# Master Informatics Eng.

2019/20

*A.J.Proença*

## The Roofline Performance Model
### *(most slides are borrowed)*

- ❖ Multicore guarantees neither good scalability nor good (attained) performance
- ❖ Performance and scalability can be extremely non-intuitive even to computer scientists

- ❖ Success of the multicore paradigm seems to be premised upon their programmability
- ❖ To that end, one must understand the limits to both scalability and efficiency.

- How can we empower programmers?

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

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, off-chip 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

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

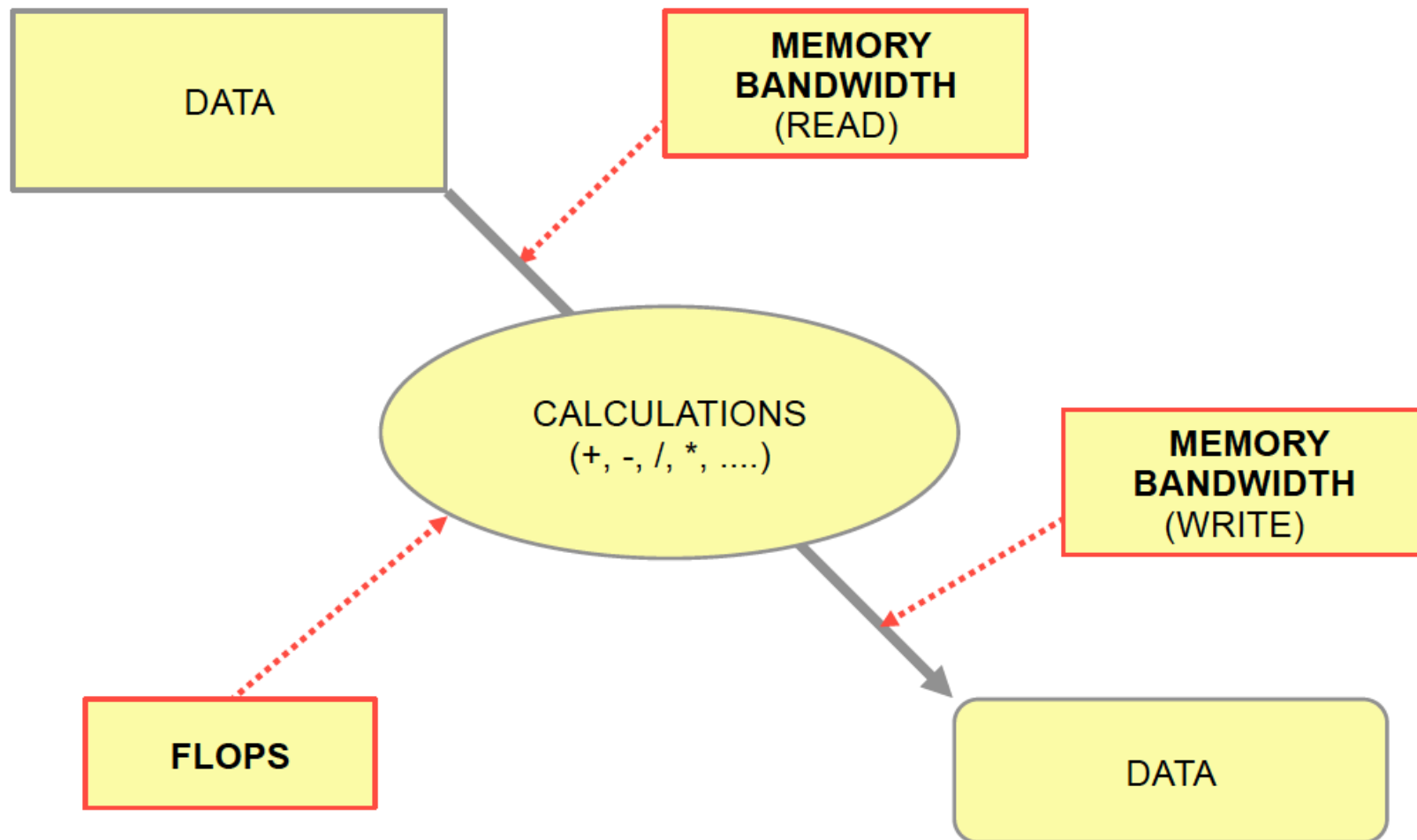BY SAMUEL WILLIAMS, ANDREW WATERMAN, AND DAVID PATTERSON

# Roofline:

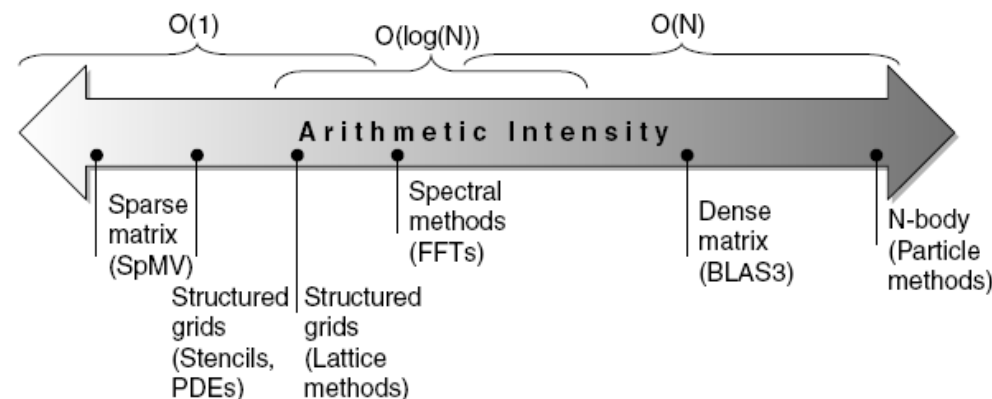# Performance Limiting Factors

Leopold Grinberg
IBM, T.J. Watson Research Center, USA

ICSC 2014, Shanghai, China

IBM

# 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

- ❖ There are three principal 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 an 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

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

# 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)

# Locality

- ❖ 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

3Cs model
for caches

- ❖ 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

❖ Temporal Locality
- 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.**
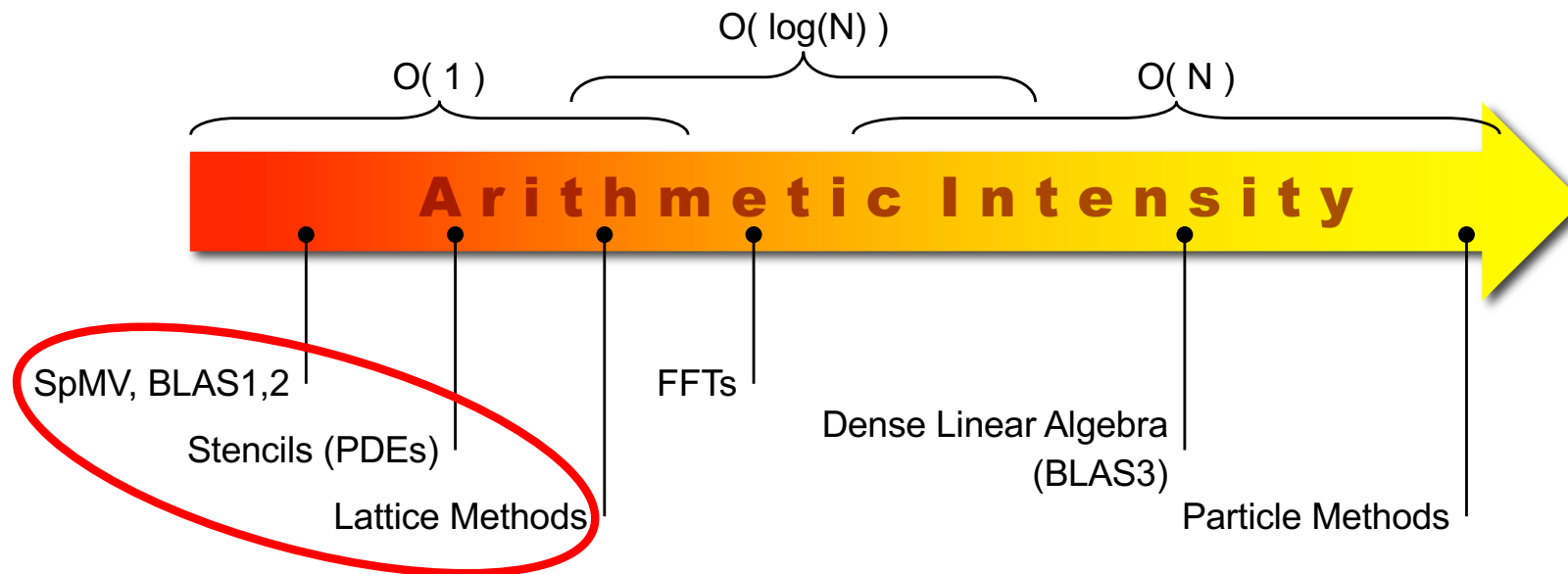
# *Preliminary notes in the Roofline Model*

