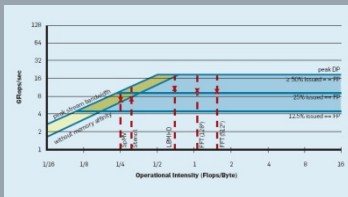# "Simple" performance modeling:
# The Roofline Model

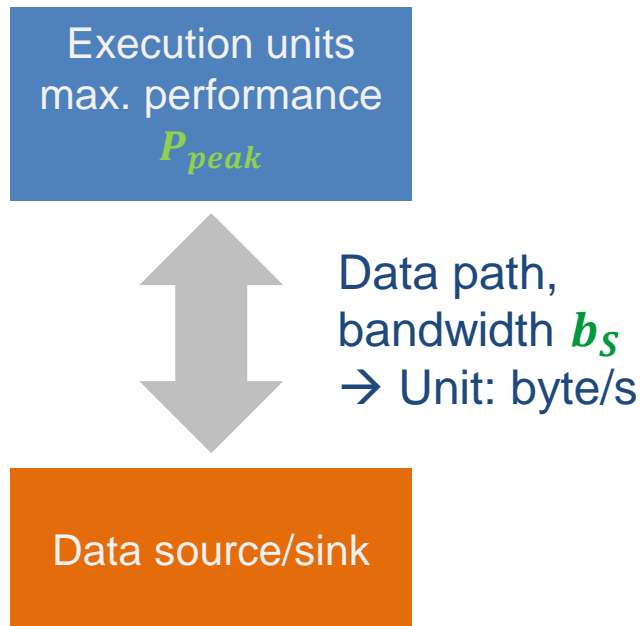Loop-based performance modeling: Execution vs. data transfer

R.W. Hockney and I.J. Curington: $f_{1/2}$: A parameter to characterize memory and communication bottlenecks.
Parallel Computing 10, 277-286 (1989).  DOI: 10.1016/0167-8191(89)90100-2

W. Schönauer: Scientific Supercomputing: Architecture and Use of Shared and Distributed Memory Parallel Computers.  Self-edition (2000)

S. Williams: Auto-tuning Performance on Multicore Computers.  UCB Technical Report No. UCB/EECS-2008-164. PhD thesis (2008)

# A simple performance model for loops

Simplistic view of the hardware:

Execution units
max. performance
$P_{peak}$

Data path,
bandwidth $b_S$
→ Unit: byte/s

Data source/sink

Simplistic view of the software:

```fortran
! may be multiple levels
do i = 1,<sufficient>
   <complicated stuff doing
    N flops causing
    V bytes of data transfer>
enddo
```

Computational intensity
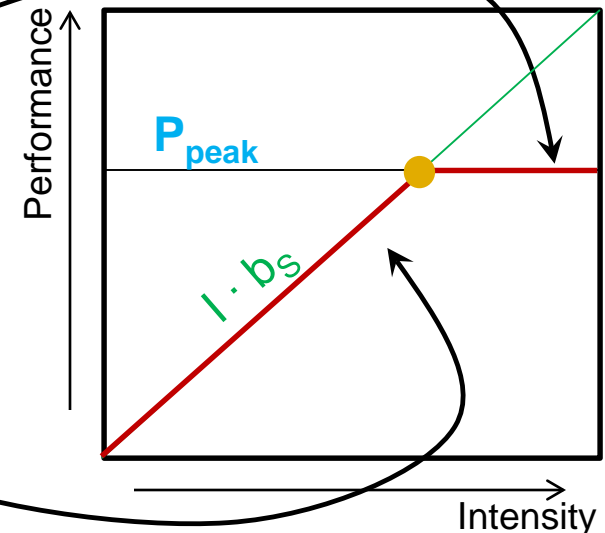$$I = \frac{N}{V}$$
→ Unit: flop/byte

# Naïve Roofline Model

How fast can tasks be processed? **$P$ [flop/s]**

The bottleneck is either

- The execution of work: $P_{\text{peak}}$    [flop/s]
- The data path: $I \cdot b_S$    [flop/byte x byte/s]

$$P = \min(P_{\text{peak}}, I \cdot b_S)$$

This is the "Naïve Roofline Model"

