# On the limits of (and opportunities for?) GPU acceleration

Aparna Chandramowlishwaran, Jee Choi, Kenneth Czechowski, Murat (Efe) Guney, Logan Moon, Aashay Shringarpure, Richard (Rich) Vuduc

HotPar'10, Berkeley — June 15, 2010

**Georgia Tech** | College of Computing
Computational Science and Engineering

**hpcgarage**

# Twelve Ways to Fool the Masses When Giving Performance Results on Parallel Computers
## David H. Bailey
## June 11, 1991
## Ref: Supercomputing Review, Aug. 1991, pg. 54--55

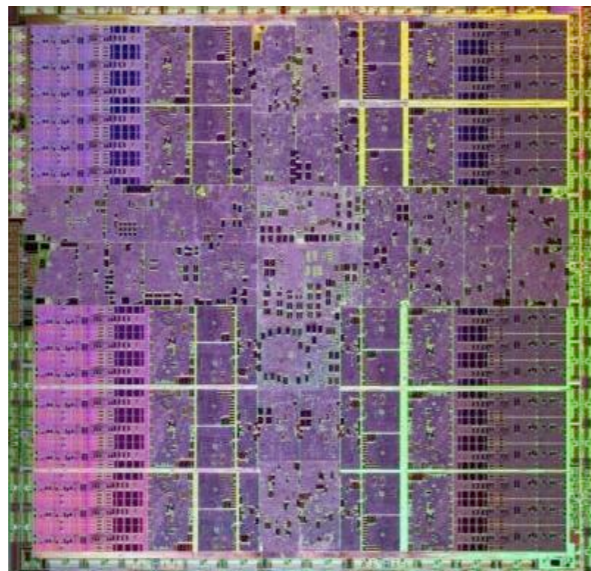6. Compare your results against scalar, unoptimized code on Crays.

It really impresses the audience when you can state that your code runs several times faster than a Cray, currently the world's dominant supercomputer. Unfortunately, with a little tuning many applications run quite fast on Crays. Therefore you must be careful not to do any tuning on the Cray code. Do not insert vectorization directives, and if you find any, remove them. In extreme cases it may be necessary to disable all vectorization with a command line flag. Also, Crays often run much slower with bank conflicts, so be sure that your Cray code accesses data with large, power-of-two strides whenever possible. It is also important to avoid multitasking and autotasking on Crays --- imply in your paper that the one processor Cray performance rates you are comparing against represent the full potential of a $25 million Cray system.
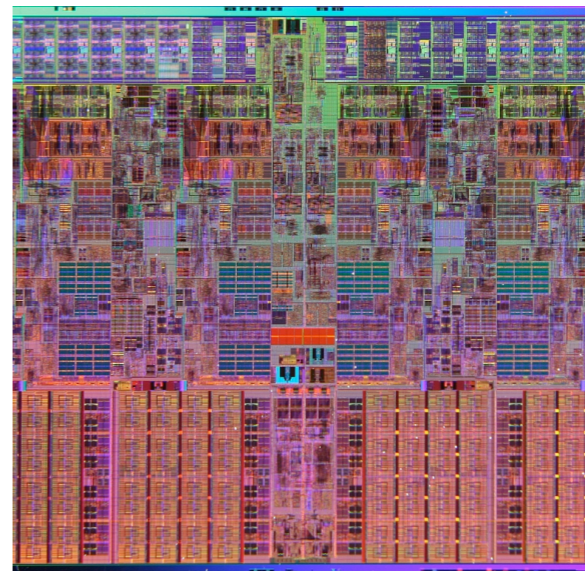
# Q: Port to GPU?

(Posed to me by Scott Klasky at ORNL)

▶ A: Given roughly same level of tuning & power*, …



GPU ≅ CPU

▶ Meta-analysis, for **semi-irregular** sci. comp. + data analytics apps (sparse iterative + direct solvers; tree-based particle methods)

# Summary of "limits"

▶ Bandwidth-bound: Aggregate bandwidth < 3x, restricted access patterns

▶ Compute-bound: 5 to 10x peak potential, but there is a multithreading granularity mismatch

▶ PCIe: Will architects replicate GPU memory system in on-die CPU/GPU?