- 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
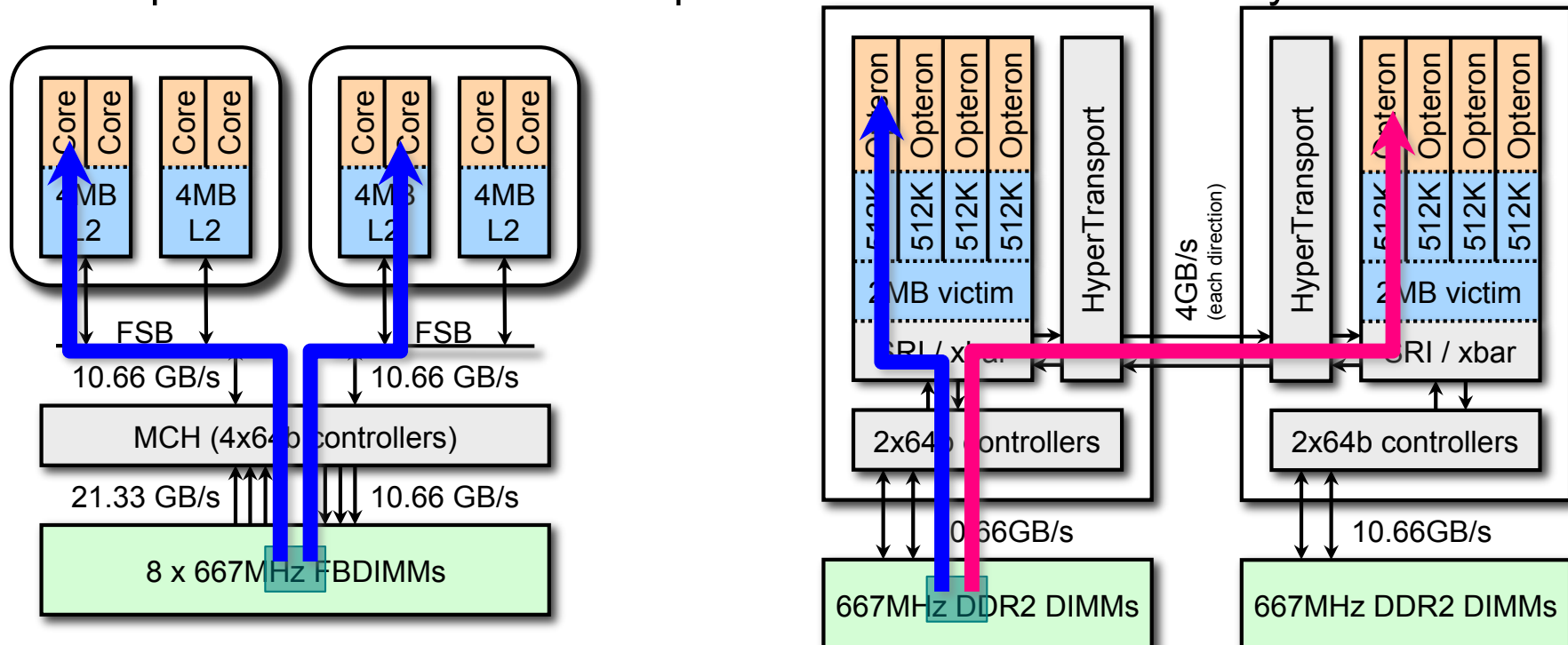
# Key elements in the Roofline Model

- *x*-axis: the "operational intensity", operations per byte of RAM traffic, Flops/byte (traffic between caches and memory)

- *y*-axis: the attainable floating-point performance, GFlops/sec includes both peak processor/memory  performance

- *peak processor FP performance*: a horizontal line computed from the processor specs

- *peak memory performance*: bounds the max FP performance of the memory system for a given operational intensity

- *for each kernel*: its performance is a point on a vertical line that crosses the *x*-axis on the kernel operational intensity
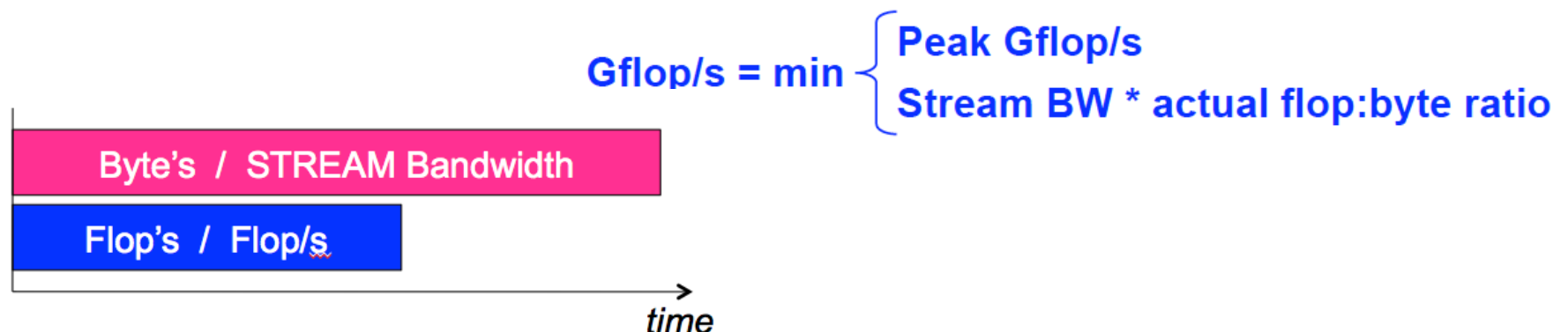
O( log(N) )

O( 1 )

O( N )

**A r i t h m e t i c   I n t e n s i t y**

SpMV, BLAS1,2

FFTs

Stencils (PDEs)

Dense Linear Algebra
(BLAS3)

Lattice Methods

Particle Methods

❖ **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.

F U T U R E    T E C H N O L O G I E S    G R O U P

- ❖ Recent multicore SMPs have integrated the memory controllers on chip.
- ❖ As a result, memory-access is non-uniform (NUMA)
- ❖ That is, the bandwidth to read a given address varies dramatically among between cores
- ❖ Exploit NUMA (affinity+first touch) when you malloc/init data.
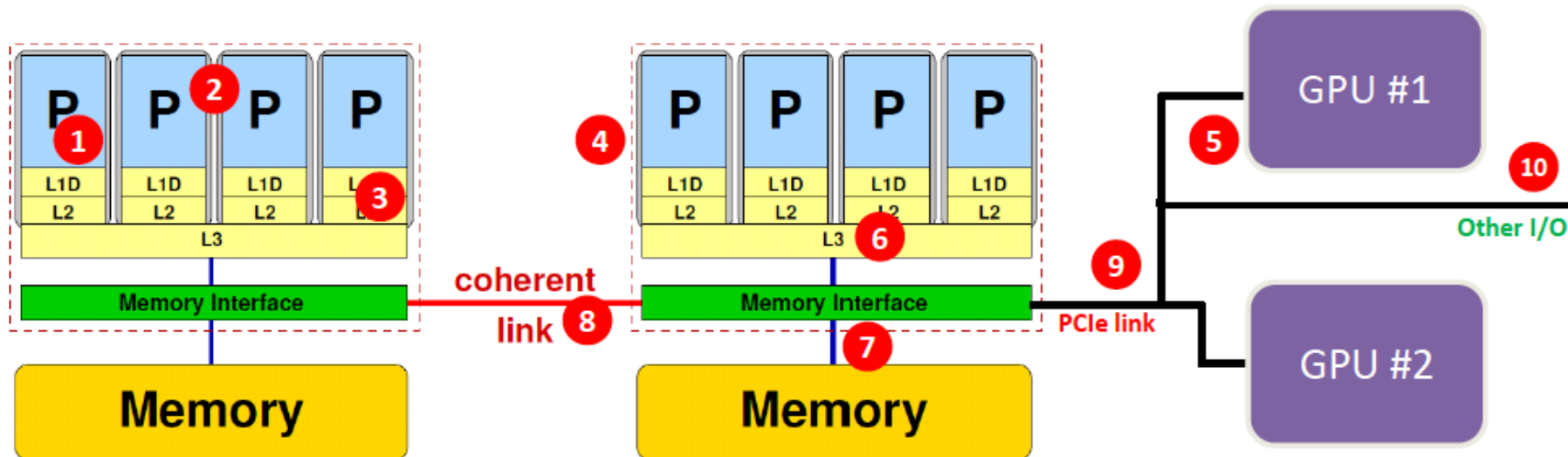- ❖ Concept is similar to data decomposition for distributed memory

- Memory bandwidth #'s collected via micro benchmarks (or the STREAM benchmark)

- Computation #'s derived from optimization manuals (pencil and paper)

- Assume complete overlap of either communication or computation =>

$$\text{Gflop/s} = \min \begin{cases} \text{Peak Gflop/s} \\ \text{Stream BW} * \text{actual flop:byte ratio} \end{cases}$$

Byte's / STREAM Bandwidth

Flop's / Flop/s

time

# Parallelism in a modern compute node

- **Parallel and shared resources within a shared-memory node**



## Parallel resources:

- Execution/SIMD units **1**
- Cores **2**
- Inner cache levels **3**
- Sockets / ccNUMA domains **4**
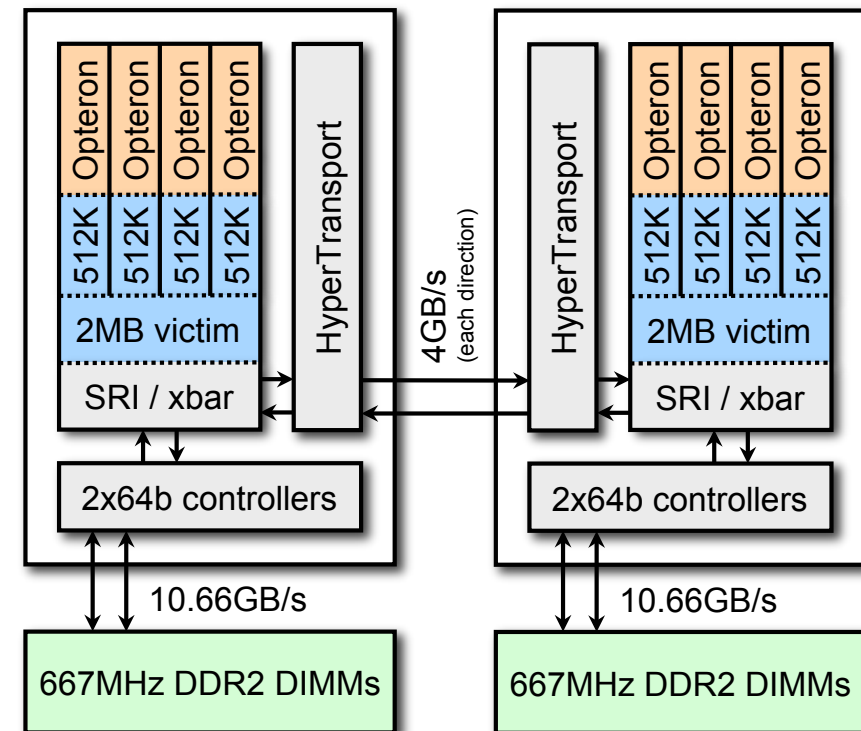- Multiple accelerators **5**

## Shared resources ("bottlenecks"):

- Outer cache level per socket **6**
- Memory bus per socket **7**
- Intersocket link **8**
- PCIe bus(es) **9**
- Other I/O resources **10**

