# How to Write Fast Numerical Code

Fall 2016
*Lecture:* Roofline model

**Instructor:** Torsten Hoefler & Markus Püschel
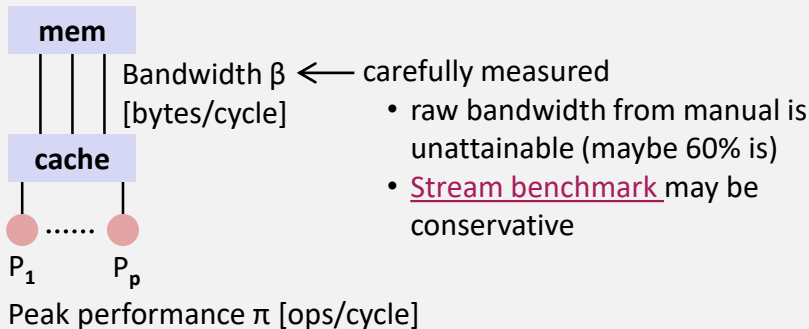
**TA:** Salvatore Di Girolamo

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Roofline model (Williams et al. 2008)

Resources in a processor that bound performance:
- peak performance [flops/cycle]
- memory bandwidth [bytes/cycle]
- <others>

**Platform model**



**mem**

Bandwidth β ← carefully measured
[bytes/cycle]
- raw bandwidth from manual is unattainable (maybe 60% is)
- Stream benchmark may be conservative

**cache**

P$_1$ ...... P$_p$

Peak performance π [ops/cycle]

**Algorithm model (n is the input size)**

Operational intensity I(n) = W(n)/Q(n) =

$$\frac{\text{number of flops (cost)}}{\text{number of bytes transferred between memory and cache}} \quad [\text{ops/bytes}]$$

Q(n): assumes empty cache;
best measured with performance counters

**Notes**
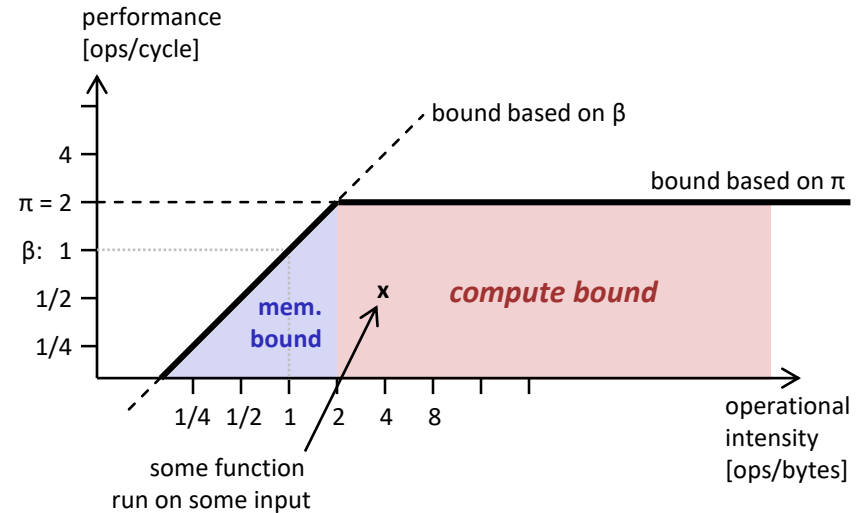In general, Q and hence W/Q depend on the cache size m [bytes].
For some functions the optimal achievable W/Q is known:
    FFT/sorting: Θ(log(m))
    Matrix multiplication: Θ(sqrt(m))

**Roofline model**
Example: one core with π = 2 and β = 1 and no SSE
ops are double precision flops



**Bound based on β?**
- assume program as operational intensity of x ops/byte
- it can get only β bytes/cycle
- hence: performance = y ≤ βx
- in log scale: $\log_2(y) \leq \log_2(\beta) + \log_2(x)$
- line with slope 1; y = β for x = 1

**Variations**
- vector instructions: peak bound goes up (e.g., 4 times for AVX)
- multiple cores: peak bound goes up (p times for p cores)
- program has uneven mix adds/mults: peak bound comes down (note: now this bound is program specific)
- accesses with little spatial locality: operational intensity decreases (because entire cache blocks are loaded)