▶ Productivity: To 0th order, tuning required on all platforms
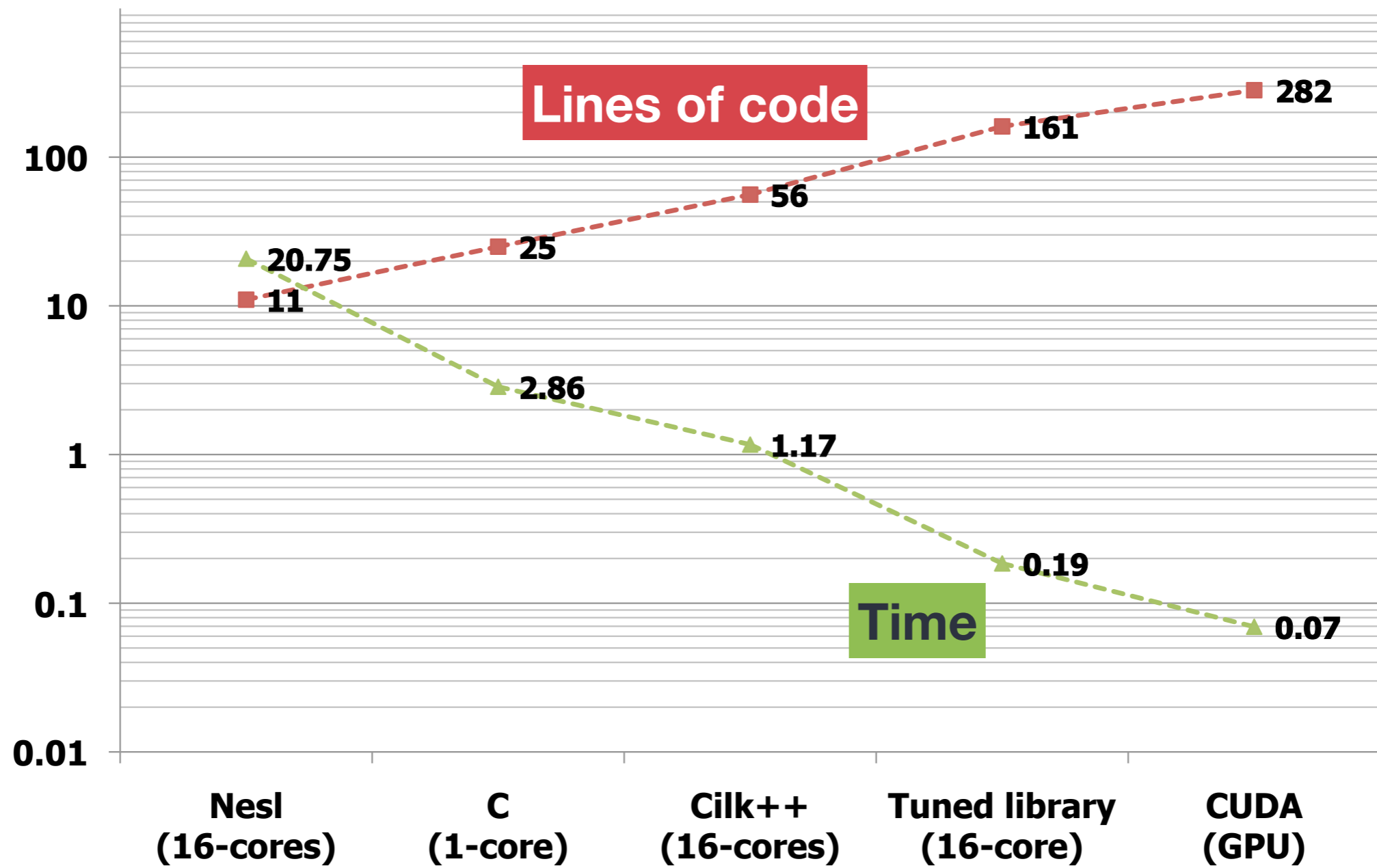I.e., Bailey's Rule #6

# Limitations of this meta-analysis

▶ Mix of partial results, very rough back-of-the envelope, apples vs. oranges

   ▶ But: "It's all fruit" — *My Big Fat Greek Wedding*

▶ Narrow: Scientific computing? (yawn)

   ▶ But: More physically realistic games & graphics, signal analysis, data analytics

▶ Limited to today's platforms, and excludes Fermi. What about tomorrow?

