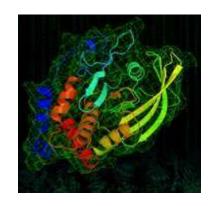
### Introduction to the Roofline Model

















Based on the slides of
Charlene Yang
Lawrence Berkeley National Laboratory
Jun 16 2019, Frankfurt





### **Performance Models**









The Maze of Performance Optimization

The Map!!!

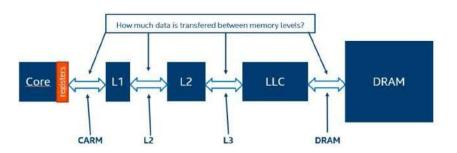




### **Performance Models**

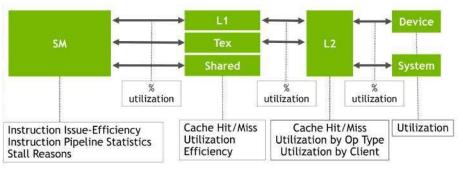


#### Modern architectures are complicated!



Intel Haswell CPU<sup>1</sup>

#### **NVIDIA Volta GPU<sup>2</sup>**





- 1. https://software.intel.com/en-us/articles/integrated-roofline-model-with-intel-advisor
- 2. http://on-demand.gputechconf.com/gtc/2016/presentation/s6659-avinash-baliga-perfworks.pdf



#### **Performance Models**



- § Many components contribute to the kernel (核心算子,例如矩阵乘法)
  run time
- § An interplay of application characteristics and machine characteristics

#FP operations FLOP/s
Cache data movement Cache GB/s
DRAM data movement DRAM GB/s
PCle data movement PCle bandwidth
MPI Message Size Network Bandwidth
MPI Send:Wait ratio Network Gap
#MPI Wait's Network Latency
IO File systems

Roofline Model
Roofline Model
Focus on one or two



#### **Roofline Performance Model**



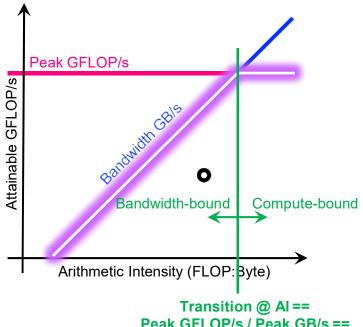
Sustainable performance is bound by

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

Arithmetic Intensity (AI) =

FLOPs /Bytes

How did this come about? **CPU DRAM example** 



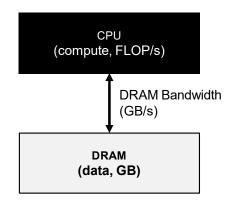








- § One could hope to always attain peak performance (FLOP/s)
- § However, finite locality (reuse) and bandwidth limit performance.
- § Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)



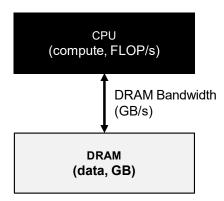






- § One could hope to always attain peak performance (FLOP/s)
- § However, finite locality (reuse) and bandwidth limit performance.
- § Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)

```
Time
#FP ops = max = 1 / Peak GFLOP/s
#Bytes / #FP ops / Peak GB/s
```

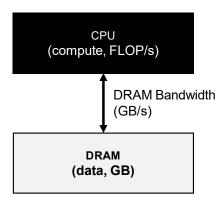








- § One could hope to always attain peak performance (FLOP/s)
- § However, finite locality (reuse) and bandwidth limit performance.
- § Assume:
  - Idealized processor/caches
  - Cold start (data in DRAM)









**DRAM Bandwidth** 

(GB/s)

CPU (compute, FLOP/s)

DRAM (data, GB)

- § One could hope to always attain peak performance (FLOP/s)
- § However, finite locality (reuse) and bandwidth limit performance.
- § Assume:
  - Cold start (data in DRAM)



