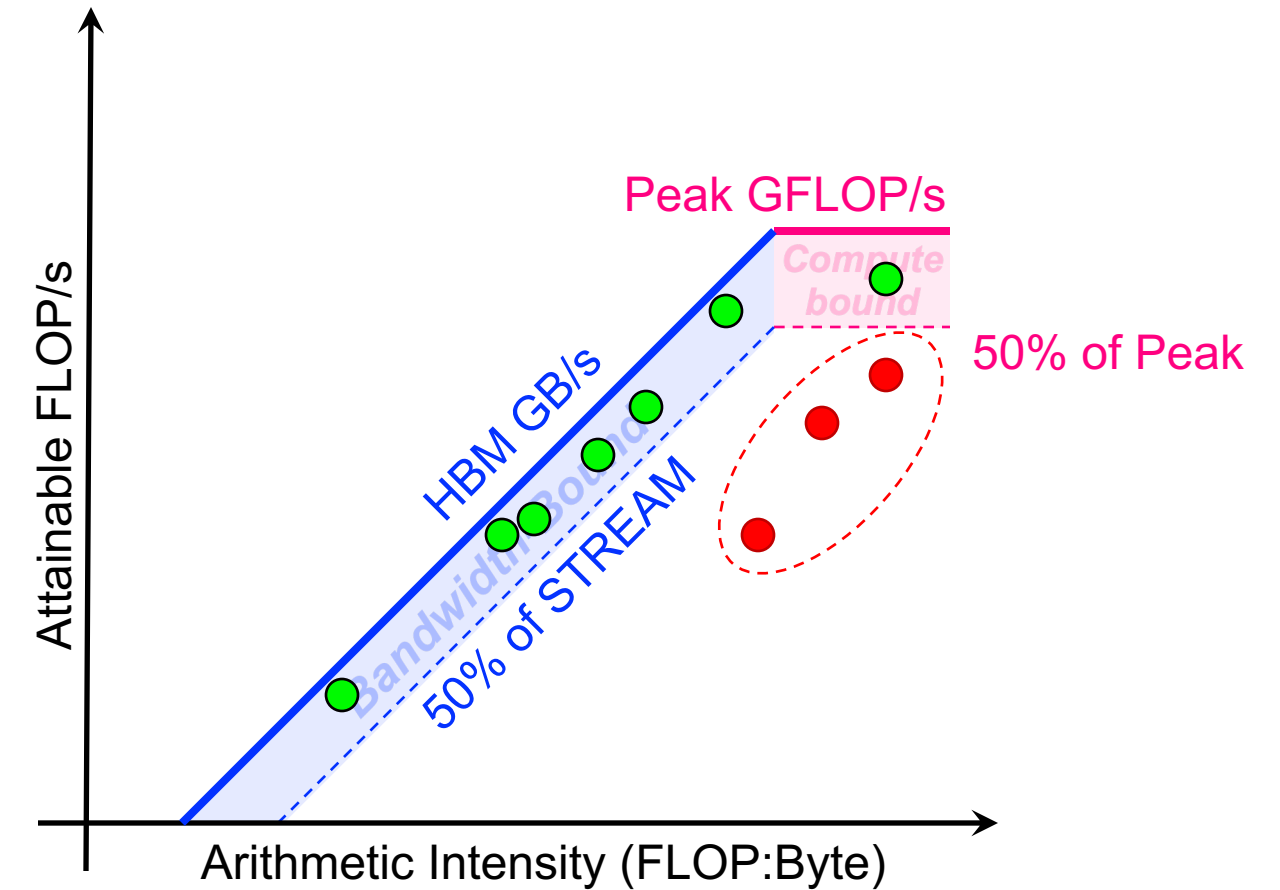
- We can sort kernels by arithmetic intensity…