- High intensity: P limited by execution
- Low intensity: P limited by data transfer
- "Knee" at $P_{max} = I \cdot b_S$:
  Best use of resources
- Roofline is an "optimistic" model ("light speed")

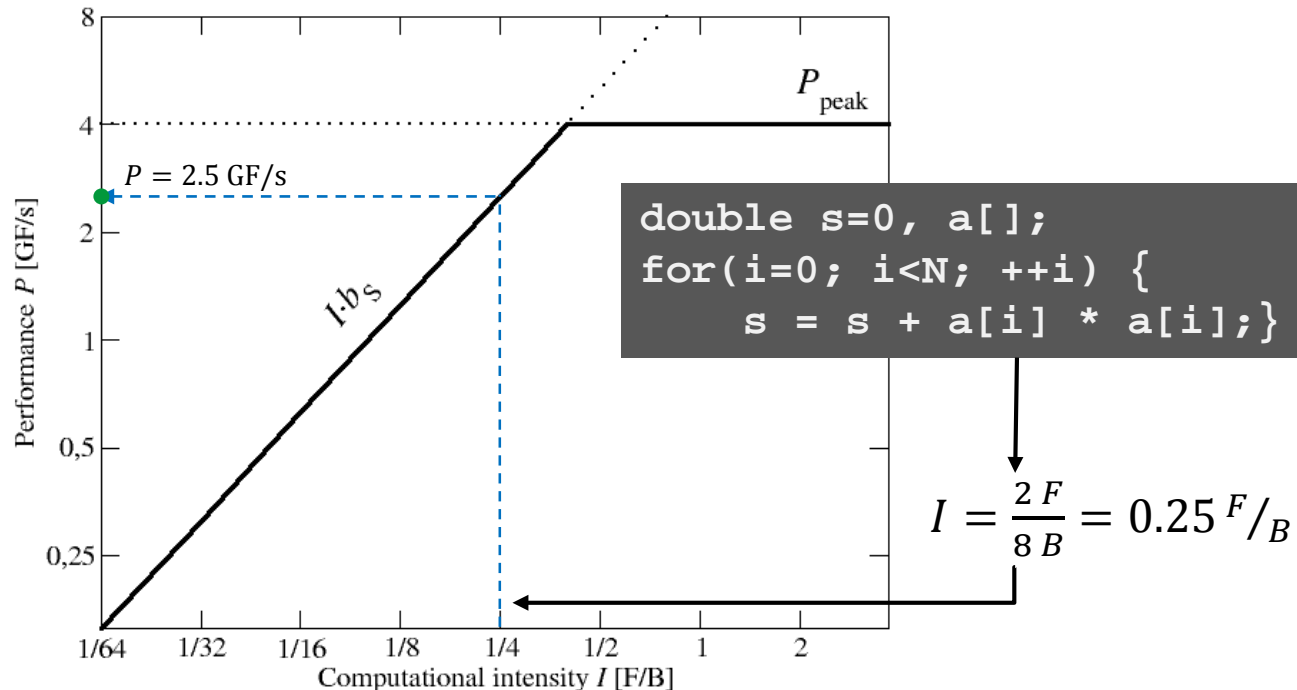## Apply the naive Roofline model in practice

- Machine parameter #1:  Peak performance:  $P_{peak}\ \left[\frac{F}{s}\right]$

- Machine parameter #2:  Memory bandwidth:  $b_S\ \left[\frac{B}{s}\right]$

- Code characteristic:  Computational intensity:  $I\ \left[\frac{F}{B}\right]$

Machine properties:

$$\boldsymbol{P_{peak}} = 4\,\frac{\text{GF}}{\text{s}}$$

$$\boldsymbol{b_S} = 10\,\frac{\text{GB}}{\text{s}}$$

Application property: $I$

```
double s=0, a[];
for(i=0; i<N; ++i) {
    s = s + a[i] * a[i];}
```

$$I = \frac{2\,F}{8\,B} = 0.25\ {}^{F}/_{B}$$

$P = 2.5$ GF/s

$P_{peak}$

$I \cdot b_S$

Performance $P$ [GF/s]

Computational intensity $I$ [F/B]