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





#### **Roofline Performance Model**



§ Thus we obtain the model as

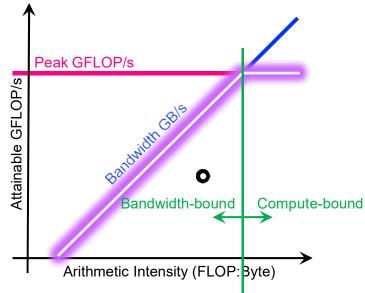
GFLOP/s = min 
$$\begin{cases} Peak GFLOP/s \\ AI * Peak GB/s \end{cases}$$

where Arithmetic Intensity (AI) is

FLOPs /Bytes

Machine Balance (FLOPs/Byte) =
 8.9 (V100, DP, HBM) or 5.1 (KNL, DP, HBM)

An application achieved FLOP/s



Transition @ Al ==
Peak GFLOP/s / Peak GB/s ==
'Machine Balance'





# More Advanced on Roofline















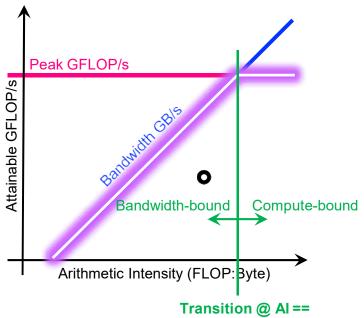


### **Roofline Performance Model**



- § This is a single Roofline
- § What about the memory hierarchy, different execution configurations, and instruction mixes?

Hierarchical Roofline
Multiple compute ceilings



Transition @ Al ==
Peak GFLOP/s / Peak GB/s ==
'Machine Balance'



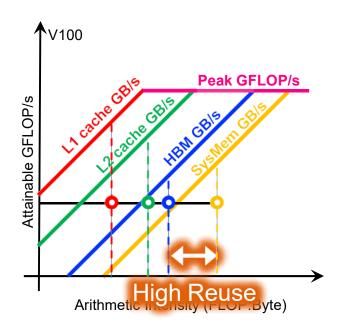


#### **Hierarchical Roofline**



- Superposition of multiple Rooflines
  - Incorporate full memory hierarchy
  - Arithmetic Intensity =FLOPs /Bytes<sub>L1/L2/HBM/SysMem</sub>

 Each kernel will have multiple Al's but one observed GFLOP/s performance



Hierarchical Roofline tells you about cache locality

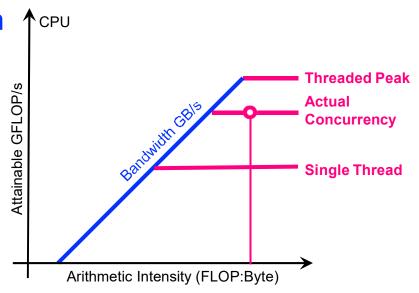




## **Multiple Compute Ceilings**



- Impact of execution configuration
- Concurrency affects your peak
  - OpenMP thread concurrency
  - SM occupancy
  - load balance
  - threadblock/thread configuration



Performance is bound by the actual concurrency ceiling

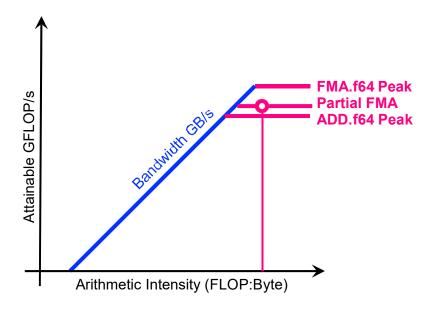




## **Multiple Compute Ceilings**



- Impact of instruction mix
- Applications are usually a mix of FMA.f64, ADD.f64, MUL.f64...
- Performance is a weighted average
   ... bound by a partial FMA ceiling







# **Roofline Drives Optimization**

















### **Roofline Performance Model**



#### The Roofline Model

- helps you identify the bottlenecks
- guides you through optimization
- tells you when to stop

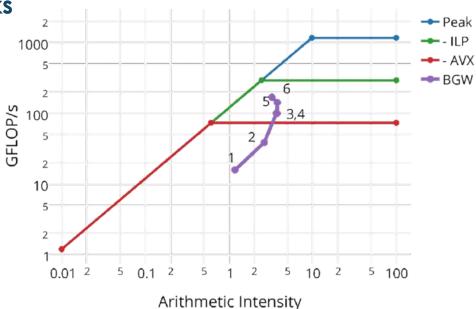
#### An example:

NESAP for Cori - BerkeleyGW

BerkeleyGW is a massively parallel computational package in Material Science that simulates electron excited state properties for a variety of material systems from bulk semiconductors and metals to nanostructured materials and molecules.

https://bitbucket.org/berkeleylab/berkeleygw-nesap/src/master/

#### Haswell Roofline Optimization Path





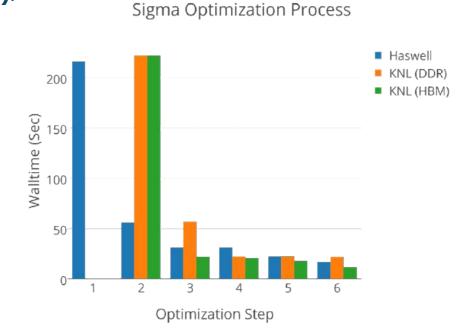


## Roofline Example: BerkeleyGW



#### Optimization Path for Kernel-C (Sigma):

- 1. Add OpenMP
- 2. Initial Vectorization
  - loop reordering
  - conditional removal
- 3. Cache-Blocking
- 4. Improved Vectorization
  - divides
- 5. Hyper-threading

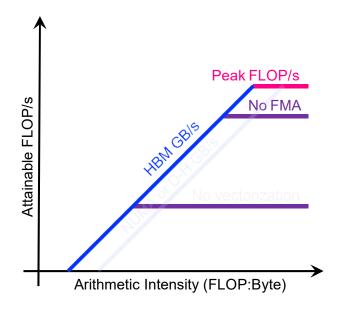








§ Broadly speaking, three approaches to improving performance:

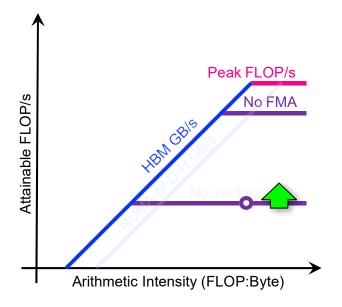








- § Broadly speaking, three approaches to improving performance:
- § Maximize compute performance
  - § multithreading
  - § vectorization
  - § increase SM occupancy
  - § utilize FMA instructions
  - § minimize thread divergence

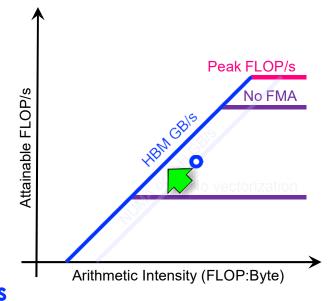








- § Broadly speaking, three approaches to improving performance:
- § Maximize compute performance
- § Maximize memory bandwidth
  - § utilize higher-level caches
  - § NUMA-aware allocation
  - § avoid H-D transfers
  - § avoid uncoalesced memory access

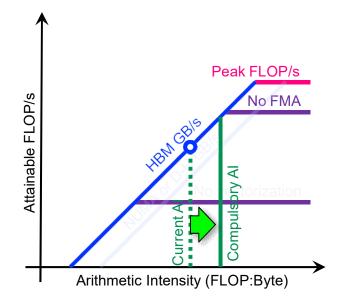








- § Broadly speaking, three approaches to improving performance:
- § Maximize compute performance
- § Maximize memory bandwidth
- § Improve Al
  - § minimize data movement
  - § exploit cache reuse







# **Roofline Data Collection**

















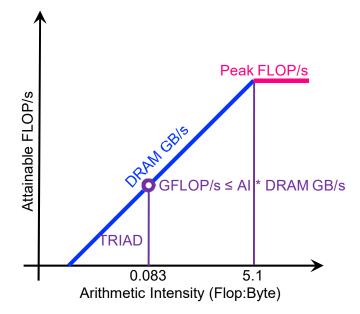
## Pen and Paper



§ Example #1: STREAM Triad

```
for(i=0;i<N;i++){
    Z[i] = X[i] + alpha*Y[i];
}</pre>
```

- 2 FLOPs per iteration
- Transfer 24 bytes per iteration
  - read X[i], Y[i], and write Z[i]
- AI = 0.083 FLOPs per byte
- Memory bound



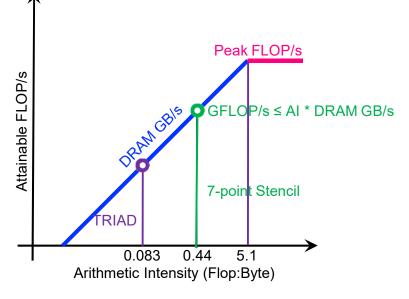




### Pen and Paper



- § Example #2: 7-pt stencil
  - 7 FLOPs; 8 memory references (7 reads, 1 store) per pt
  - Cache can filter all but 1 read and 1 write per pt
  - AI = 0.44 FLOPs per byte
  - Memory bound, but 5x the GFLOP/s rate







## Pen and Paper



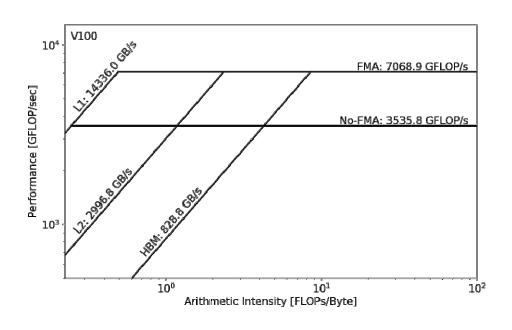
- Not scalable for real-life applications
- Millions of lines of code; mix of different languages
- Complicated modern architecture
  - memory hierarchy, caching effects
  - ISA
- Different problem sizes





### We Need Tools!





#### Roofline ceilings

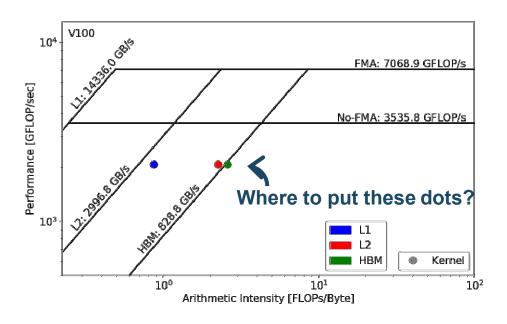
- vendor specifications
- empirical measurements
  - ERT
  - https://bitbucket.org/be rkeleylab/cs-rooflinetoolkit





### We Need Tools!



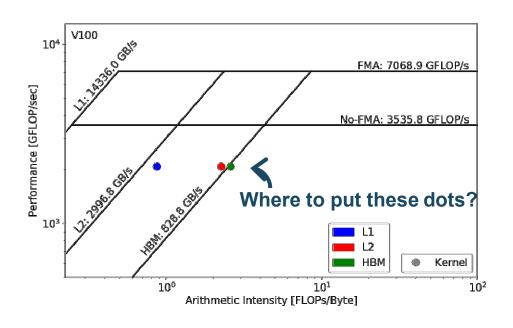






### We Need Tools!





#### Require three raw measurements:

- Runtime
- FLOPs
- Bytes (on each cache level)

In order to calculate AI and GFLOP/s:

Arithmetic Intensity = 
$$\frac{FLOPs}{Data Movement}$$

$$\frac{FLOPs}{Performance} = \frac{FLOPs}{Runtime}$$





## Methodology to Construct Roofline



- 1. Collect Roofline ceilings
  - **compute** (FMA/no FMA) and **bandwidth** (DRAM, L2, ...)
  - ERT: https://bitbucket.org/berkeleylab/cs-roofline-toolkit
- 2. Collect application performance
  - FLOPs, bytes (DRAM, L2, ...), runtime
  - SDE, VTune, LIKWID, Advisor, nvprof, ...
- 3. Plot Roofline with Python Matplotlib (or other tools of your preference)
  - arithmetic intensity, GFLOP/s performance, ceilings
  - example scripts: https://github.com/cyanguwa/nersc-roofline





# **Automated Data Collection**

















### **Data Collection on Intel CPUs**



#### The not-so-automated way 1:

- Intel SDE for FLOPs (emulation)
- Intel VTune for DRAM bytes (HW counters)
- Runtime
- DRAM Roofline only
- Used by NESAP for Cori
  - **NERSC Exascale Science Application Program**
  - http://www.nersc.gov/users/application-performance/measuring-arithmetic-intensity/



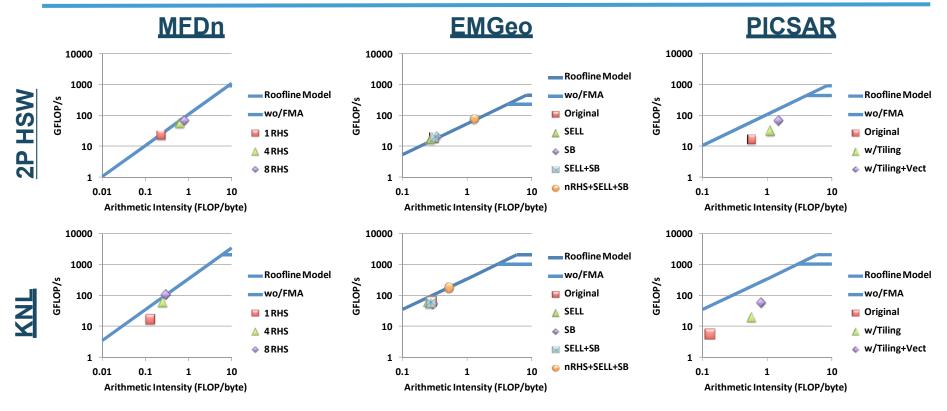






### **Data Collection on Intel CPUs**







**DRAM Rooflines of NESAP Codes** 

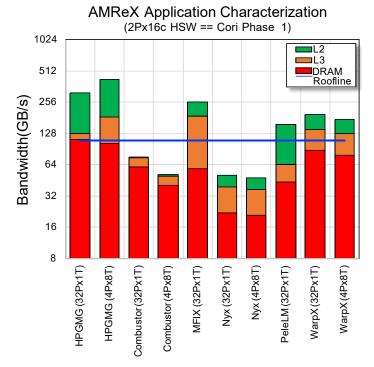


### **Data Collection on Intel CPUs**



#### The not-so-automated way 2:

- LIKWID for FLOPs and bytes
  - Both are based on HW counters
- Runtime
- Hierarchical Roofline
- Limited by quality of HW counters
- High-level characterization, no callstack



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









#### The fully automated way:

- Intel Advisor, Roofline feature
- Instrument applications automatically
  - one dot per loop nest/function
- FLOPs, bytes and runtime
- Hierarchical Roofline
- Integrates with other Advisor capabilities
- Benchmarks target system



N Q = 1 = B + | □ Use Single-Trimmled Rooks\*





### Data Collection on NVIDIA GPUs



- Still very manual at this stage, but...
- Runtime:
  - Internal timers or nvprof --print-gpu-trace
- FLOPs:
  - DP/SP/HP counters and metrics, nvprof --metrics
     'flop\_count\_dp/sp/hp' Or `tensor\_precision\_fu\_utilization'
- Bytes for different cache levels:
  - Bytes = (read transactions + write transactions) x transaction size
  - nvprof --metrics 'metric\_name' e.g. gld/gst\_transactions
- Hierarchical Roofline





## Summary



- The Roofline Model formulizes the interaction between machine characteristics and application characteristics, and guides optimization
  - Peak computational throughput and bandwidth
  - Arithmetic intensity, cache locality, instruction mix...
- Automate Roofline data collection
  - Intel CPUs
    - Intel SDE + Intel VTune, Intel Advisor
  - NVIDIA GPUs
    - nvprof, Nsight Compute



More in the next few talks!





#### Reference



- S. Williams, A. Waterman and D. Patterson, "Roofline: An Insightful Visual Performance Model for Multicore Architectures," Communications of the ACM, vol. 52, no. 4, pp. 65–76, 2009
- LBNL CRD Roofline Research:
   <a href="https://crd.lbl.gov/departments/computer-science/PAR/research/roofline">https://crd.lbl.gov/departments/computer-science/PAR/research/roofline</a>
- Empirical Roofline Toolkit (ERT):
   https://bitbucket.org/berkeleylab/cs-roofline-toolkit
- Python scripts for plotting manually-collected Roofline:
   <a href="https://github.com/cyanguwa/nersc-roofline/tree/master/Plotting">https://github.com/cyanguwa/nersc-roofline/tree/master/Plotting</a>