# What is "Good" Performance?

- We can sort kernels by arithmetic intensity…

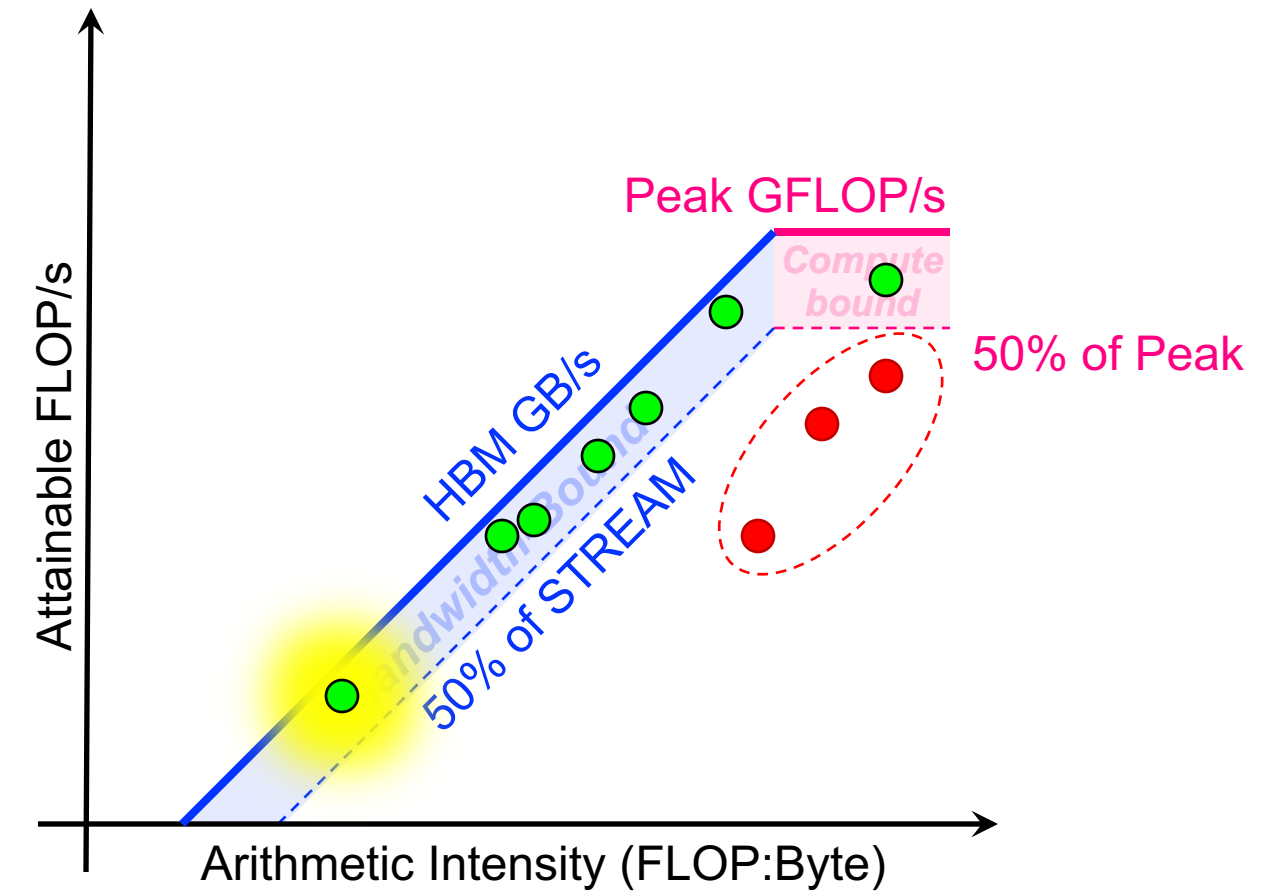- … and compare performance relative to machine capabilities