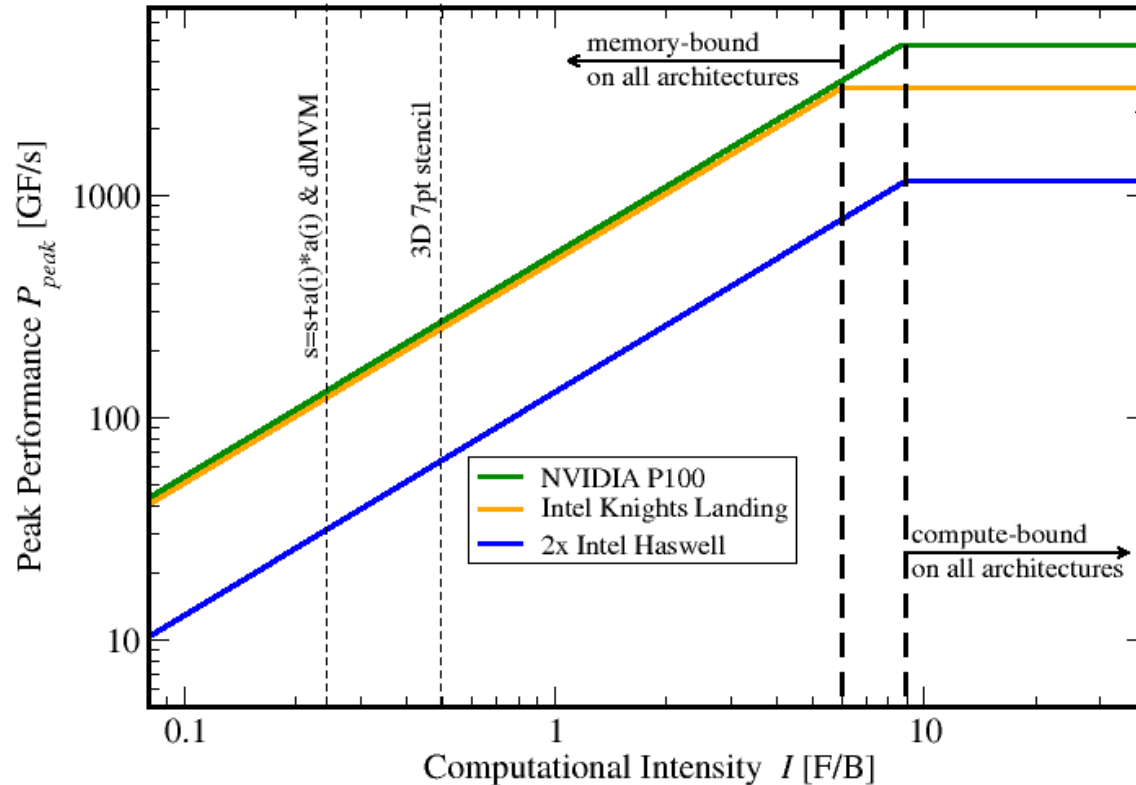# Prerequisites for the Roofline Model

- **The roofline formalism is based on some (crucial) prerequisites:**
  - There is a clear concept of "work" vs. "traffic"
    - "work" = flops, updates, iterations…
    - "traffic" = required data to do "work"

  - Machine input parameters: Peak Performance and Peak Bandwidth
    Application/kernel is expected to achieve is limits theoretically

- **Assumptions behind the model:**
  - Data transfer and core execution overlap perfectly!
    - **Either** the limit is core execution **or** it is data transfer
    - **Slowest limiting factor "wins";** all others are assumed to have no impact
  - Latency effects are ignored, i.e., perfect streaming mode
  - "Steady-state" code execution (no wind-up/-down effects)

Compare capabilities of different machines:



Assuming double precision –
for single precision:
$P_{peak} \rightarrow 2 \cdot P_{peak}$

- Roofline always provides upper bound – but is it realistic?
- If code is not able to reach this limit (e.g., contains add operations only), machine parameters need to be redefined (e.g., $P_{peak} \rightarrow P_{peak}/2$)

# A refined Roofline Model

1. *$P_{max}$ = Applicable peak performance* of a loop, assuming that data comes from the level 1 cache (this is not necessarily $P_{peak}$)
   → e.g., $P_{max}$ = 176 GFlop/s