**Where is the bottleneck for your application?**

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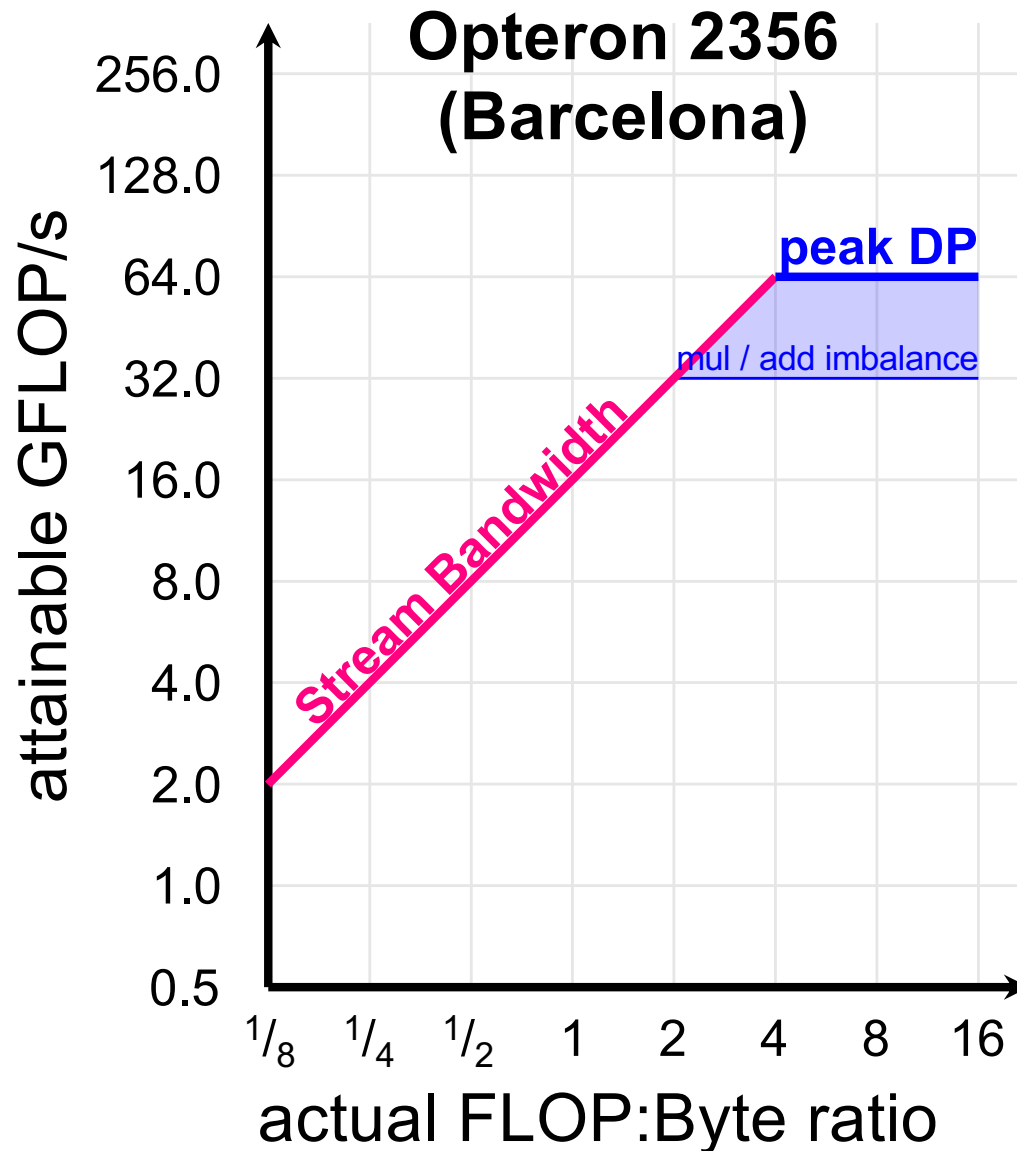Roofline model and beyond

Georg Hager, Jan Treibig, Gerhard Wellein

- ❖ Consider the Opteron 2356:
  - ▪ Dual Socket (NUMA)
  - ▪ limited HW stream prefetchers
  - ▪ quad-core (8 total)
  - ▪ 2.3GHz
  - ▪ 2-way SIMD (DP)
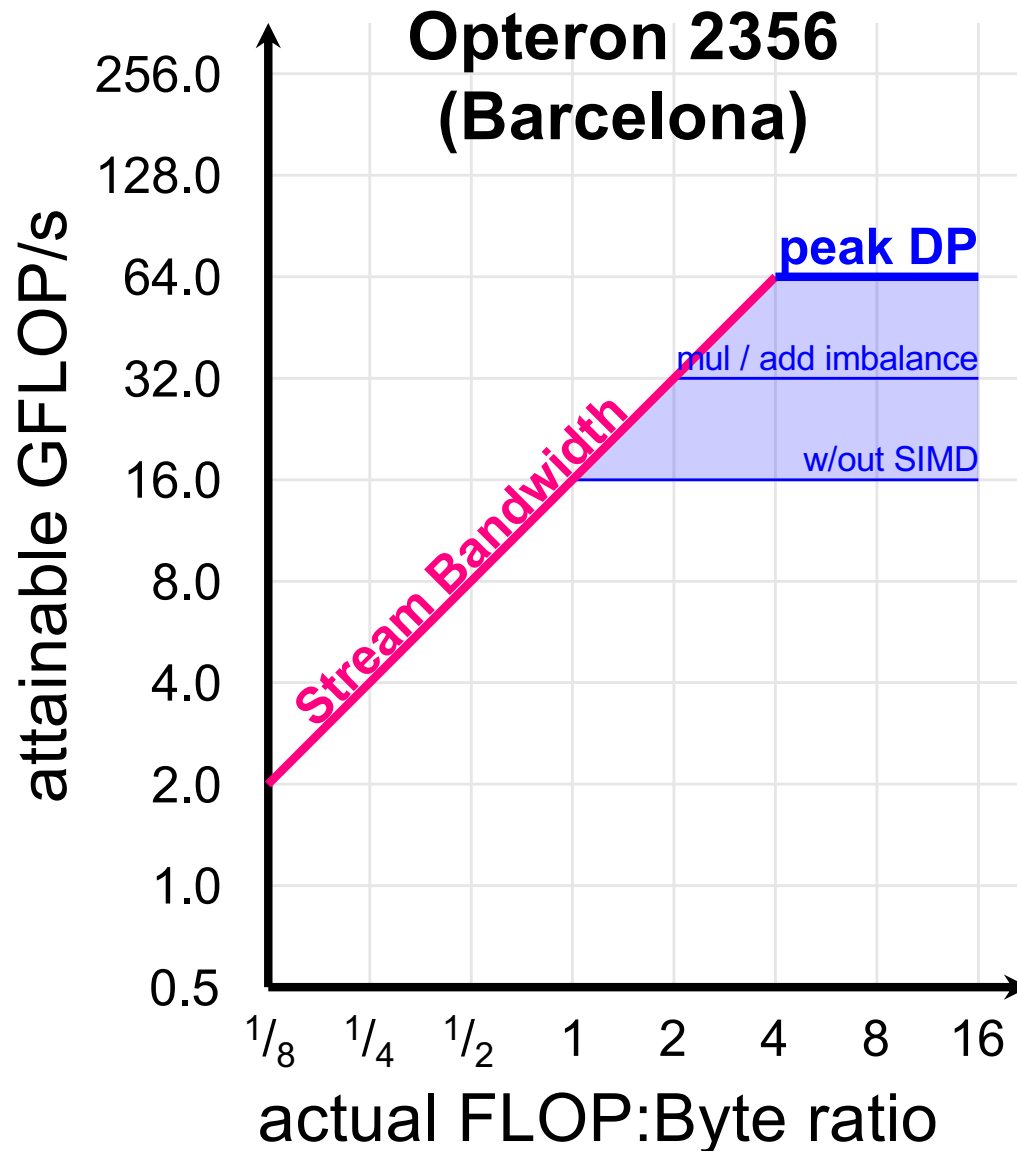  - ▪ separate FPMUL and FPADD datapaths
  - ▪ 4-cycle FP latency



- ❖ Assuming **expression of parallelism** is the challenge on this architecture, what would the roofline model look like ?
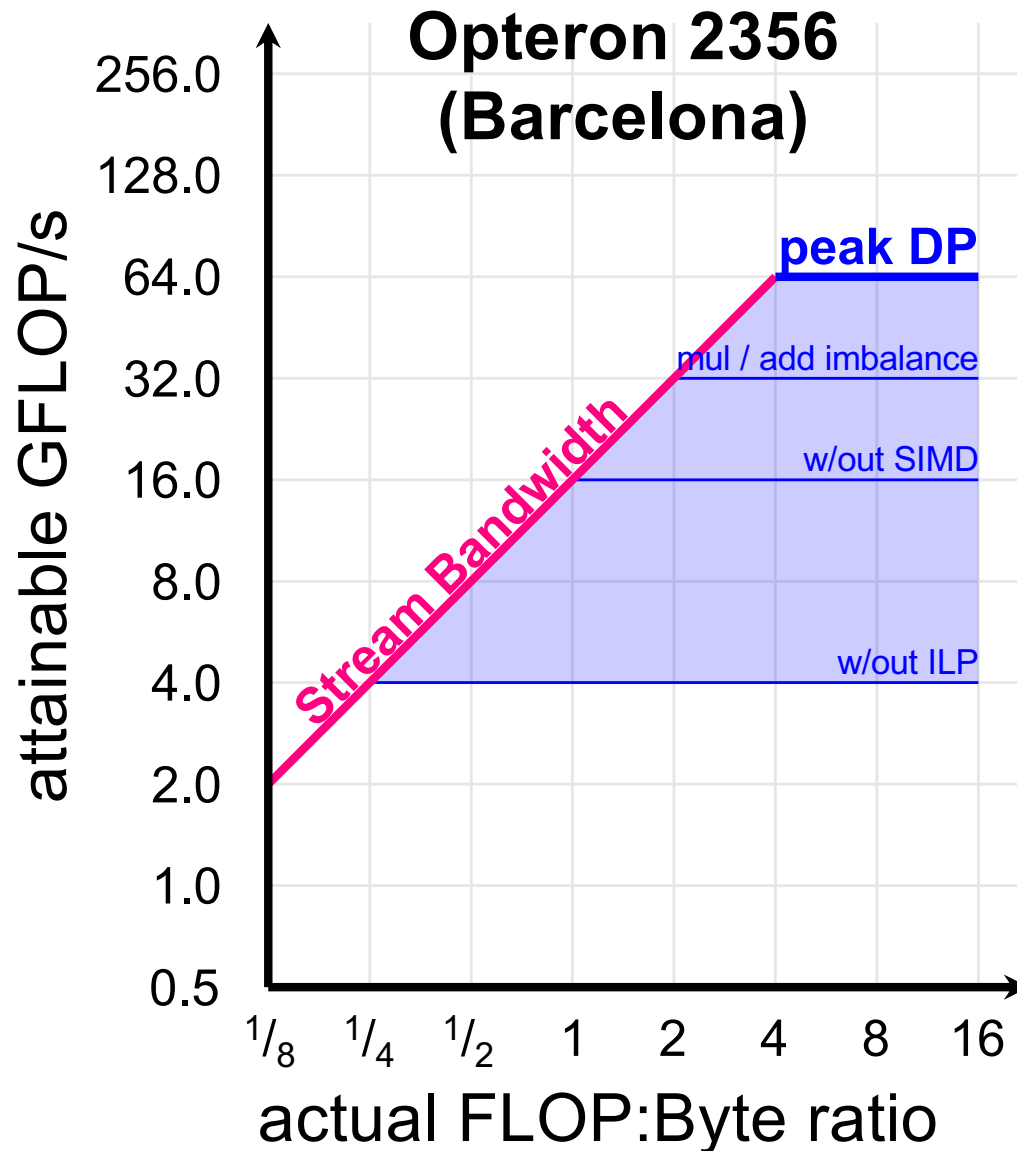
**Opteron 2356 (Barcelona)**