BERKELEY LAB

# What is "Good" Performance?

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

BERKELEY LAB

# What is "Good" Performance?

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

  ➢ kernels can have **low performance** (GFLOP/s), but make **good use** (%STREAM) of a machine
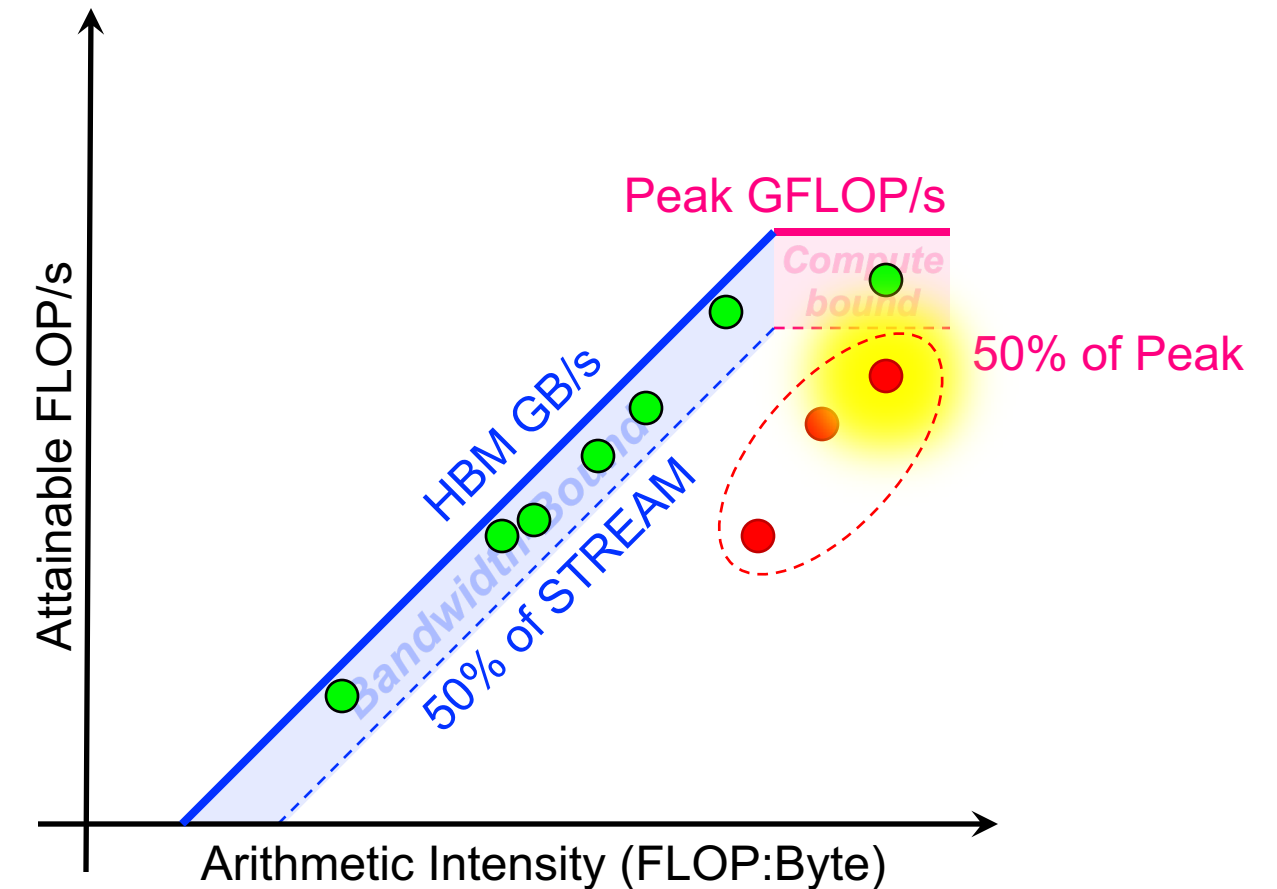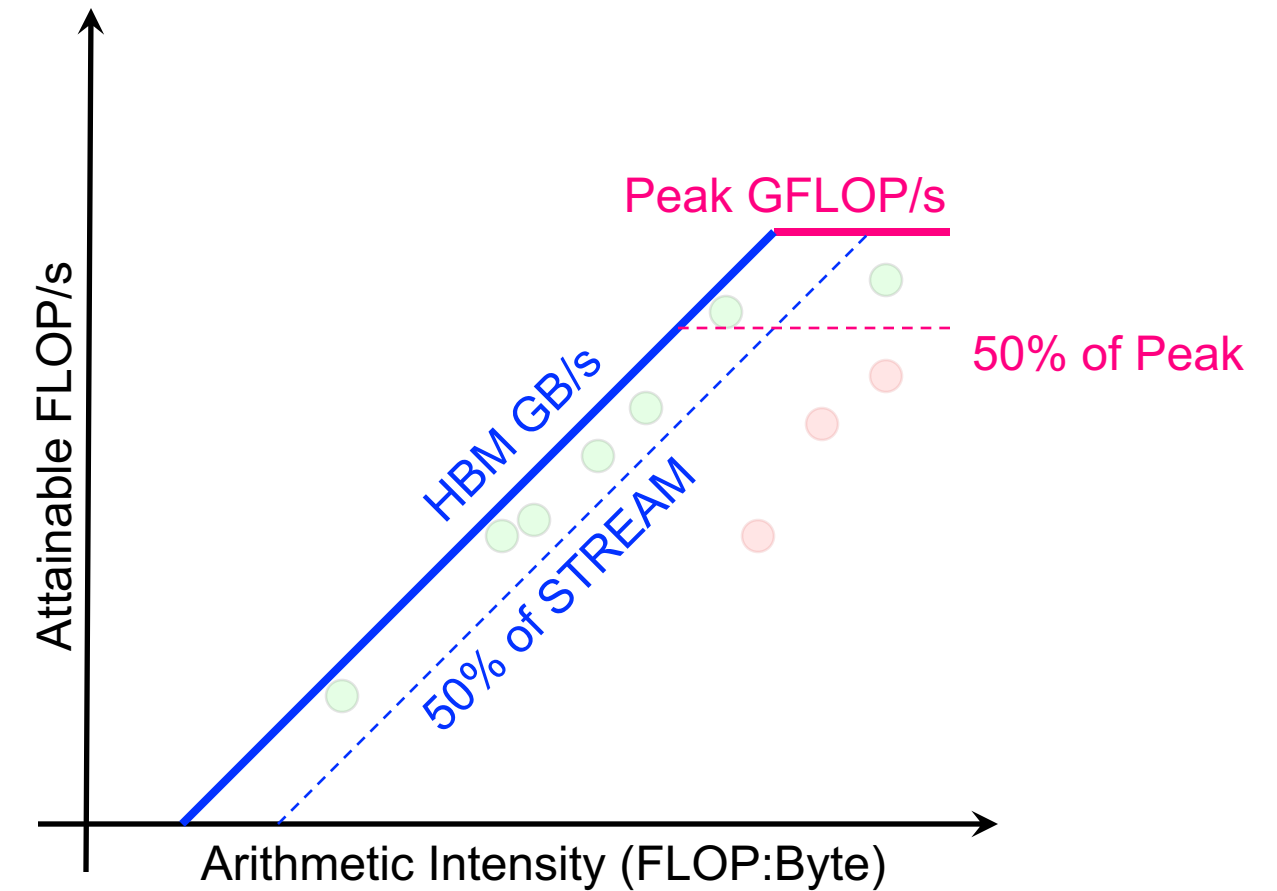
# What is "Good" Performance?