# Roofline Measurements

- **Tool developed in our group**
  *(G. Ofenbeck, R. Steinmann, V. Caparros-Cabezas, D. Spampinato)*
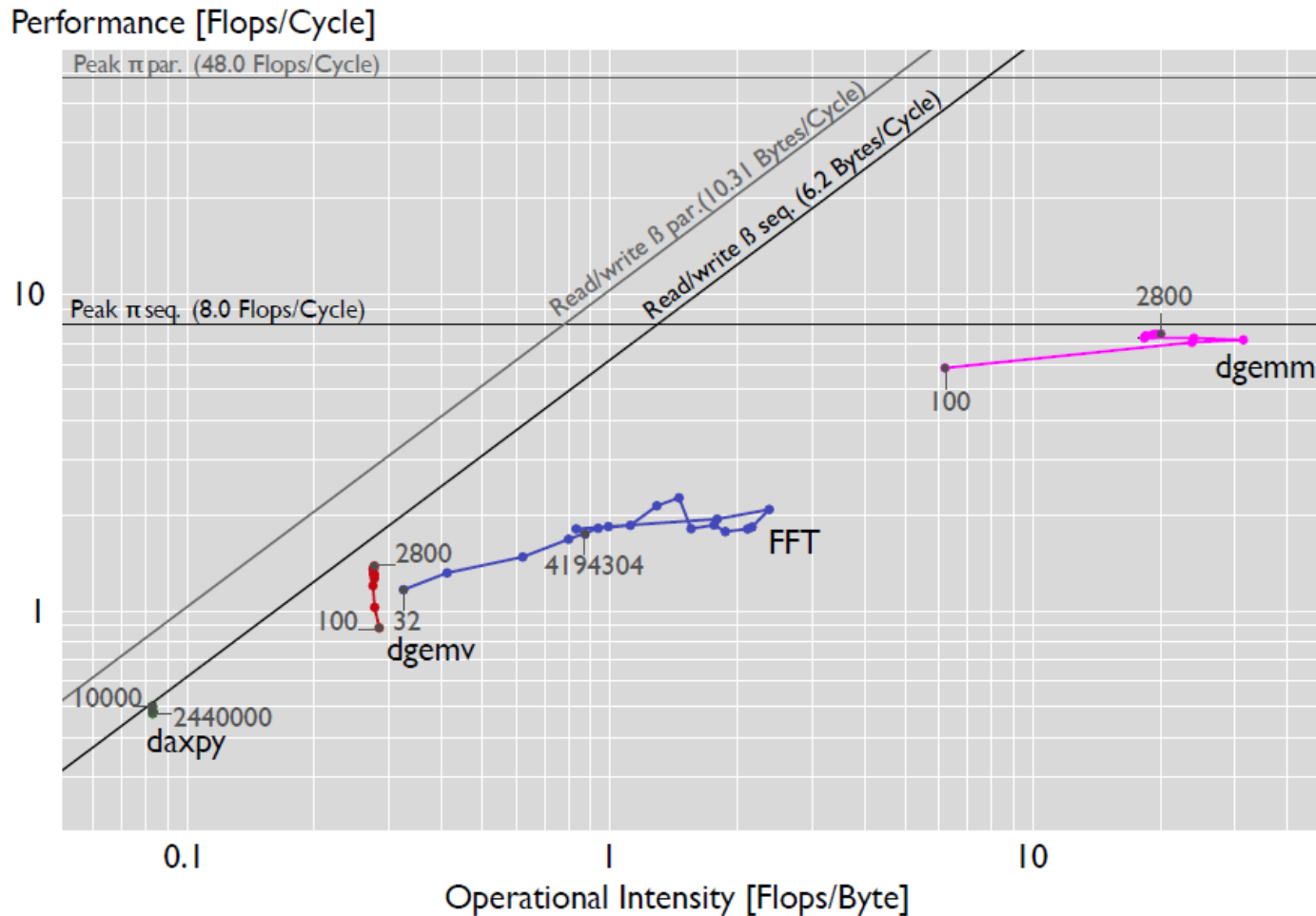  *http://www.spiral.net/software/roofline.html*

