# Design of Parallel and High-Performance Computing

Fall 2013
*Lecture:* Roofline

**Instructor:** Torsten Hoefler & Markus Püschel

**TA:** Timo Schneider

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

# Operational Intensity

- **Definition: Given a program P, assume cold (empty) cache**

  *Operational intensity:* $I(n) = \dfrac{W(n)}{Q(n)}$

  #flops (input size n)

  #bytes transferred cache $\leftrightarrow$ memory (for input size n)

- **Examples: Determine asymptotic bounds on I(n)**
  - Vector sum: y = x + y          **O(1)**
  - Matrix-vector product: y = Ax   **O(1)**
  - Fast Fourier transform          **O(log(n))**
  - Matrix-matrix product: C = AB + C   **O(n)**
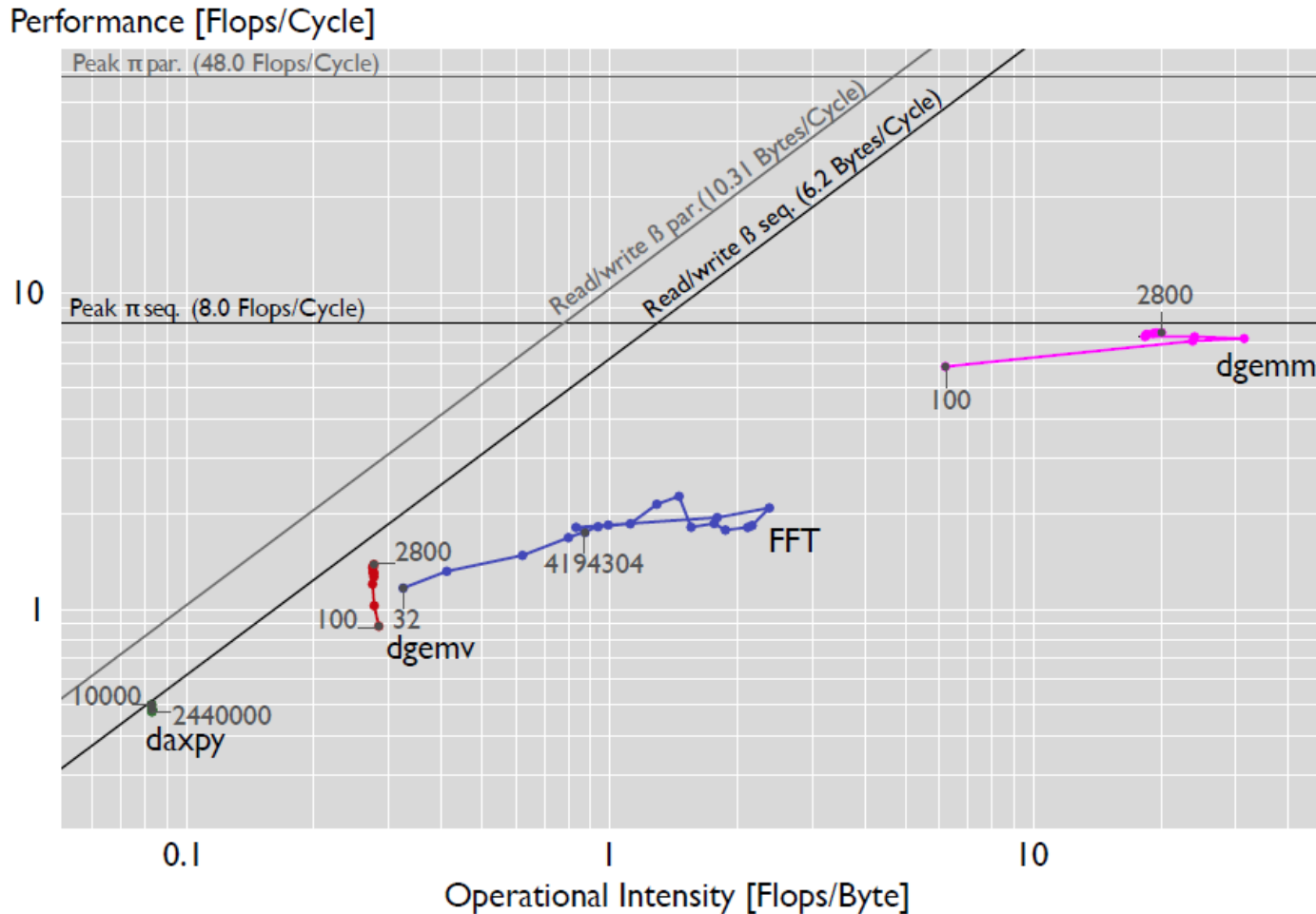
# Example MVM: y = Ax + y

