# Using the Roofline Model and Intel Advisor

## Samuel Williams

SWWilliams@lbl.gov

**Computational Research Division**
**Lawrence Berkeley National Lab**

## Tuomas Koskela

TKoskela@lbl.gov

**NERSC**
**Lawrence Berkeley National Lab**

**BERKELEY LAB**
LAWRENCE BERKELEY NATIONAL LABORATORY

**U.S. DEPARTMENT OF ENERGY**

# Acknowledgements

Introduction
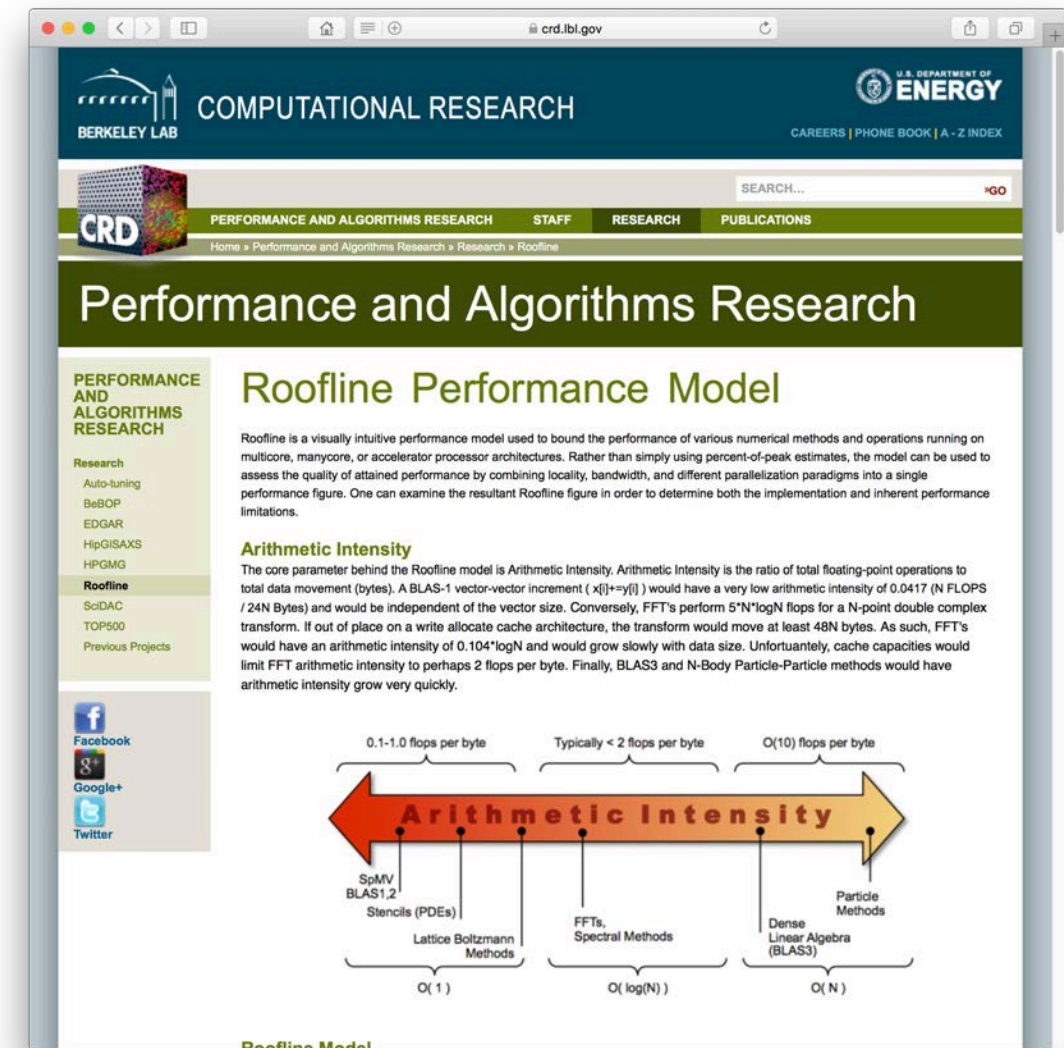
# Performance Models and Tools