2. *$I$ = Computational intensity ("work" per byte transferred)* over the slowest data path utilized (code balance $B_C = I^{-1}$)
   → e.g., $I$ = 0.167 Flop/Byte → $B_C$ = 6 Byte/Flop

3. *$b_S$ = Applicable (saturated) peak bandwidth* of the slowest data path utilized (measure attainable bandwidth using, e.g. STREAM)
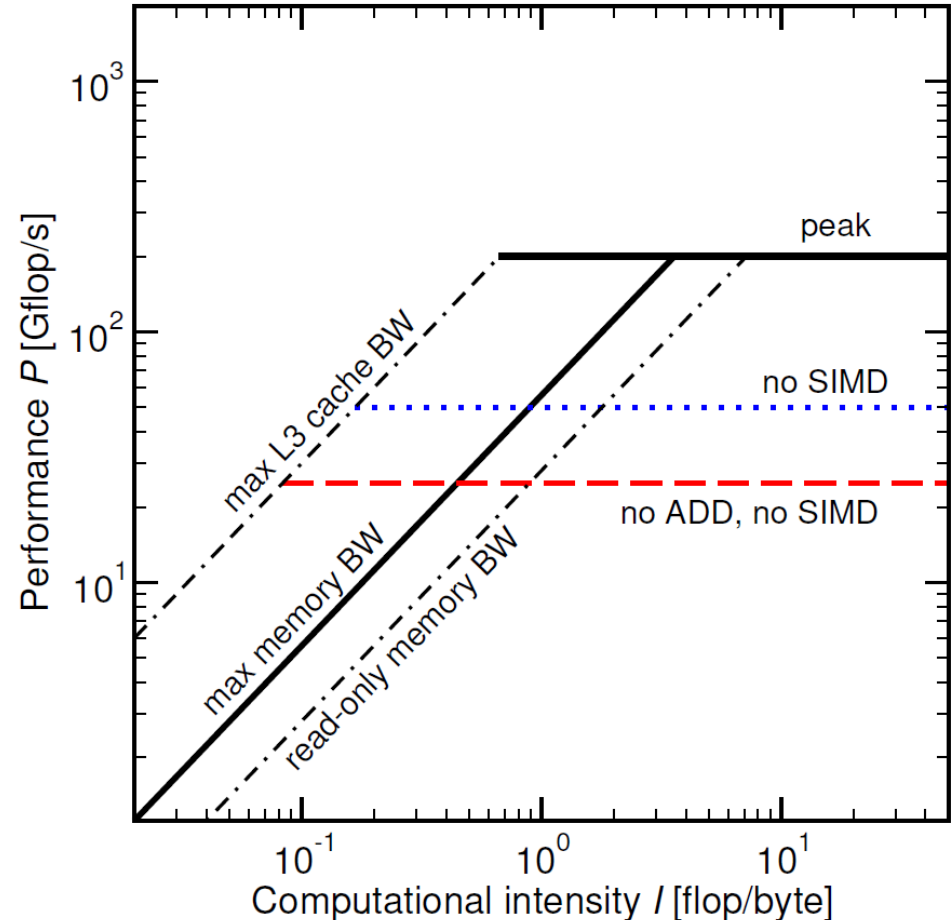   → e.g., $b_S$ = 56 GByte/s

Expected performance:

$$P = \min(P_{\max}, I \cdot b_S) = \min\left(P_{\max}, \frac{b_S}{B_C}\right)$$
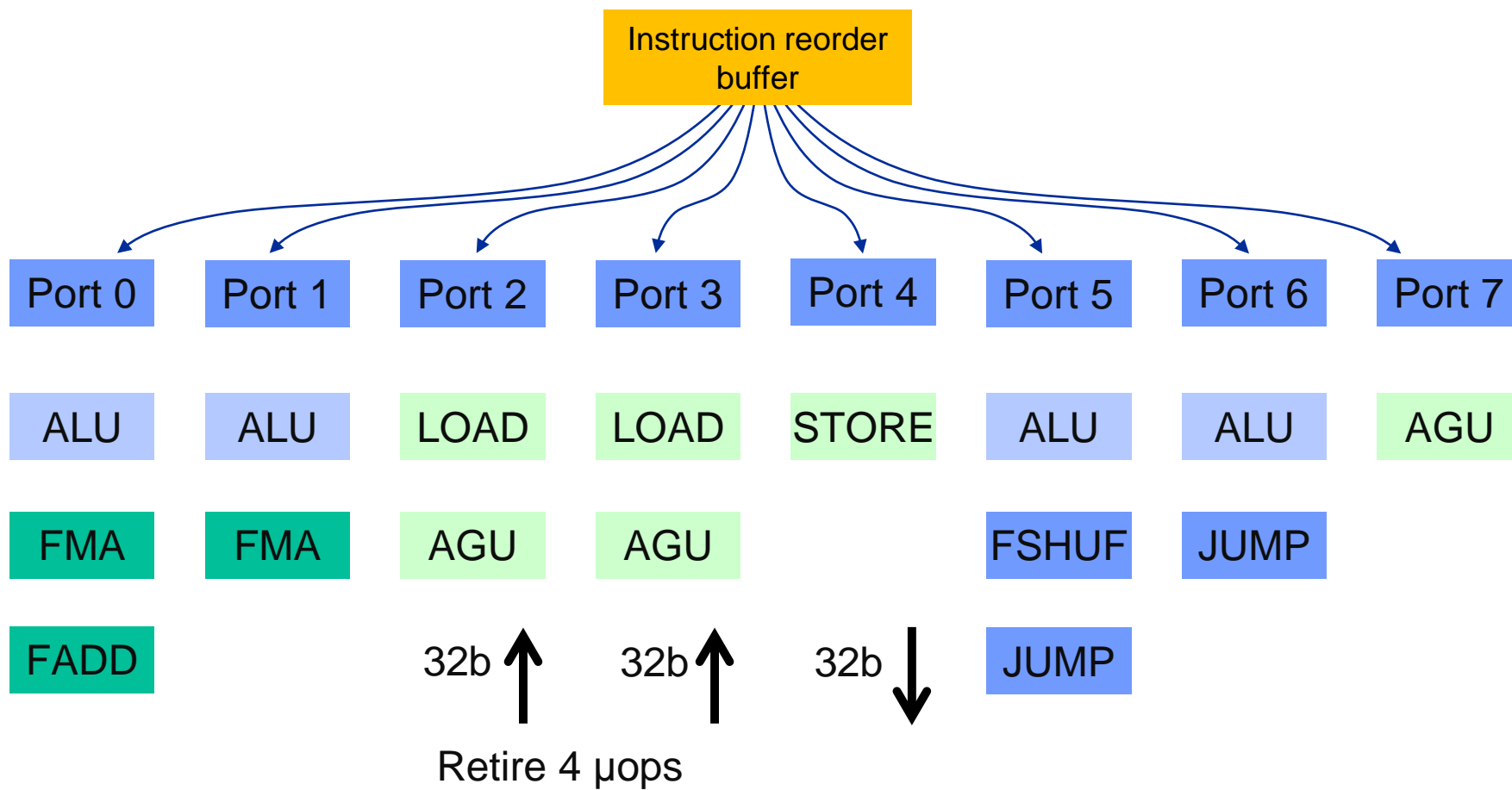
[Byte/s]

[Byte/Flop]

Multiple ceilings may apply

- Different bandwidths /data paths
  → different inclined ceilings

- Different $P_{max}$
  → different flat ceilings

  In fact, $P_{max}$ should always come from code analysis; generic ceilings are usually impossible to attain

Haswell/Broadwell port scheduler model:

| Instruction reorder buffer | | | | | | | |
|---|---|---|---|---|---|---|---|
| Port 0 | Port 1 | Port 2 | Port 3 | Port 4 | Port 5 | Port 6 | Port 7 |
| ALU | ALU | LOAD | LOAD | STORE | ALU | ALU | AGU |
| FMA | FMA | AGU | AGU | | FSHUF | JUMP | |
| FADD | | 32b ↑ | 32b ↑ | 32b ↓ | JUMP | | |

Retire 4 µops

Haswell/Broadwell

```
double  *A, *B, *C, *D;
for (int i=0; i<N; i++) {
    A[i] = B[i] + C[i] * D[i];
}
```

Minimum number of cycles to process **one AVX-vectorized iteration** (equivalent to 4 scalar iterations) on one core?