- ❖ Plot on log-log scale
- ❖ Given AI, we can easily bound performance
- ❖ But architectures are much more complicated

- ❖ We will bound performance as we eliminate specific forms of in-core parallelism

# Roofline Model
## computational ceilings

**Opteron 2356 (Barcelona)**

*attainable GFLOP/s* vs *actual FLOP:Byte ratio*

- **peak DP** (64.0)
- **Stream Bandwidth**
- mul / add imbalance

Y-axis values: 256.0, 128.0, 64.0, 32.0, 16.0, 8.0, 4.0, 2.0, 1.0, 0.5
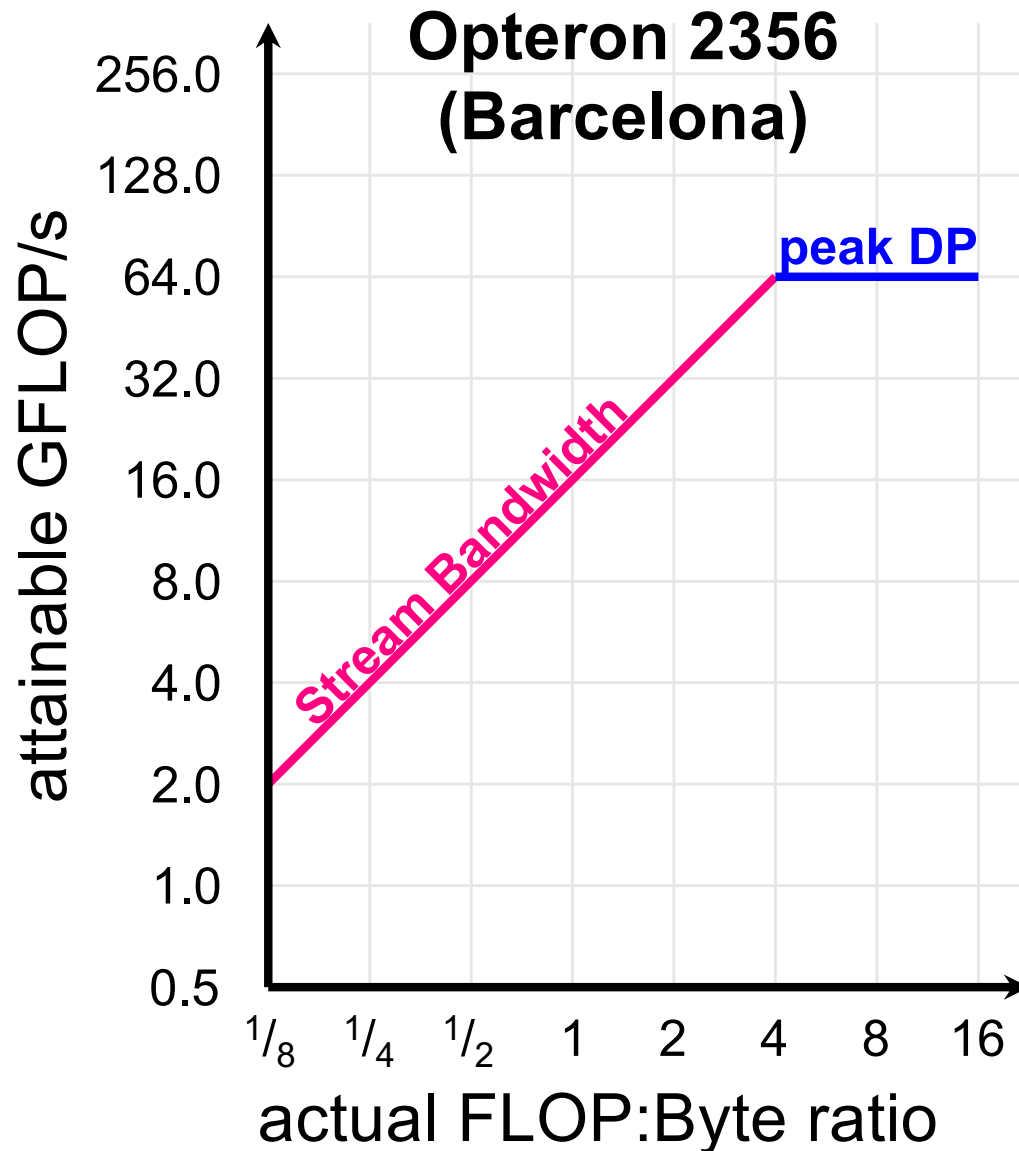
X-axis values: 1/8, 1/4, 1/2, 1, 2, 4, 8, 16

- ❖ Opterons have dedicated multipliers and adders.
- ❖ If the code is dominated by adds, then attainable performance is half of peak.
- ❖ We call these **Ceilings**
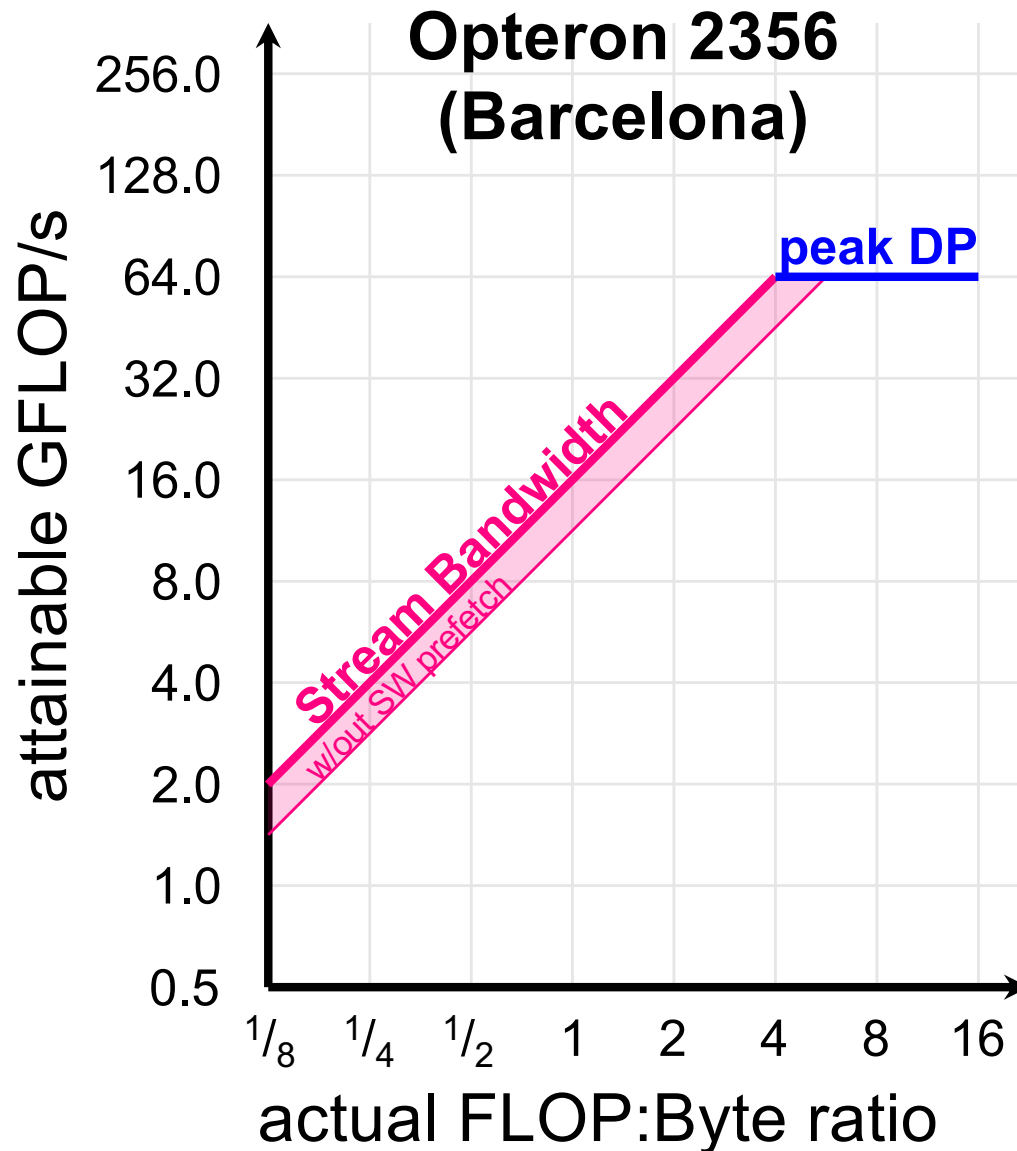- ❖ They act like constraints on performance