- Identify performance bottlenecks
- Motivate software optimizations
- **Determine when we're done optimizing**
  - Assess performance relative to machine capabilities
  - Motivate need for algorithmic changes
- Predict performance on future machines / architectures
  - Sets realistic expectations on performance for future procurements
  - Used for HW/SW Co-Design to ensure future architectures are well-suited for the computational needs of today's applications.

BERKELEY LAB

# Performance Models / 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 time

  - Augmented with Little's Law

    (concurrency = latency*bandwidth)

  - Independent of ISA and architecture

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

- Three main components:

  - Machine Characterization (realistic performance potential of the system)

  - Monitoring (characterize application's execution)

  - Application Models (how well could my kernel perform with perfect compilers, procs, …)



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

- Ideally, we could always attain peak Flop/s

- However, finite locality (reuse) limits performance.

- Plot the performance bound using Arithmetic Intensity (AI) as the x-axis…
  - **Perf Bound = min ( peak Flop/s, peak GB/s * AI )**
  - AI = Flops / Bytes presented to DRAM
  - Log-log makes it easy to doodle, extrapolate performance, etc…
  - Kernels with AI less than machine balance are ultimately memory bound.

# Roofline Examples

- 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 = 0.166 flops per byte == Memory bound**



Roofline plot: Attainable Flop/s (y-axis) vs Arithmetic Intensity (Flop:Byte) (x-axis). Horizontal line labeled "Peak Flop/s". Diagonal line labeled "DRAM GB/s". Point labeled "TRIAD".

# Roofline Examples

- 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 = 0.43 flops per byte == memory bound,**

    **but 3x 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)

- We may measure a bandwidth and define an AI for each level
  - A given application / kernel / loop nest will thus have multiple AI's
  - A kernel could be DDR-limited…

# 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)

- **We may measure a bandwidth and define an AI for each level**
  - A given application / kernel / loop nest will thus have multiple AI's
  - A kernel could be DDR-limited…
  - **or MCDRAM-limited depending on relative bandwidths and AI's**

# 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)

# Basic Roofline Modeling

## Machine Characterization

*Potential of my target system*

- How does my system respond to a lack of FMA, DLP, ILP, TLP?

- How does my system respond to reduced AI (i.e. memory/cache bandwidth)?

- How does my system respond to NUMA, strided, or random memory access patterns?

- ...

## Application Instrumentation

*Properties of my app's execution*

- What is my app's real AI?

- How does AI vary with memory level ?

- How well does my app vectorize?

- Does my app use FMA?

- ...

# How Fast is My Target System?

- Challenges:
  - Too many systems; new ones each year
  - Voluminous documentation on each
  - Real performance often less than **"Marketing Numbers"**
  - **Compilers can "give up" on big loops**

- 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/

# Application Instrumentation Can Be Hard...

- Flop counters **can be broken/missing in production HW** (Haswell)

- Counting Loads and Stores is a poor proxy for data movement as they **don't capture reuse**

- Counting L1 misses is a poor proxy for data movement as they **don't account for HW prefetching**.

- DRAM counters are accurate, but **are privileged and thus nominally inaccessible in user mode**

- OS/kernel changes must be approved by vendor (e.g. Cray) and the center (e.g. NERSC)

# Application Instrumentation

- ## NERSC/CRD (==NESAP/SUPER) collaboration…

  - Characterize applications running on NERSC production systems

  - Use **Intel SDE** (binary instrumentation) to create software Flop counters (could use Byfl as well)

  - Use **Intel VTune** performance tool (NERSC/Cray approved) to access uncore counters

  - **Produced accurate measurement of Flop's and DRAM data movement on HSW and KNL**

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

BERKELEY LAB

# Use by NESAP

- NESAP is the NERSC KNL application readiness project.