- **Example plots follow**

- **Get (non-asymptotic) bounds on I:**
  - daxpy:     $y = \alpha x + y$
  - dgemv:     $y = Ax + y$
  - dgemm:     $C = AB + C$
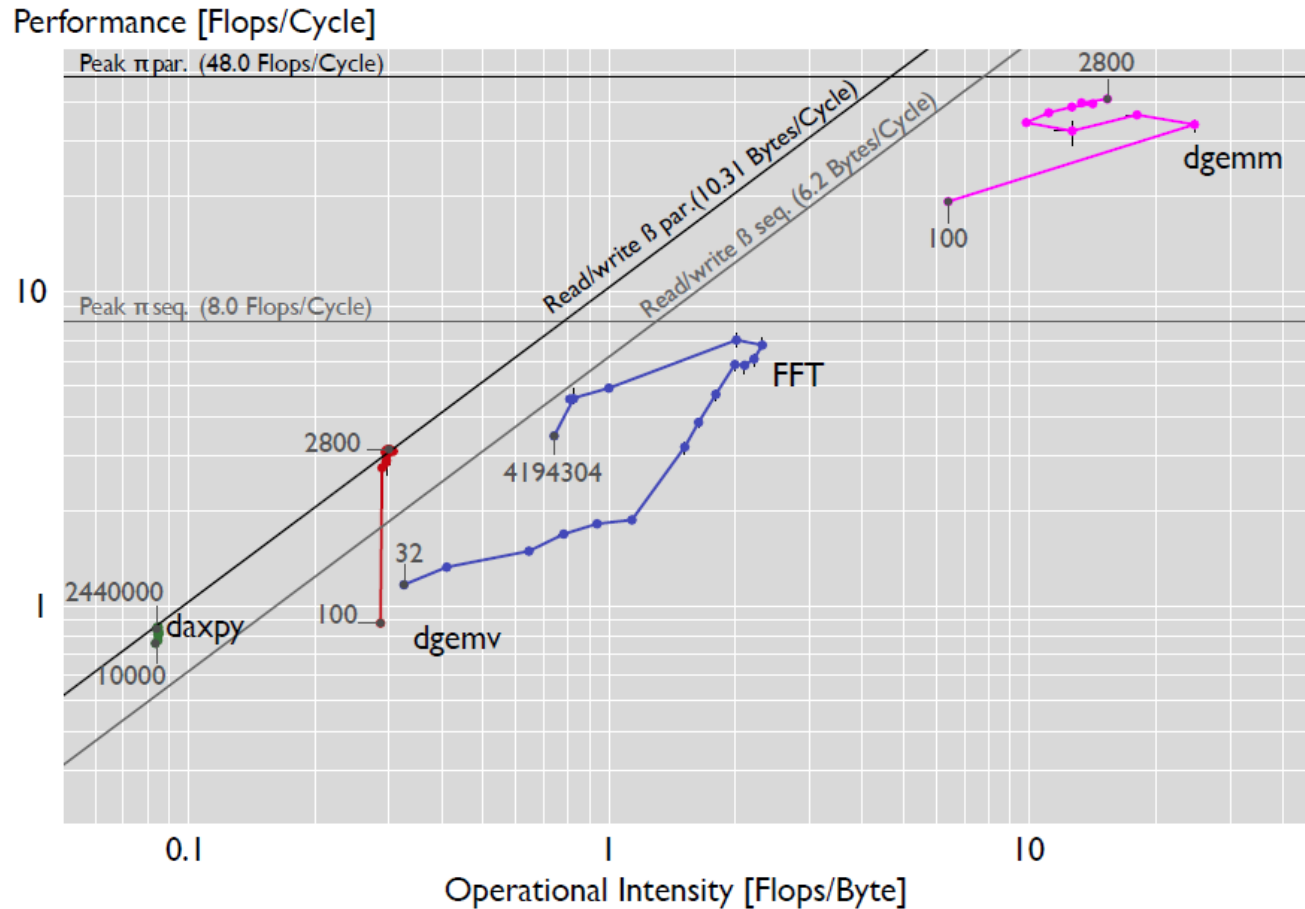  - FFT

# Roofline Measurements

**Performance [Flops/Cycle]**

Peak π par. (48.0 Flops/Cycle)

Read/write β par.(10.31 Bytes/Cycle)

Read/write β seq. (6.2 Bytes/Cycle)

Peak π seq. (8.0 Flops/Cycle)

2800

dgemm

100

FFT

2800

4194304

100

32

dgemv

10000

2440000

daxpy

10

1

0.1

1

10

**Operational Intensity [Flops/Byte]**

*What happens when we go to parallel code?*

# Roofline Measurements

*What happens when we go to warm cache?*

# Roofline Measurements

# Roofline Measurements

**MMM: Try to guess the basic shapes**

# Summary

- **Roofline plots distinguish between memory and compute bound**

- **Can be used on paper**

- **Measurements difficult (performance counters) but doable**

- **Interesting insights: *use in your project!***

# References

- Samuel Williams, Andrew Waterman, David Patterson
  **Roofline: an insightful visual performance model for multicore architectures**
  Communications ACM 55(6): 121-130 (2012)

- Georg Ofenbeck, Ruedi Steinmann, Victoria Caparros, Daniele G. Spampinato and Markus Püschel
  **Applying the Roofline Model**
  Proc. IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2014, pp. 76-85

- Victoria Caparros and Markus Püschel
  **Extending the Roofline Model: Bottleneck Analysis with Microarchitectural Constraints**
  Proc. IEEE International Symposium on Workload Characterization (IISWC), pp. 222-231, 2014