→ Assuming full throughput:

Cycle 1:  LOAD + LOAD + STORE
Cycle 2:  LOAD + LOAD + FMA + FMA
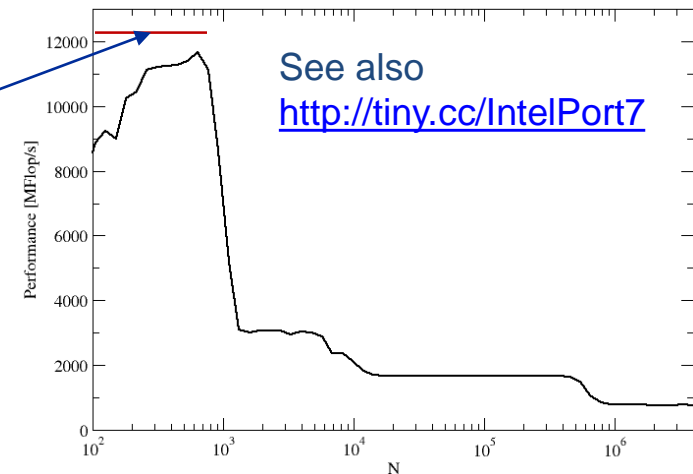Cycle 3:  LOAD + LOAD + STORE          **Answer:  1.5 cycles**

```
double  *A, *B, *C, *D;
for (int i=0; i<N; i++) {
   A[i] = B[i] + C[i] * D[i];
}
```

What is the **performance in GFlops/s per core** and the bandwidth in GBytes/s?

One AVX iteration (1.5 cycles) does 4 x 2 = 8 flops:

$$2.3 \cdot 10^9 \text{ cy/s} \cdot \frac{8 \text{ flops}}{1.5 \text{ cy}} = \mathbf{12.27} \frac{\textbf{Gflops}}{\textbf{s}}$$

$$12.27 \frac{\text{Gflops}}{s} \cdot 16 \frac{\text{bytes}}{\text{flop}} = 196 \frac{\text{Gbyte}}{s}$$



See also
http://tiny.cc/IntelPort7

**Vector triad `A(:)=B(:)+C(:)*D(:)` on a 2.3 GHz 14-core Haswell chip**

Consider full chip (14 cores):

Memory bandwidth: $b_S$ = **50 GB/s**

Code balance (incl. write allocate):
$B_c$ = (4+1) Words / 2 Flops = 20 B/F → $I$ **= 0.05 F/B**

→ $I \cdot b_S$ **= 2.5 GF/s** (0.5% of peak performance)

$P_{peak}$ / core = 36.8 Gflop/s ((8+8) Flops/cy x 2.3 GHz)
$P_{max}$ / core = 12.27 Gflop/s (see prev. slide)

→ $P_{max}$ **= 14 \* 12.27 Gflop/s =172 Gflop/s** (33% peak)