- NESAP used Roofline to drive optimization and analysis on KNL
  - Bound performance expectations (ERT)
  - Quantify DDR and MCDRAM data movement
  - Compare KNL data movement to Haswell (sea of private/coherent L2's vs. unified L3)
  - Understand importance of vectorization

  - Doerfer et al., **"Applying the Roofline Performance Model to the Intel Xeon Phi Knights Landing Processor",** *Intel Xeon Phi User Group Workshop (IXPUG),* June 2016.
  - Barnes et al. **"Evaluating and Optimizing the NERSC Workload on Knights Landing",** *Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS),* November 2016.

BERKELEY LAB

# Roofline for NESAP Codes

# **Need a integrated solution…**

- Having to compose VTune, SDE, and graphing tools worked correctly and benefitted NESAP, but …

- …placed a very high burden on users…
  - **forced to learn/run multiple tools**
  - **forced to instrument each routine in their application**
  - **forced to manually parse/compose/graph the output**

- …still lacked integration with compiler/debugger/disassembly


- CRD/NERSC wanted a more integrated solution…

# Break / Questions

# Roofline vs. "Cache-Aware" Roofline

# There are two Major Roofline Formulations:

- ## Original / DRAM / Hierarchical Roofline…

  - **Williams, et al, "Roofline: An Insightful Visual Performance Model for Multicore Architectures", CACM, 2009**
  - Defines multiple bandwidth ceilings and multiple AI's per kernel
  - Performance bound is the minimum of the intercepts and flops

- ## "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 original Roofline, some use cache-aware **== Users need to understand the differences**
  - **Intel Advisor uses the Cache-Aware Roofline Model** (alpha/experimental DRAM Roofline being evaluated)
  - CRD/NERSC prefer the hierarchical Roofline as it provides greater insights into the behavior of the memory hierarchy

# 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 *directly observed*
- Requires *performance counters* to measure AI

# "Cache-Aware" Roofline

- Captures cache effects
- AI is Flop:Bytes *as presented to the L1 cache*
- Single Arithmetic Intensity
- AI *independent* of problem size
- Memory/Cache/Locality effects are *indirectly observed*
- 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

## Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

Actual Performance is the minimum of the two intercepts

Multiple AI's….
1) based on flop:DRAM bytes
2) Based on flop:L1 bytes (same)

## "Cache-Aware" Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

Observed performance is correlated with DRAM bandwidth

Single AI based on flop:L1 bytes

# 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)

## Roofline



- Peak Flop/s
- L1 GB/s
- DRAM GB/s
- Actual Performance is the minimum of the two
- Multiple AI's….
  1) flop:DRAM ~ 0.44
  2) flop:L1 ~ 0.11
- Attainable Flop/s
- Arithmetic Intensity (Flop:Byte)

## "Cache-Aware" Roofline



- Peak 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
- Attainable Flop/s
- Arithmetic Intensity (Flop:Byte)

# Example: 7-point Stencil (Large Problem)

## Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

Capacity misses reduce DRAM AI and performance

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

## "Cache-Aware" Roofline



Peak Flop/s

L1 GB/s

DRAM GB/s

Attainable Flop/s

Arithmetic Intensity (Flop:Byte)

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

Single AI based on flop:L1 bytes

BERKELEY LAB

Break / Questions

# Intel Advisor:
## Introduction and General Usage

*DRAM Roofline and OS/X Advisor GUI:* These are preview features that may or may not be included in mainline product releases. They may not be stable as they are prototypes incorporating very new functionality. Intel provides preview features in order to collect user feedback and plans further development and productization steps based on the feedback.