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

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

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

BERKELEY LAB

# 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



Peak GFLOP/s

50% of Peak

HBM GB/s

50% of STREAM

Attainable FLOP/s

Arithmetic Intensity (FLOP:Byte)

BERKELEY LAB

# Roofline is made of two components

- **Machine Model**
  - Lines defined by peak GB/s and GF/s (**Benchmarking**)
  - Unique to each architecture
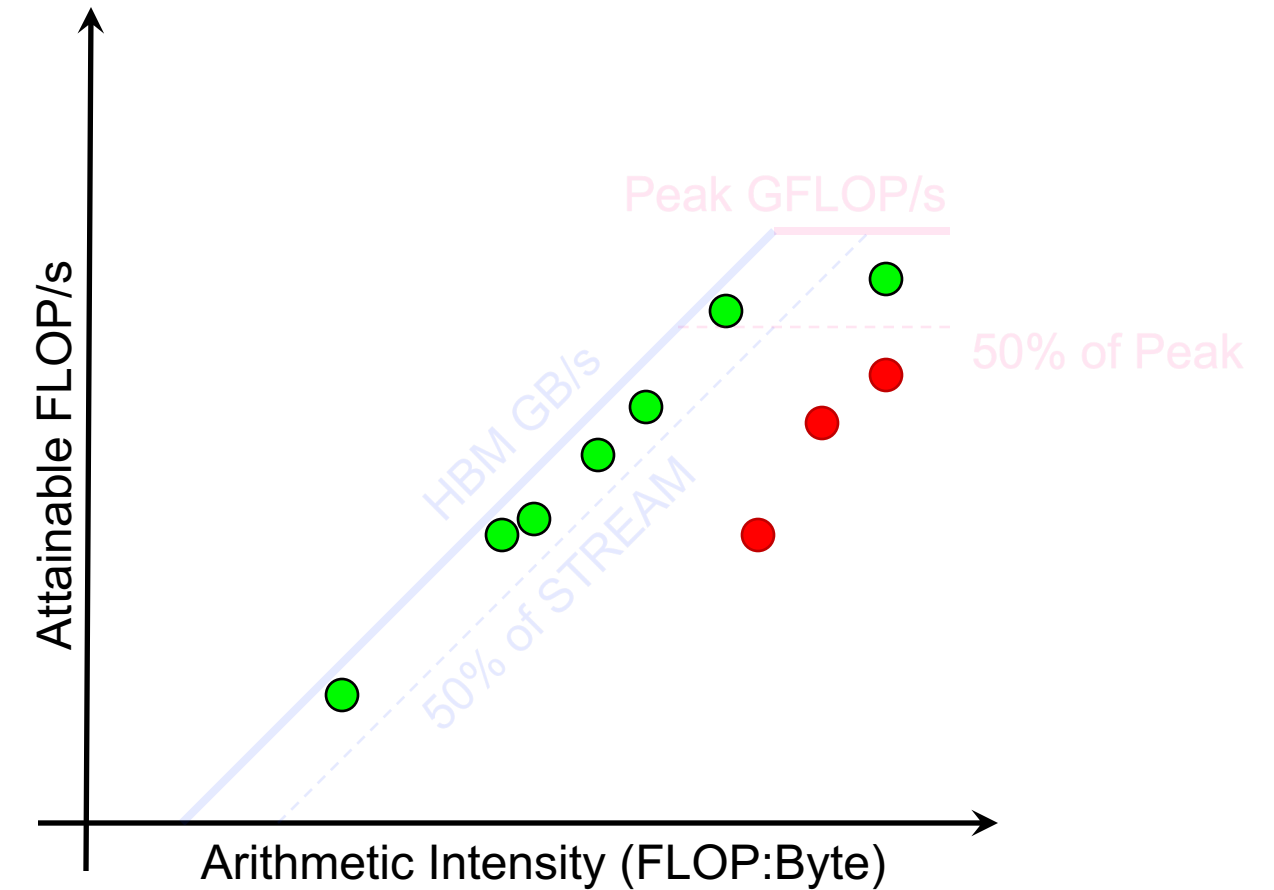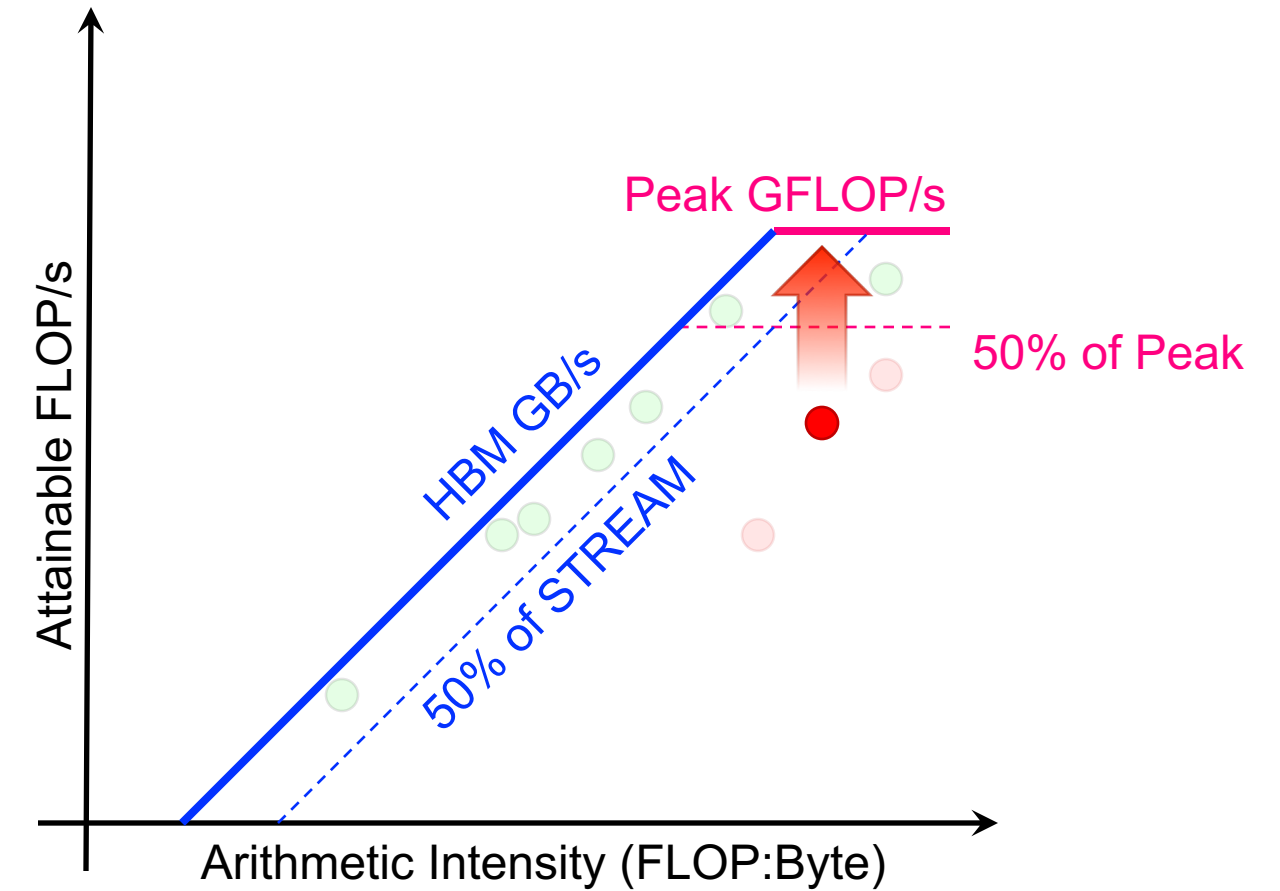  - Common to all apps on that architecture