# Roofline Model
## computational ceilings



Opteron 2356 (Barcelona)

- ❖ Opterons have 128-bit datapaths.
- ❖ If instructions aren't SIMDized, attainable performance will be halved

# Roofline Model
## computational ceilings

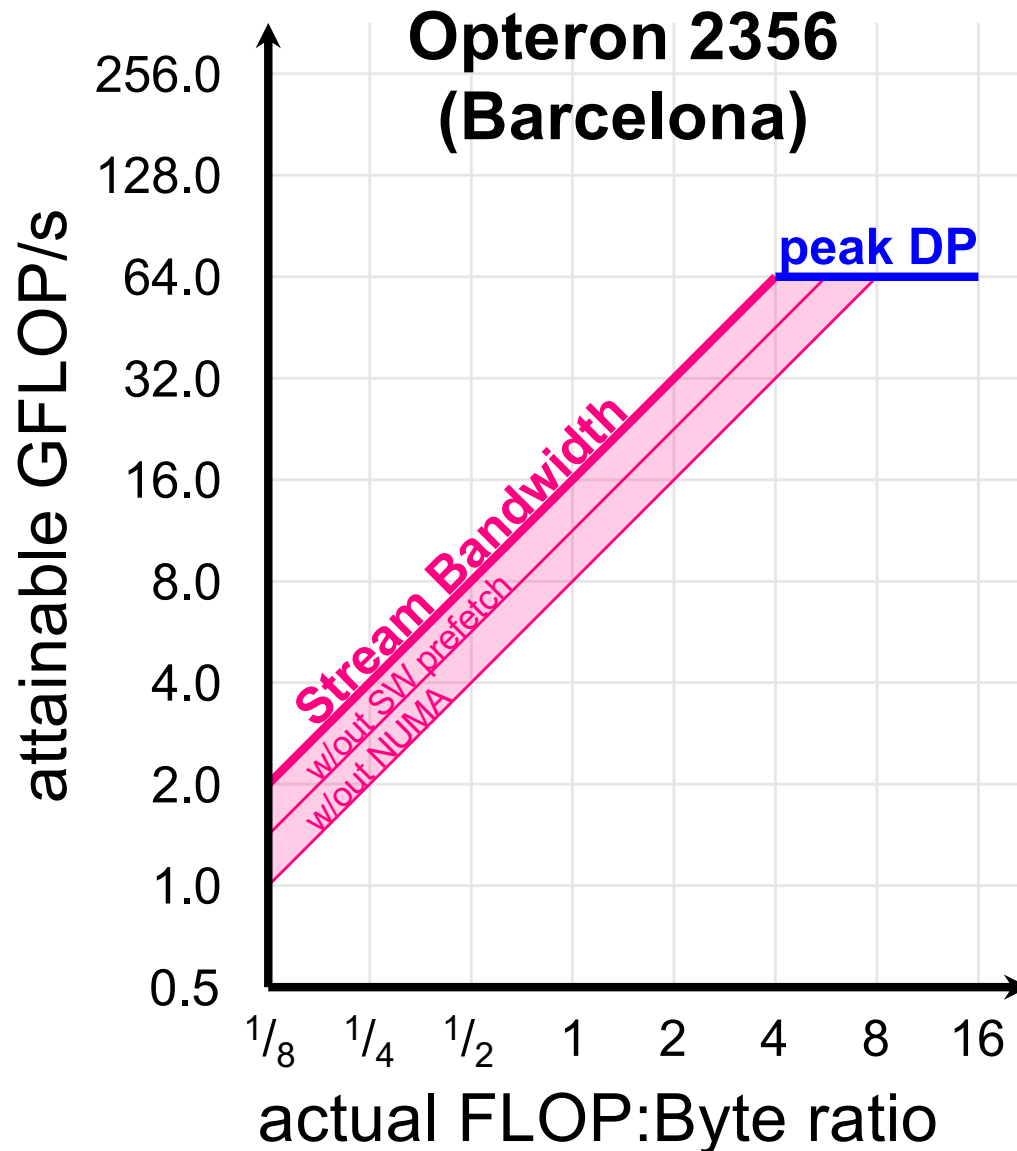F U T U R E     T E C H N O L O G I E S     G R O U P



**Opteron 2356 (Barcelona)**

- ❖ On Opterons, floating-point instructions have a 4 cycle latency.

- ❖ If we don't express 4-way ILP, performance will drop by as much as 4x

**Opteron 2356 (Barcelona)**

❖ We can perform a similar exercise taking away parallelism from the memory subsystem

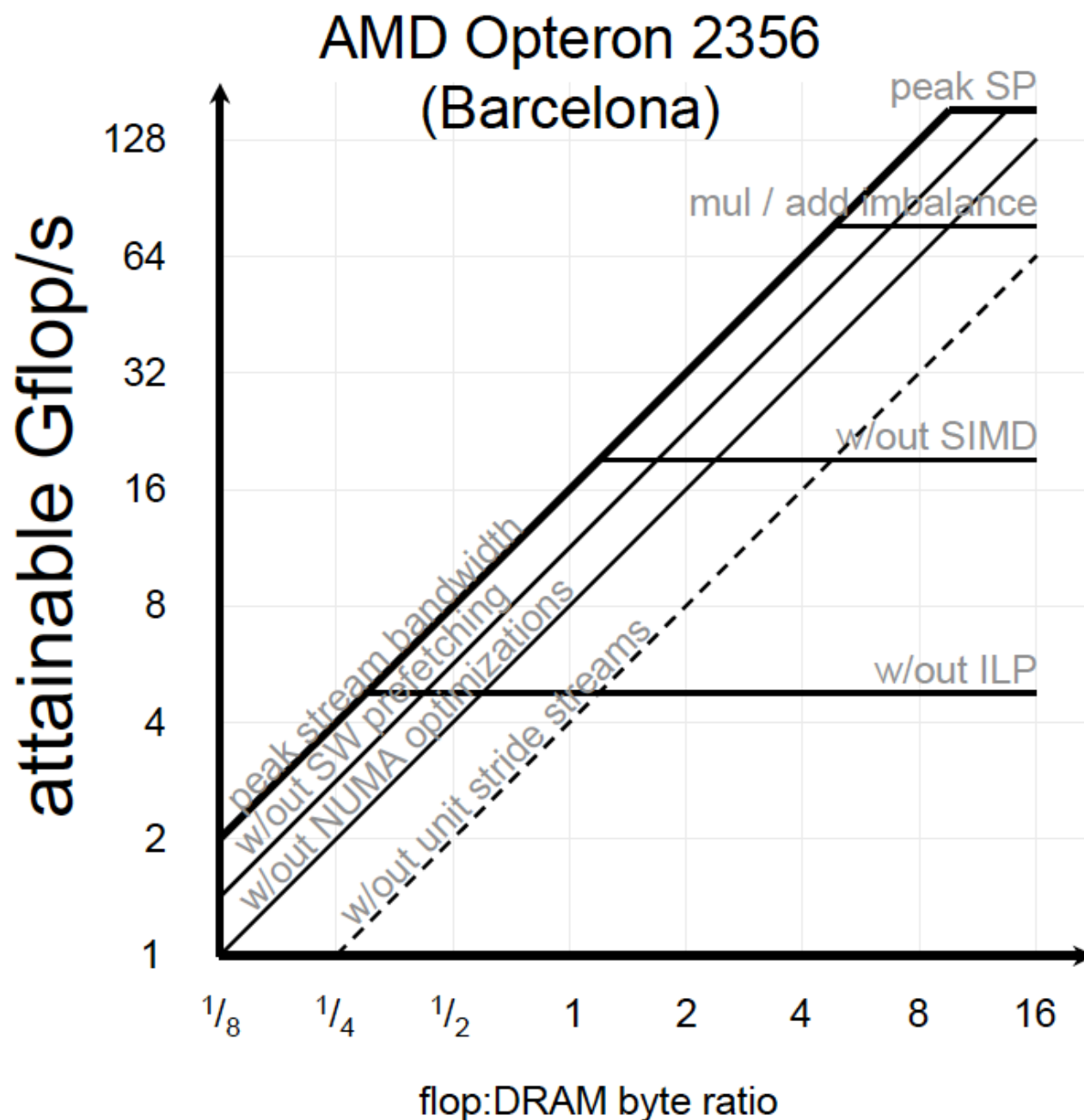❖ Explicit software prefetch instructions are required to achieve peak bandwidth

# Roofline Model
## communication ceilings

**Opteron 2356 (Barcelona)**

attainable GFLOP/s vs actual FLOP:Byte ratio

peak DP

Stream Bandwidth
w/out SW prefetch
w/out NUMA

❖ Opterons are NUMA

❖ As such memory traffic must be correctly balanced among the two sockets to achieve good Stream bandwidth.