# Intel Advisor

- **Integrated Performance Analysis Tool**
  - Performance information including timings, flops, and trip counts
  - Vectorization Tips
  - Memory footprint analysis
  - **Uses the Cache-Aware Roofline Model**
  - All connected back to source code

- **CRD/NERSC began a collaboration with Intel**
  - Ensure Advisor runs on Cori in user-mode
  - Push for **Hierarchical Roofline**
  - Make it functional/scalable to many MPI processes across multiple nodes
  - Validate results on NESAP, SciDAC, and ECP codes

NESAP is NERSC's KNL application readiness project
SciDAC is the DOE Office of Science's Scientific Discovery thru Advanced Computing program
ECP is the DOE's Exascale Computing Project

32

BERKELEY LAB

# Intel Advisor (Useful Links)

## Background

- https://software.intel.com/en-us/intel-advisor-xe

- https://software.intel.com/en-us/articles/getting-started-with-intel-advisor-roofline-feature

- https://www.youtube.com/watch?v=h2QEM1HpFgg

## Running Advisor on NERSC Systems

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

# Using Intel Advisor at NERSC

- ## Compile…

  ```
  use '-g' when compiling
  ```

- ## Submit Job…

  ```
  % salloc –perf=vtune            <<< interactive sessions; --perf only needed for DRAM Roofline
           –or–
  #SBATCH –perf=vtune             <<< batch submissions; --perf only needed for DRAM Roofline
  ```

  ## Benchmark…

  ```
  % module load advisor
  % export ADVIXE_EXPERIMENTAL=roofline_ex        <<< only needed for DRAM Roofline
  % srun [args] advixe-cl -collect survey -no-stack-stitching -project-dir $DIR -- ./a.out [args]
  % srun [args] advixe-cl -collect tripcounts -flops-and-masks -callstack-flops -project-dir $DIR -- ./a.out [args]
  ```

- ## Use Advisor GUI…

  ```
  % module load advisor
  % export ADVIXE_EXPERIMENTAL=roofline_ex        <<< only needed for DRAM Roofline
  % advixe-gui $DIR
  ```

BERKELEY LAB

NoMachine - NERSC

/global/cscratch1/sd/tkoskela/dram_roofline/stencil/advi.stencil.aug2.16 - Intel Advisor <@cori05> <2>

File  View  Help

Welcome | e000 (read-only) ✕

Elapsed time: 50.50s  | Vectorized | Not Vectorized |  OFF  Smart Mode

INTEL ADVISOR 2017

FILTER: | All Modules ▼ | All Sources ▼ | Loops And Functions ▼ | All Threads ▼ |

🌳 Summary  🌿 Survey & Roofline  🍎 Refinement Reports

| Function Call Sites and Loops | 🔥 | 💡 Vector Issues | Self Time | Total Time | Type | FLOPS | |
|---|---|---|---|---|---|---|---|
| | | | | | | GFLOPS ▼ | AI |
| ⊡🔄 [loop in bench_stencil_ver4 ... | ☐ | | 159.595s■ | 159.595s■ | Vectorized (Body) | 23.083 | 0.117 |
| ⊞🔄 [loop in bench_stencil_ver3 ... | ☐ | 💡 1 Ineffective peeled/... | 159.953s■ | 159.953s■ | Vectorized (Body; Re... | 16.274 | 0.117 |
| ⊞🔄 [loop in bench_stencil_ver2 ... | ☐ | 💡 1 Ineffective peeled/... | 160.035s■ | 160.035s■ | Vectorized (Body; Peel... | 15.662 | 0.117 |
| ⊞🔄 [loop in bench_stencil_ver1 ... | ☐ | 💡 1 Ineffective peeled/... | 159.307s■ | 159.307s■ | Vectorized (Body; Peel... | 10.218 | 0.117 |
| ⊡🔄 **[loop in bench_stencil_v ...** | ☐ | **💡 1 Potential under ...** | **157.994s**☐ | **157.994s**☐ | **Scalar** | **9.009** | **0.117** |

| Source | Top Down | Code Analytics | Assembly | 💡 Recommendations | 🔲 Why No Vectorization? |