- **Application Characteristics**
  - Dots defined by application GFLOP's and GB's (**Application Instrumentation**)
  - Unique to each application
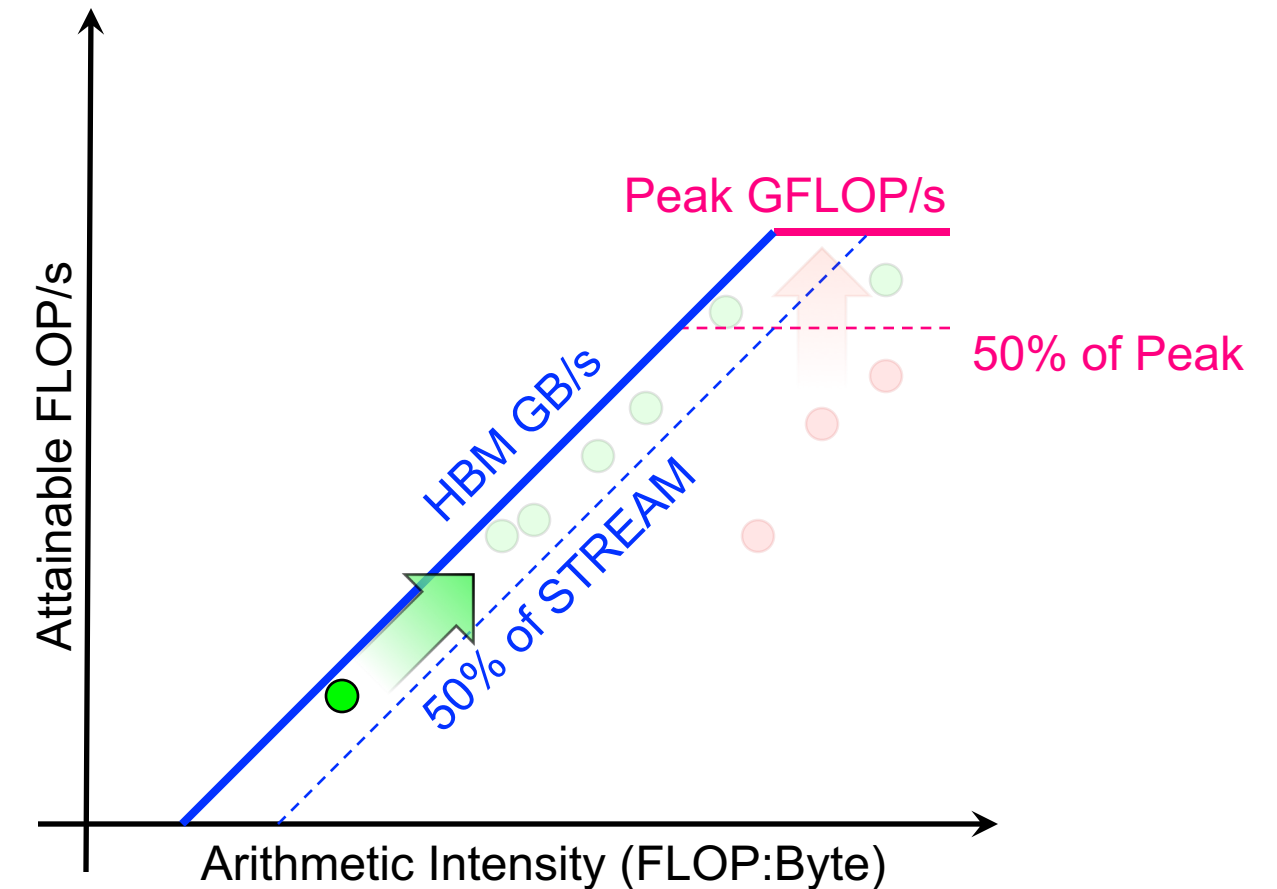  - Unique to each architecture