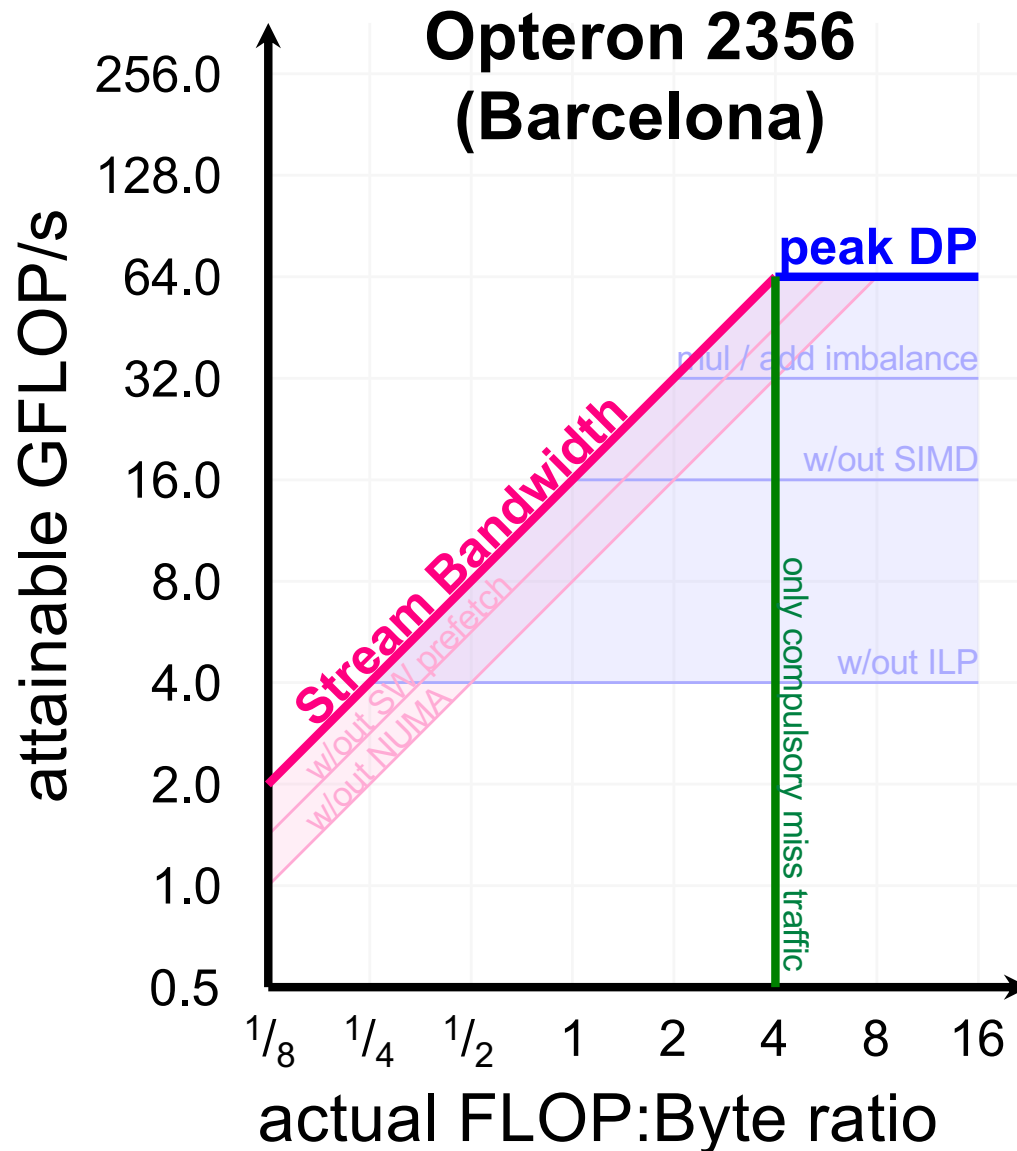❖ We could continue this by examining strided or random memory access patterns

❖ We may bound performance based on the combination of expressed in-core parallelism and attained bandwidth.

F U T U R E    T E C H N O L O G I E S    G R O U P



Opteron 2356 (Barcelona) roofline chart showing attainable GFLOP/s vs actual FLOP:Byte ratio, with Stream Bandwidth, w/out SW prefetch, w/out NUMA lines, peak DP, mul / add imbalance, w/out SIMD, w/out ILP walls, and "only compulsory miss traffic" line.
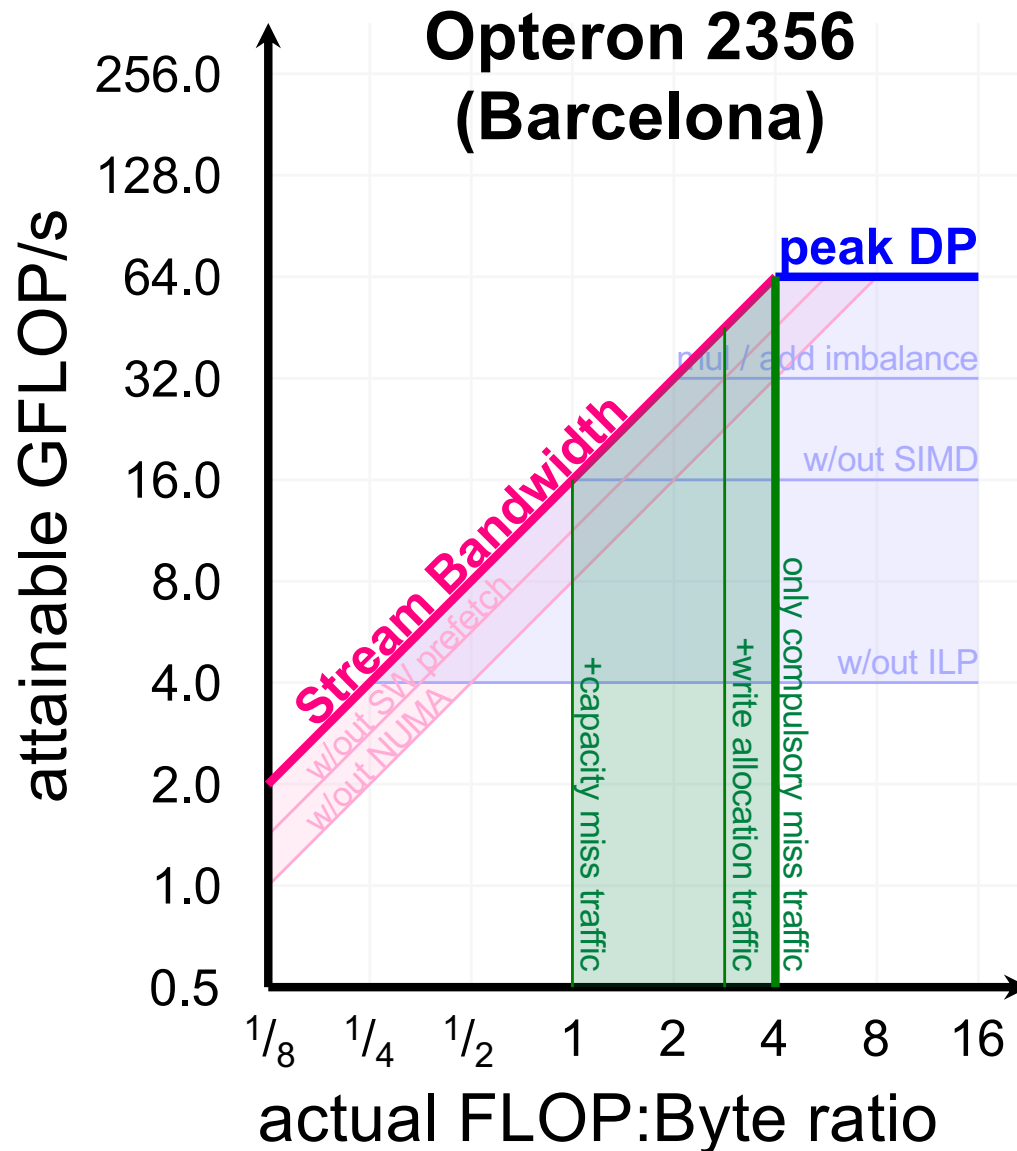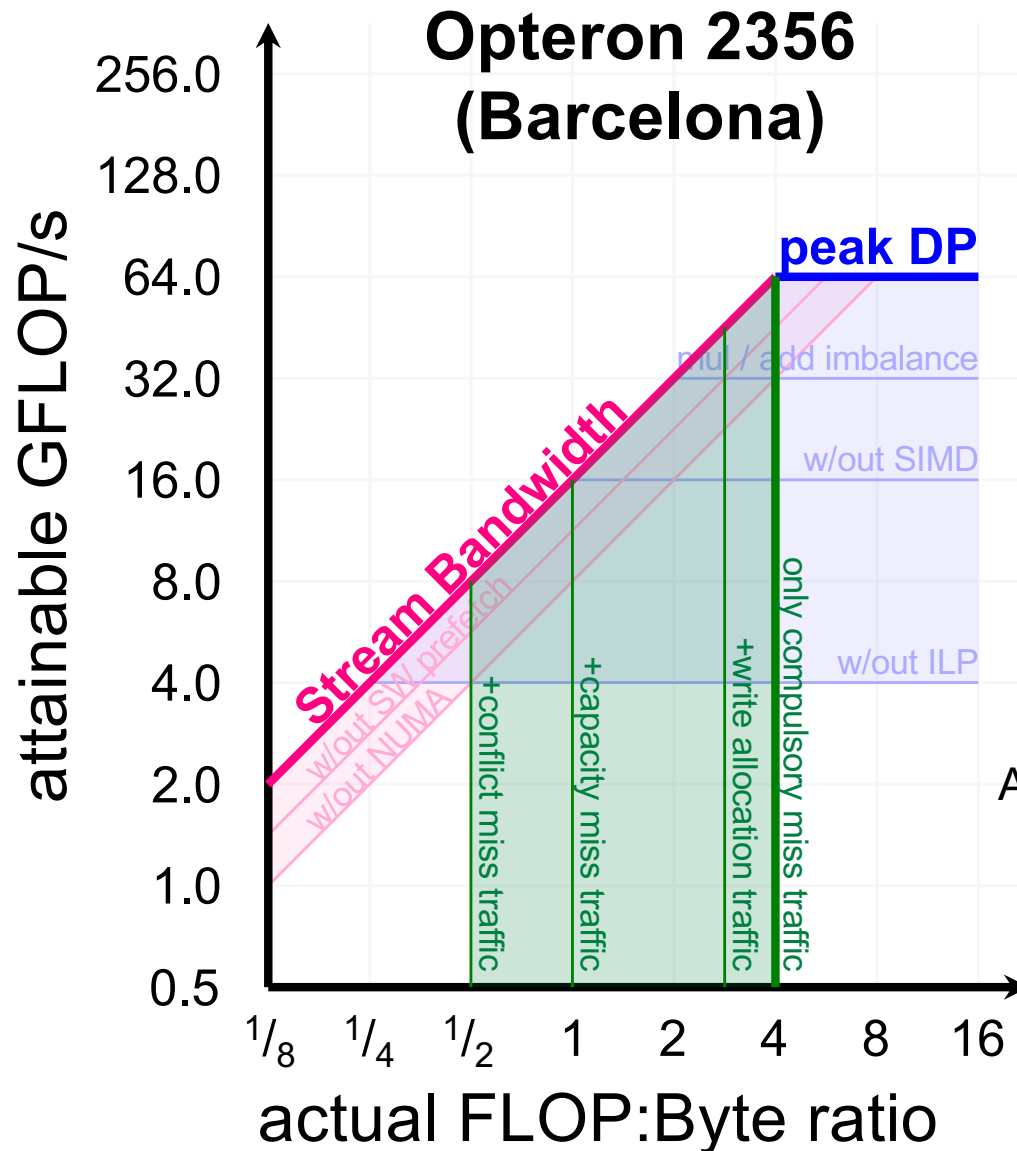
- ❖ Remember, memory traffic includes more than just compulsory misses.
- ❖ As such, actual arithmetic intensity may be substantially lower.
- ❖ Walls are unique to the architecture-kernel combination