| Function Call Sites and Loops | Total Time % | Total Time | Self Time | Type | FLOPS | |
|---|---|---|---|---|---|---|
| | | | | | GFLOPS | |
| ⊟_INTERNAL_26_____src_z_Linux_util_cpp_2d702c13::[OpenM | 93.8%■ | 755.843s | 0.000s[ | Function | 0.998 | |
| ⊟🔄 [loop in _INTERNAL_26_____src_z... | 93.8% | 755.843s | 0.000s[ | Scalar | 0.998 | |
| ⊟__kmp_launch_thread | 93.8% | 755.843s | 0.000s[ | Function | 0.998 | |
| ⊟🔄 [loop in __kmp_launch_thread at kmp_runtime.cpp:565 | 93.8%■ | 755.843s | 0.000s[ | Scalar | 0.998 | |
| ⊟[OpenMP dispatcher] | 93.3%■ | 752.000s | 0.000s[ | Function | 1.003 | |
| ⊞bench_stencil_ver2$omp$parallel_for@102 | 18.7%■ | 150.903s | 0.120s[ | Function | 1.083 | |
| ⊞bench_stencil_ver3$omp$parallel_for@146 | 18.7%■ | 150.471s | 0.000s[ | Function | 1.097 | |
| ⊞bench_stencil_ver4$omp$parallel_for@193 | 18.6%■ | 149.989s | 0.000s[ | Function | 1.540 | |
| ⊞bench_stencil_ver1$omp$parallel_for@62 | 18.6%■ | 149.743s | 0.000s[ | Function | 0.696 | |

_INTERNAL_26_____src_z_Linux_util_cpp_2d702c13::[OpenMP worker]

🅰 Advixe-gui   3  | 🌐 emacs-gtk@cori05-bond0.224   | 🖥 cori :

BERKELEY LAB

# Break / Questions

# Intel Advisor:
## Stencil Roofline Demo*

*DRAM Roofline and OS/X Advisor GUI:* These are preview features that may or may not be included in mainline product releases. They may not be stable as they are prototypes incorporating very new functionality. Intel provides preview features in order to collect user feedback and plans further development and productization steps based on the feedback.

# 7-point, Constant-Coefficient Stencil

- Apply to a $512^3$ domain on a single NUMA node (single HSW socket)
- Create 5 code variants to highlight effects (as seen in advisor)

ver0.      Baseline: thread over outer loop (k), but prevent vectorization

```
#pragma novector                                    // prevent simd
int ijk = i*iStride + j*jStride + k*kStride; // variable iStride to confuse the compiler
```

ver1.      Enable vectorization

```
int ijk = i + j*jStride + k*kStride;           // unit-stride inner loop
```

ver2.      Eliminate capacity misses

*2D tiling of j-k iteration space*                     *// working set had been O(6MB) per thread*

ver3.      Improve vectorization

*Provide aligned pointers and strides*

ver4.      Force vectorization / cache bypass

```
__assume(jstride%8 == 0);                       // stride by variable is still aligned
#pragma omp simd, vector nontemportal          // force simd; force cache bypass
```

Cache-Aware Roofline

Cache-Aware Roofline

# Cache-Aware Roofline

Cache-Aware Roofline

DRAM Roofline*

DRAM Roofline*

DRAM Roofline*

# DRAM Roofline*

DRAM Roofline*

Wrap up / Questions

# Roofline/Advisor Tutorial at SC'17

- Sunday November 12th
- 8:30am-12pm (half-day tutorial)
- multi-/manycore focus

BERKELEY LAB

# Intel Advisor (Useful Links)

## Background

- https://software.intel.com/en-us/intel-advisor-xe

- https://software.intel.com/en-us/articles/getting-started-with-intel-advisor-roofline-feature

- https://www.youtube.com/watch?v=h2QEM1HpFgg

## Running Advisor on NERSC Systems

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

# Acknowledgements