BERKELEY LAB

# General Performance Optimization Strategy

- Get to the Roofline

BERKELEY LAB

# General Performance Optimization Strategy

- **Get to the Roofline**

- **Increase Arithmetic Intensity when bandwidth-limited**

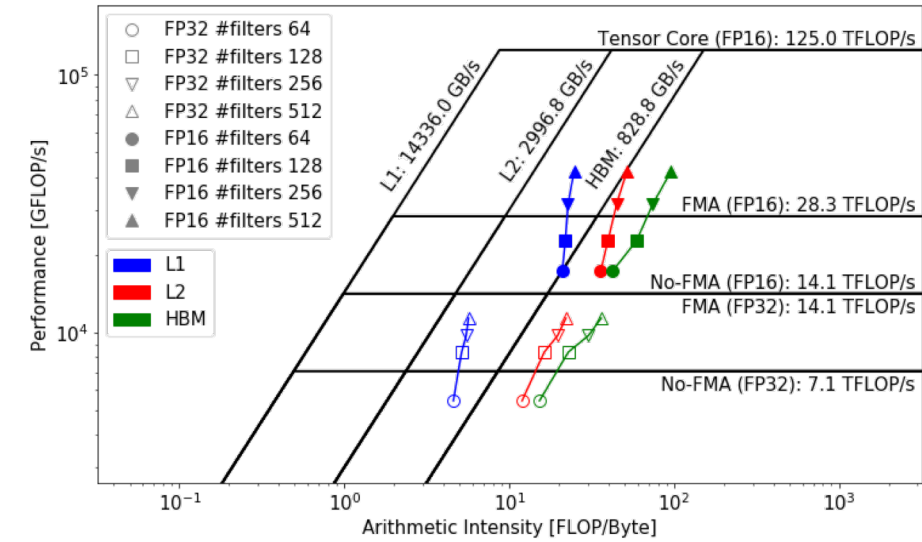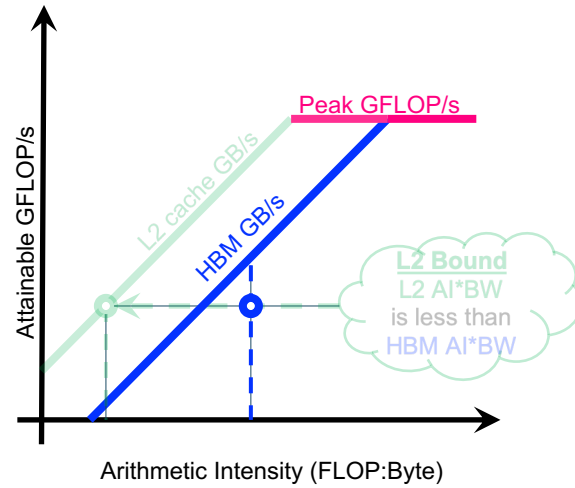  o Reducing data movement increases AI