$$P = \min(P_{max}, I \cdot b_S) = \min(172, 2.5) \, \text{GFlop/s} = \textcolor{red}{2.5 \, \text{GFlop/s}}$$
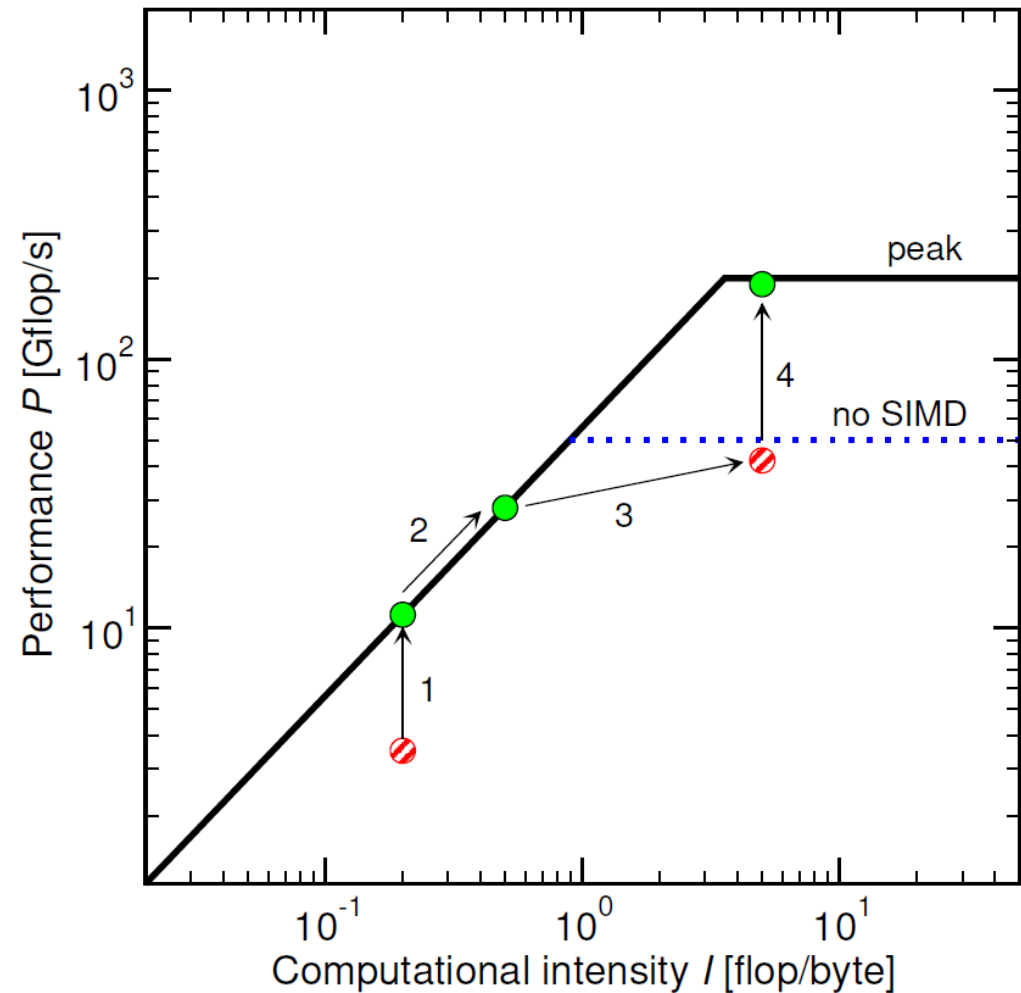
# A not so simple Roofline example

**Example:**    `do i=1,N; s=s+a(i); enddo`

in single precision on an 8-core 2.2 GHz Sandy Bridge socket @ "large" N

$$P = \min(P_{\max}, I \cdot b_S)$$



Machine peak (ADD+MULT) Out of reach for this code — 282 GF/s

ADD peak (best possible code) — 141 GF/s

no SIMD — 17.6 GF/s

3-cycle latency per ADD if not unrolled — 5.9 GF/s

See architecture intro

$b_s$ = 40 GB/s

**P** (better loop code)

**P** (worst loop code)

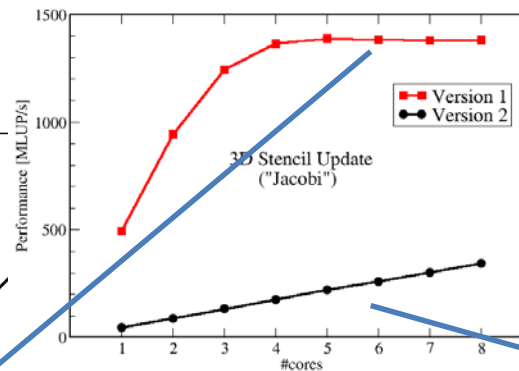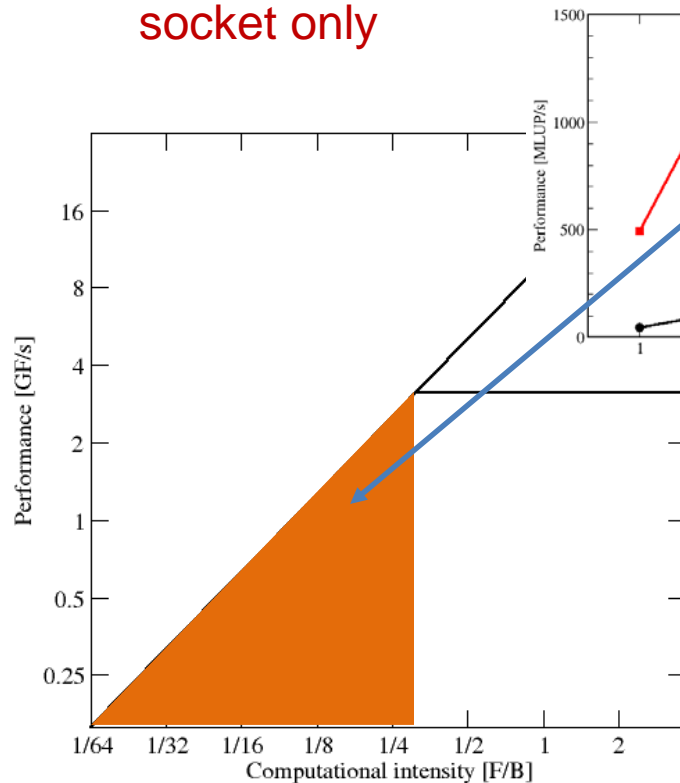$I$ = 1 flop / 4 byte (SP!)