$$AI = \frac{FLOPs}{Compulsory\ Misses}$$

- ❖ Remember, memory traffic includes more than just compulsory misses.
- ❖ As such, actual arithmetic intensity may be substantially lower.
- ❖ Walls are unique to the architecture-kernel combination

$$AI = \frac{FLOPs}{Allocations + Compulsory\ Misses}$$

FUTURE TECHNOLOGIES GROUP



**Opteron 2356 (Barcelona)**
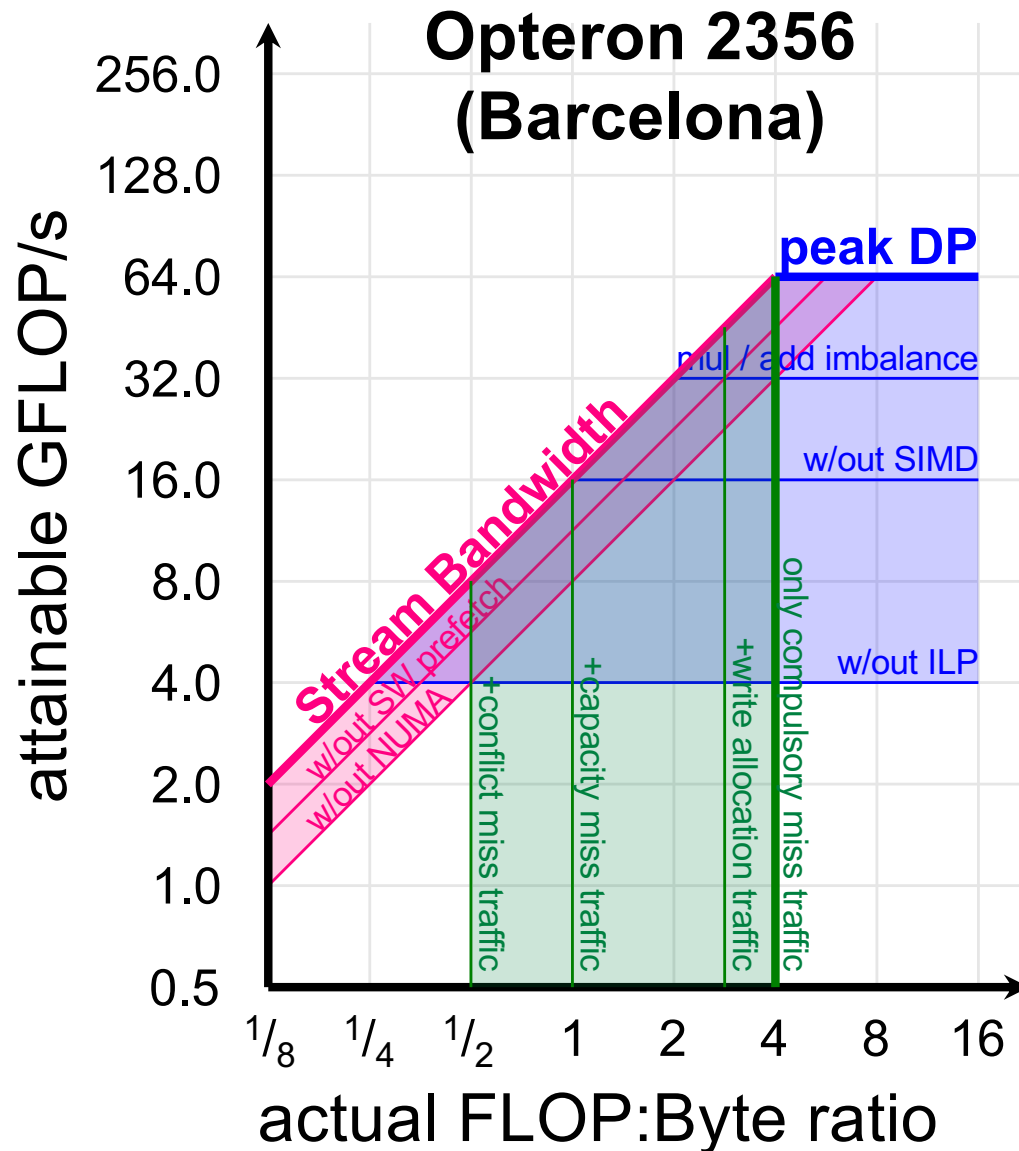
❖ Remember, memory traffic includes more than just compulsory misses.

❖ As such, actual arithmetic intensity may be substantially lower.

❖ Walls are unique to the architecture-kernel combination

$$AI = \frac{FLOPs}{Capacity + Allocations + Compulsory}$$

Opteron 2356 (Barcelona)

- ❖ Remember, memory traffic includes more than just compulsory misses.
- ❖ As such, actual arithmetic intensity may be substantially lower.
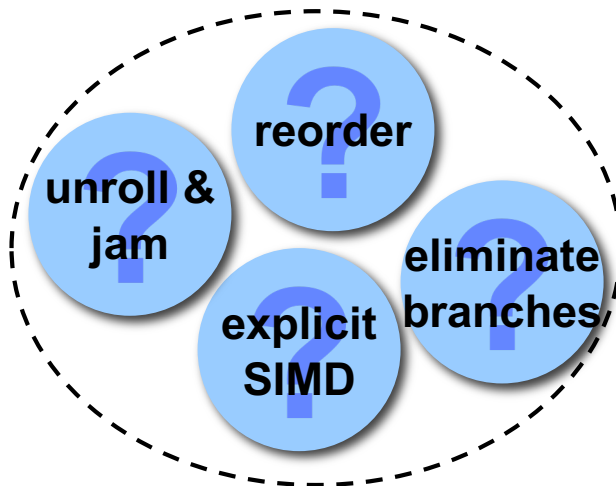- ❖ Walls are unique to the architecture-kernel combination

$$AI = \frac{FLOPs}{Conflict + Capacity + Allocations + Compulsory}$$

**Opteron 2356 (Barcelona)**

❖ Optimizations remove these walls and ceilings which act to constrain performance.
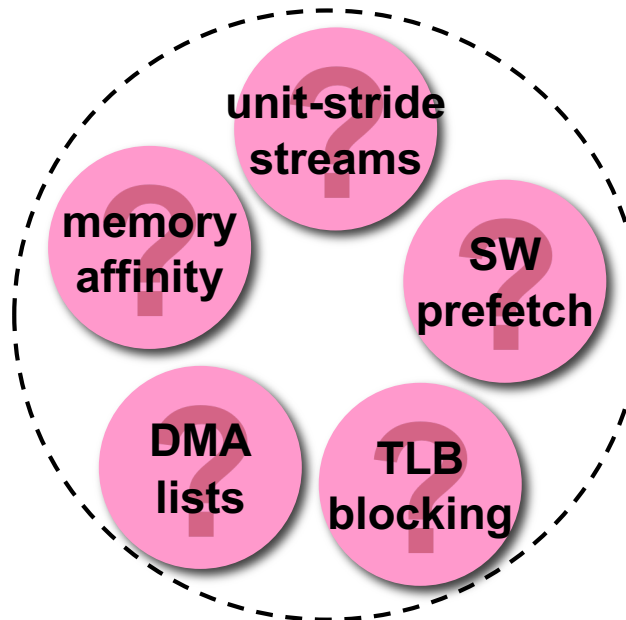
# Optimization Categorization

## Maximizing In-core Performance

- **Exploit in-core parallelism (ILP, DLP, etc…)**
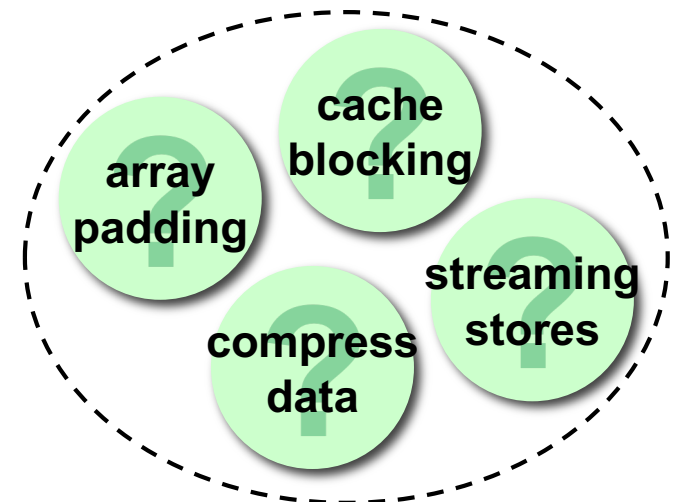
- **Good (enough) floating-point balance**