BERKELEY LAB

# How can performance ever be below the Roofline?
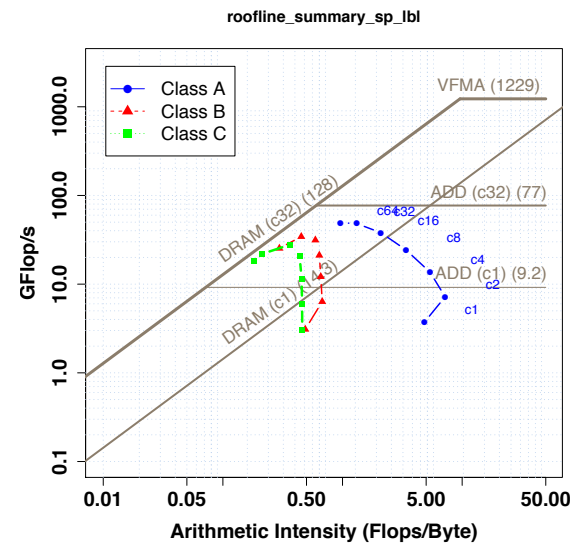
# Performance Below the Roofline?

## Hierarchical Roofline Model



Charlene Yang, Thorsten Kurth, Samuel Williams, "Hierarchical Roofline analysis for GPUs: Accelerating performance optimization for the NERSC-9 Perlmutter system", Concurrency and Computation: Practice and Experience (CCPE), August 2019.

## Additional FP Ceilings



Charlene Yang, Thorsten Kurth, Samuel Williams, "Hierarchical Roofline analysis for GPUs: Accelerating performance optimization for the NERSC-9 Perlmutter system", Concurrency and Computation: Practice and Experience (CCPE), August 2019.

## Roofline Scaling Trajectories



Khaled Ibrahim, Samuel Williams, Leonid Oliker, "Performance Analysis of GPU Programming Models using the Roofline Scaling Trajectories", International Symposium on Benchmarking, Measuring and Optimizing (Bench), BEST PAPER AWARD, November 2019.

## Instruction Roofline Model



Nan Ding, Samuel Williams, "An Instruction Roofline Model for GPUs", Performance Modeling, Benchmarking, and Simulation (PMBS), BEST PAPER AWARD, November 2019.

BERKELEY LAB

Summary

# Why We Use Roofline…

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

2. Identify performance bottlenecks & motivate software optimizations

3. 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?

4. Predict performance on future machines / architectures
   - Set realistic performance expectations
   - Drive for Architecture-Computer Science-Applied Math Co-Design

BERKELEY LAB

# Take away

- Roofline helps understand application performance relative to machine capabilities
    - just the beginning of the optimization process
    - Other bottleneck- or architecture-specific tools can be used to refine the process

- Roofline helps frame the conversation between…
    - Application Developers
    - Computer Scientists
    - Applied Mathematicians
    - Processor Vendors

…providing a common mental model and optimization language

BERKELEY LAB

Questions

BERKELEY LAB
LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY