Master Informatics Eng.

2016/17 A.J.Proença

The Roofline Performance Model

(most slides are borrowed)

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

XX



Goals of the Roofline Model

\sim

CONVENTIONAL WISDOM IN COMPUter architecture produced similar designs. Nearly every desktop and server computer uses caches, pipelining, superscalar instruction issue, and out-of-order execution. Although the instruction sets varied, the microprocessors were all from the same school of

Roofline Model

For the foreseeable future, <u>off-chip</u> memory bandwidth will often be the constraining resource in system performance.²³ Hence, we want a model that relates processor performance to off-chip memory traffic. Toward this

DOI:10.1145/1498765.1498785

3

The Roofline model offers insight on how to improve the performance of software and hardware.

BY SAMUEL WILLIAMS, ANDREW WATERMAN, AND DAVID PATTERSON



AJProença, Advanced Architectures, MiEl, UMinho, 2016/17

Performance Limiting Factors



AJProença, Advanced Architectures, MiEl, UMinho, 2016/17

5

Roofline Performance Model

- Basic idea:
 - Plot peak floating-point throughput as a function of arithmetic intensity
 - Ties together floating-point performance and memory performance for a target machine
- Arithmetic intensity
 - Floating-point operations per byte read





Copyright © 2012, Elsevier Inc. All rights reserved.

Components to performance:
Computation
Communication
Locality
Each architecture has a different balance between these
Each kernel has a different balance between these
Performance is a question of how well ap kernel's characteristics map to an architecture's characteristics

The Roofline Model: A pedagogical tool for program analysis and optimization





- For us, floating point performance (Gflop/s) is the metric of interest (typically double precision)
 ... but we could also consider SP or int
- Peak in-core performance can only be attained if:
 - fully exploit ILP, DLP, FMA, etc...
 - non-FP instructions don't sap instruction bandwidth
 - threads don't diverge (GPUs)
 - transcendental/non pipelined instructions are used sparingly
 - branch mispredictions are rare
- To exploit a form of in-core parallelism, it must be:
 - Inherent in the algorithm
 - Expressed in the high level implementation
 - Explicit in the generated code



ParLab Summer Retreat Samuel Williams, David Patterson



Communication

✤ For us, DRAM bandwidth (GB/s) is the metric of interest

Peak bandwidth can only be attained if certain optimizations are employed:

- Few unit stride streams
- NUMA allocation and usage
- SW Prefetching
- Memory Coalescing (GPU)

The Roofline Model: A pedagogical tool for program analysis and optimization 6

BERKELEY PAR LAI



Locality



3Cs model

for caches

- Computation is free, Communication is expensive.
- Maximize locality to minimize communication
- * There is a lower limit to communication: compulsory traffic
- Hardware changes can help minimize communication
 - Larger cache capacities minimize capacity misses
 - Higher cache associativities minimize conflict misses
 - Non-allocating caches minimize compulsory traffic
- Software optimization can also help minimize communication
 - Padding avoids conflict misses
 - Blocking avoids capacity misses
 - Non-allocating stores minimize compulsory traffic

The Roofline Model: A pedagogical tool for program analysis and optimization

> ParLab Summer Retreat Samuel Williams, David Patterson

> > Three Classes of Locality

Temporal Locality

UTURE

- reusing data (either registers or cache lines) multiple times
- amortizes the impact of limited bandwidth.
- transform loops or algorithms to maximize reuse.
- Spatial Locality
 - data is transferred from cache to registers in words.
 - However, data is transferred to the cache in 64-128Byte lines
 - using every word in a line maximizes spatial locality.
 - transform data structures into structure of arrays (SoA) layout
- Sequential Locality
 - Many memory address patterns access cache lines sequentially.
 - CPU's hardware stream prefetchers exploit this observation to hide speculatively load data to memory latency.
 - Transform loops to generate (a few) long, unit-stride accesses.
 LAWRENCE BERKELEY NATIONAL LABORATORY

8

\sim

- goal: integrate <u>in-core performance</u>, <u>memory bandwidth</u>, and <u>locality</u> into a single readily understandable <u>performance figure</u>
- graphically show the penalty associated with <u>not including</u> certain software optimizations
- Roofline model will be unique to each architecture

AJProença, Advanced Architectures, MiEl, UMinho, 2016/17

11

Key elements in the Roofline Model

 <u>x-axis</u>: the "operational intensity", operations per byte of RAM traffic, <u>Flops/byte</u> (traffic between caches and memory)
 <u>y-axis</u>: the attainable floating-point performance, <u>GFlops/sec</u> includes both peak <u>processor/memory</u> performance
 <u>peak processor FP performance</u>: a horizontal line computed from the processor specs
• <u>peak memory performance</u> : bounds the max FP performance of the memory system for a given operational intensity
 <u>for each kernel</u>: its performance is a point on a vertical line that crosses the x-axis on the kernel operational intensity



Arithmetic Intensity

Samuel Williams



- * True Arithmetic Intensity (AI) ~ Total Flops / Total DRAM Bytes
- Some HPC kernels have an arithmetic intensity that scales with problem size (increased temporal locality)
- Others have constant intensity
- * Arithmetic intensity is ultimately limited by compulsory traffic
- Arithmetic intensity is diminished by conflict or capacity misses.

LAWRENCE BERKELEY NATIONAL LABORATORY



Additional notes



time

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

Parallelism in a modern compute node



15

Parallel and shared resources within a shared-memory node



Basics of performance modeling fo numerical applications: Roofline model and beyond



- Dual Socket (NUMA)
- limited HW stream prefetchers
- quad-core (8 total)
- 2.3GHz
- 2-way SIMD (DP)
- separate FPMUL and FPADD datapaths
- 4-cycle FP latency



Assuming expression of parallelism is the challenge on this architecture, what would the roofline model look like ?

LAWRENCE BERKELEY NATIONAL LABORATORY

















flop:DRAM byte ratio

ParLab Summer Retreat Samuel Williams, David Patterson















- Time is the sum of communication time and computation time.
- The result is that flop/s grows asymptotically.

LAWRENCE BERKELEY NATIONAL LABORATORY





TECHNOLOGIES

Samuel Williams

GROU

- Thus far, we assumed a synergy between streaming applications and bandwidth (proxied by the STREAM benchmark)
- STREAM is NOT a good proxy for short stanza/random cacheline access patterns as memory latency (instead of just bandwidth) is being exposed.
- Thus one might conceive of alternate memory benchmarks to provide a bandwidth upper bound (ceiling)
- Similarly, if data is primarily local in the LLC cache, one should construct rooflines based on LLC bandwidth and flop:LLC byte ratios.
- For GPUs/accelerators, PCIe bandwidth can be an impediment. Thus one can construct a roofline model based on PCIe bandwidth and the flop:PCIe byte ratio.

LAWRENCE BERKELEY NATIONAL LABORATORY



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

The Roofline model: Hardware vs. Software



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17



Some more examples

NTERNATIONAL CONFERENCE ON PARALLEL

Some more examples





- Arising from HPC kernels, its no surprise roofline use DP Flop/s.
- Of course, it could use
 - SP flop/s,
 - integer ops,
 - bit operations,
 - pairwise comparisons (sorting),
 - graphics operations,
 - etc...