## Maximizing Memory Bandwidth

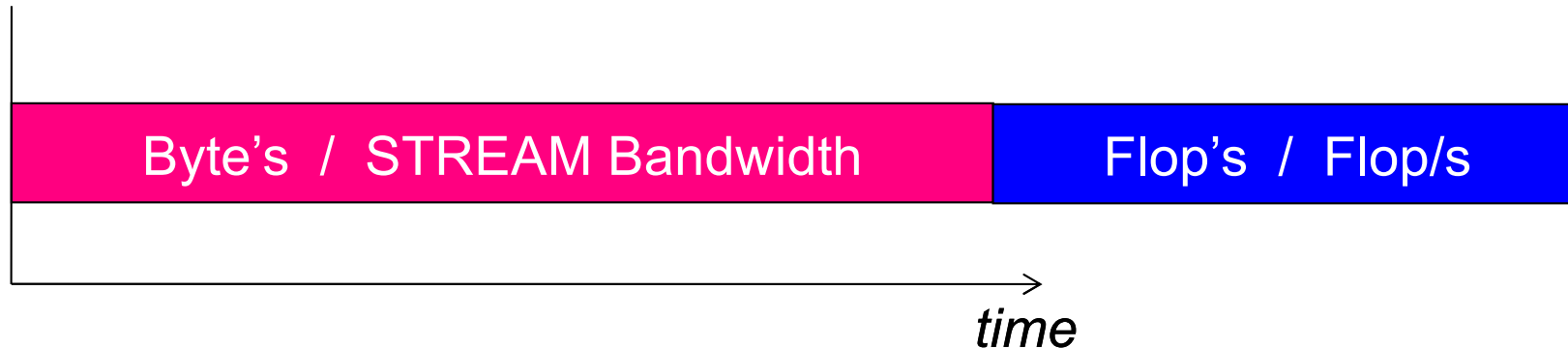- **Exploit NUMA**

- **Hide memory latency**

- **Satisfy Little's Law**

## Minimizing Memory Traffic

**Eliminate:**
- **Capacity misses**
- **Conflict misses**
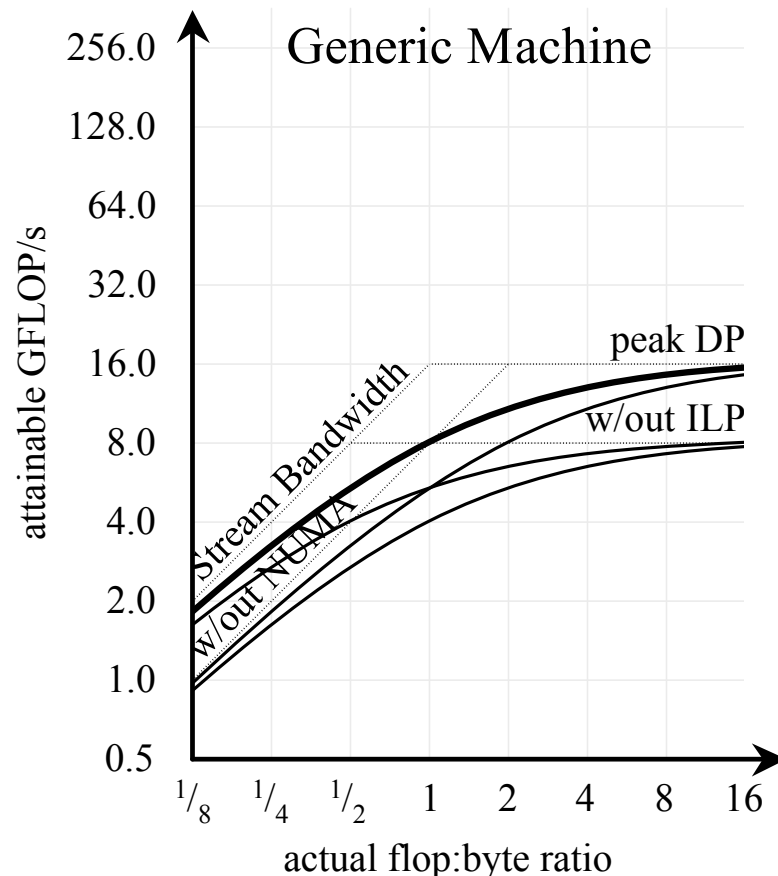- **Compulsory misses**
- **Write allocate behavior**

❖ Previously, we assumed perfect overlap of communication or computation.

❖ What happens if there is a dependency (either inherent or by a lack of optimization) that serializes communication and computation ?

| Byte's / STREAM Bandwidth | Flop's / Flop/s |
|---|---|

*time*

❖ Time is the sum of communication time and computation time.

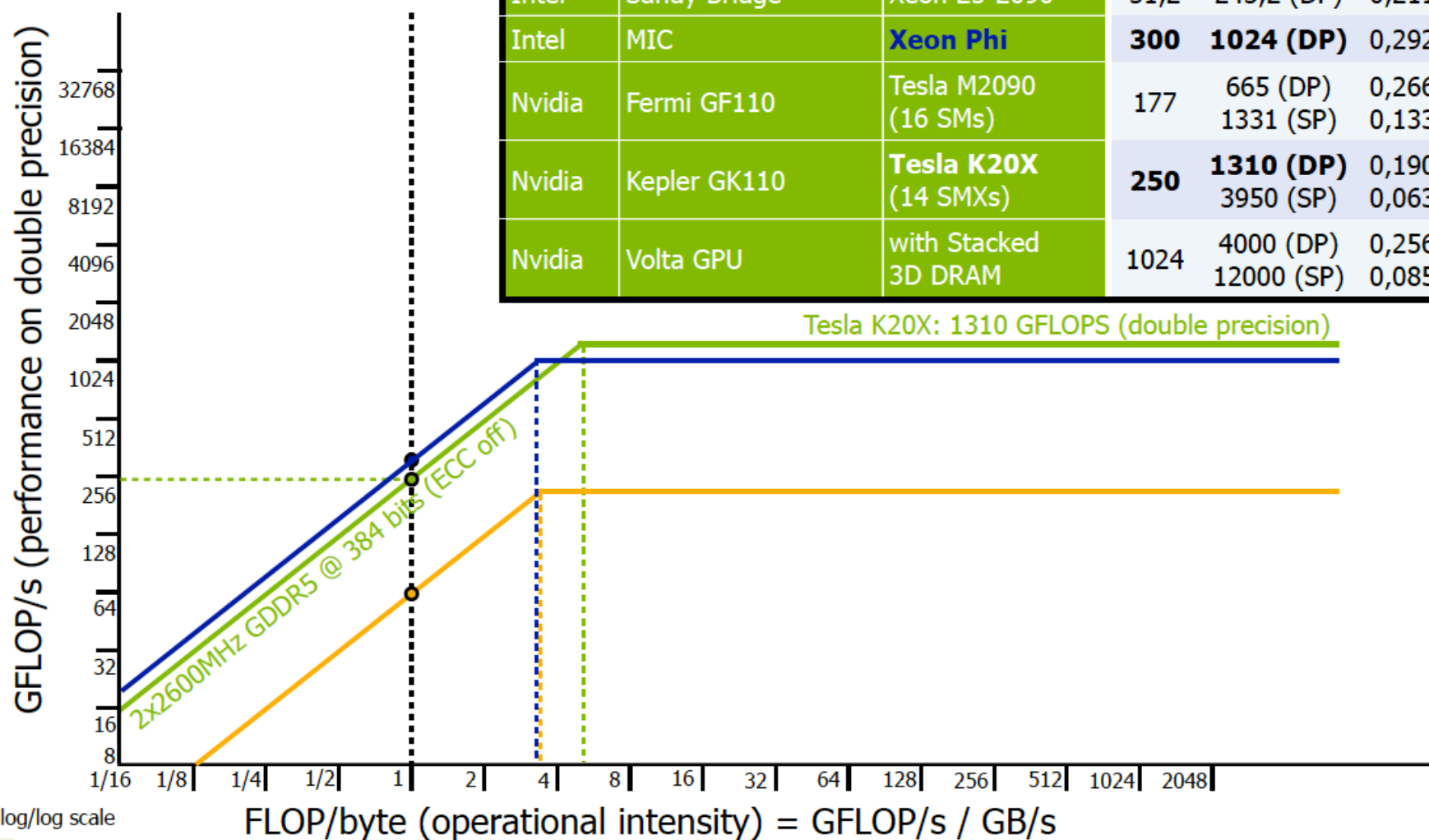❖ **The result is that flop/s grows asymptotically.**

Generic Machine

❖ Consider a generic machine

❖ If we can perfectly decouple and overlap communication with computation, the roofline is sharp/angular.

❖ However, without overlap, the roofline is smoothed, **and attainable performance is degraded by up to a factor of 2x.**

❖ 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.

# Platforms to compare

| Vendor | Microarchitecture | Model | GB/s. | GFLOP/s. | Byte/FLOP |
|--------|-------------------|-------|-------|----------|-----------|
| AMD | Bulldozer | **Opteron 6284** | **59,7** | **217,6 (DP)** | 0,235 |
| AMD | Souther Islands | Radeon HD7970 | 288 | 1010 (DP) | 0,285 |
| Intel | Sandy Bridge | Xeon E5-2690 | 51,2 | 243,2 (DP) | 0,211 |
| Intel | MIC | **Xeon Phi** | **300** | **1024 (DP)** | 0,292 |
| Nvidia | Fermi GF110 | Tesla M2090 (16 SMs) | 177 | 665 (DP) / 1331 (SP) | 0,266 / 0,133 |
| Nvidia | Kepler GK110 | **Tesla K20X** (14 SMXs) | **250** | **1310 (DP)** / 3950 (SP) | 0,190 / 0,063 |
| Nvidia | Volta GPU | with Stacked 3D DRAM | 1024 | 4000 (DP) / 12000 (SP) | 0,256 / 0,085 |



Roofline plot: GFLOP/s (performance on double precision) vs FLOP/byte (operational intensity) = GFLOP/s / GB/s, log/log scale. Tesla K20X: 1310 GFLOPS (double precision). 2x2600MHz GDDR5 @ 384 bits (ECC off)

*AJProença, Advanced Architectures, MiEI, UMinho, 2019/20*

# The Roofline model: Hardware vs. Software

# Some more examples

# *Some more examples*

## STREAM Benchmark Results on INTEL Hardware



Legend:
- E5-2670 v3 (Haswell)
- E5-2650 v2 (Ivy Bridge)
- E5-2620    (Sandy Bridge)
- Xeon Phi 7120 (Knights Corner)
- Xeon Phi 7250 (Knights Landing), DDR4
- Xeon Phi 7250 (Knights Landing), MCDRAM

*STREAM triad benchmark results for Knights Landing in comparison to three Intel Xeon generations (dual socket) and Knights Corner. A single KNL core achieves about the same memory bandwidth as Ivy Bridge and Sandy Bridge Xeons for small thread counts. Peak bandwidth is obtained with one thread per core; oversubscription reduces memory bandwidth slightly.*

*AJProença, Advanced Architectures, MiEI, UMinho, 2019/20*

F U T U R E     T E C H N O L O G I E S     G R O U P

- ❖ 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…