1. **Hit the BW bottleneck by good serial code**
   (e.g., Ninja C++ → Fortran)

2. **Increase intensity to make better use of BW bottleneck**
   (e.g., spatial loop blocking [see later])

3. **Increase intensity and go from memory bound to core bound**
   (e.g., temporal blocking)

4. **Hit the core bottleneck by good serial code**
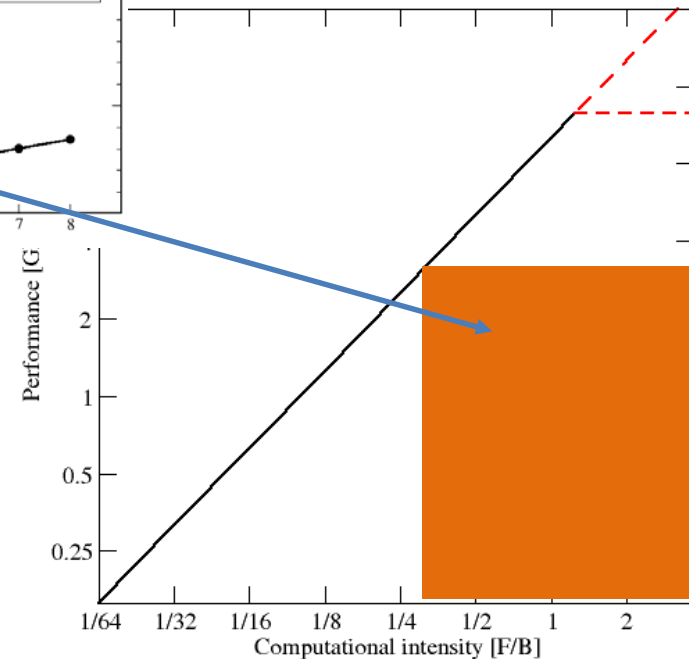   (e.g., `-fno-alias` [see later])

## Bandwidth-bound (simple case)

1. Accurate traffic calculation (write-allocate, strided access, …)

2. Practical ≠ theoretical BW limits

3. Saturation effects → consider full socket only

## Core-bound (may be complex)

1. **Multiple bottlenecks**: LD/ST, arithmetic, pipelines, SIMD, execution ports

2. Limit is linear in # of cores

# Shortcomings of the roofline model

- **Saturation effects in multicore chips are not explained**
  - Reason: "saturation assumption"
  - Cache line transfers and core execution do sometimes not overlap perfectly
  - It is not sufficient to measure single-core STREAM to make it work
  - Only increased "pressure" on the memory interface can saturate the bus
    → need more cores!

- **In-cache performance is not correctly predicted**


- **The ECM performance model gives more insight:**

  G. Hager, J. Treibig, J. Habich, and G. Wellein: Exploring performance and power properties of modern multicore chips via simple machine models. Concurrency and Computation: Practice and Experience (2013).
  DOI: 10.1002/cpe.3180 Preprint: arXiv:1208.2908

`A(:)=B(:)+C(:)*D(:)`