- **Number of flops?**

- **Number of compulsory misses (cold cache)?**

- **Upper bound on the operational intensity?**

# Roofline Measurements

- **Tool developed in our group**
  *(G. Ofenbeck, R. Steinmann, V. Caparros-Cabezas, D. Spampinato)*

- **Example plots follow**

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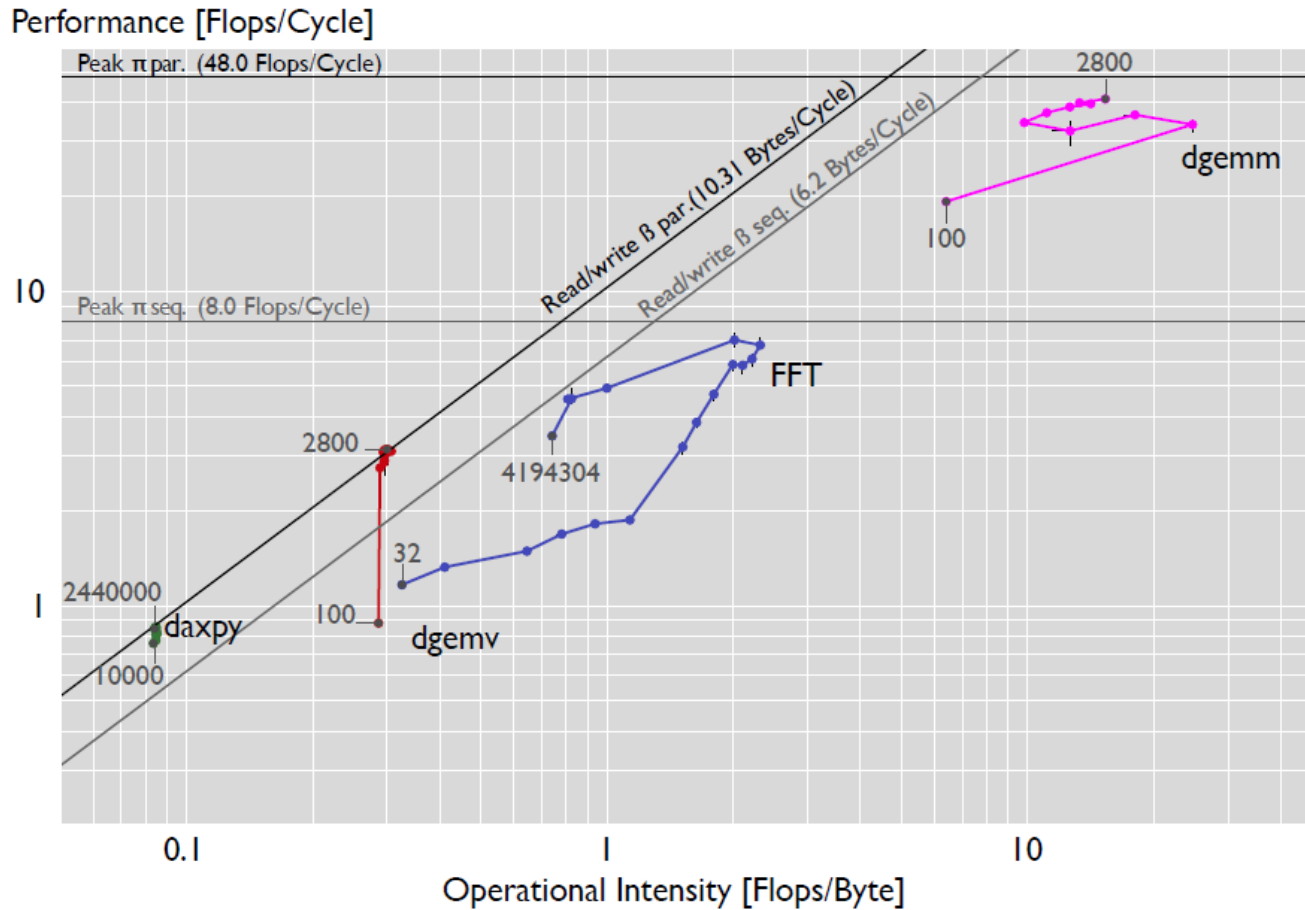  - FFT

# Roofline Measurements

**What happens when we go to parallel code?**

# 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)

2800

dgemm

100

10

Peak π seq. (8.0 Flops/Cycle)

FFT

2800

4194304

32

1

2440000

100

daxpy

dgemv

10000

0.1          1          10

*Operational Intensity [Flops/Byte]*

***What happens when we go to warm cache?***

6

# Roofline Measurements

# Summary

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

- **Can be used on paper**

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