   ▶ But: "Prediction is very difficult, especially if it's about the future." — N. Bohr

# Other voices

▶ Vishkin, et al., HotPar'10 poster!

▶ Bordawekar, et al. (IBM). "Believe it or Not! Multi-core CPUs Can Match GPU Performance for FLOP-intensive Application!" IBM Technical Report RC24982, April 2010.

▶ Lee, et al. (Intel). "Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU." ISCA, June 2010.
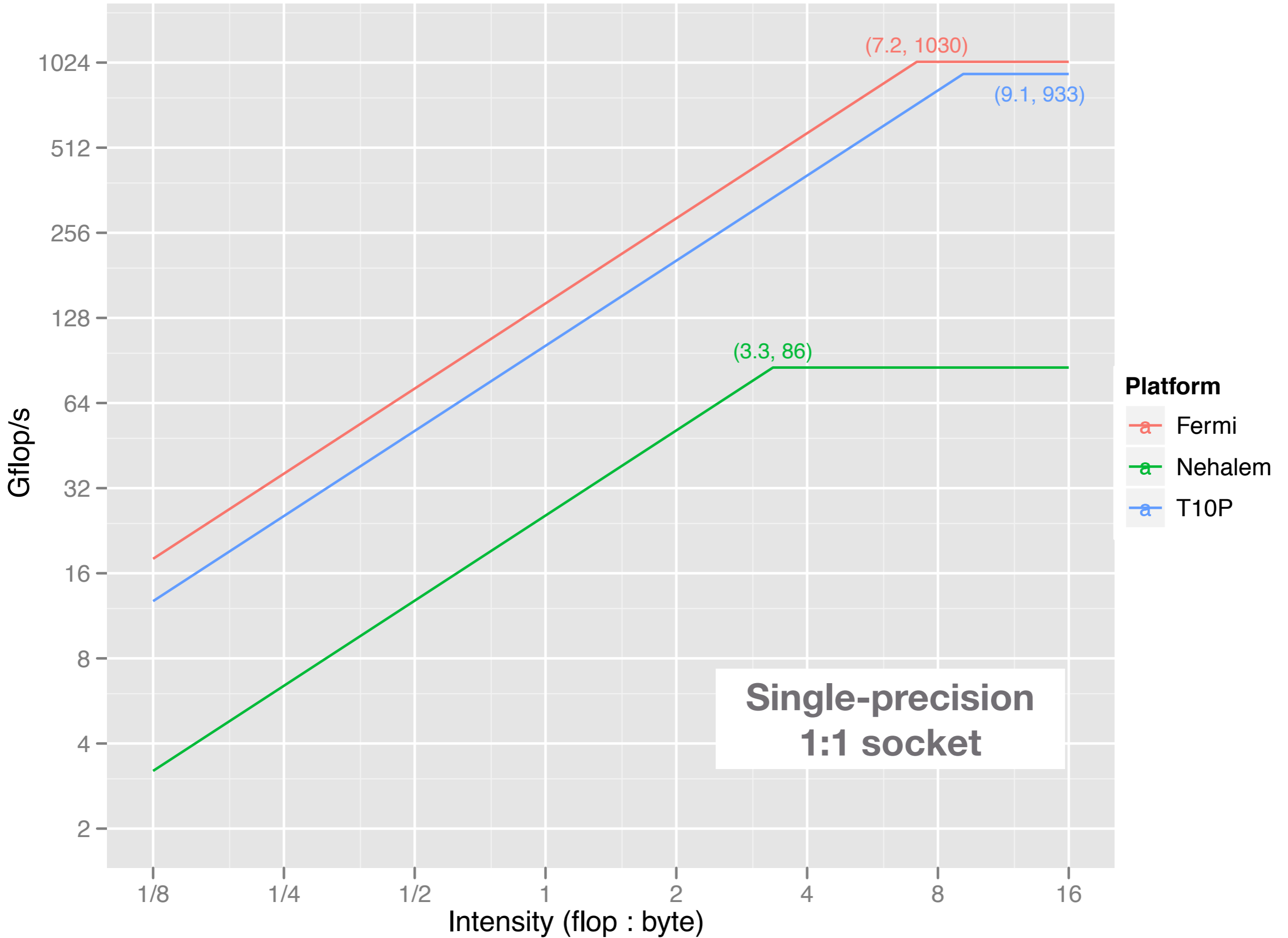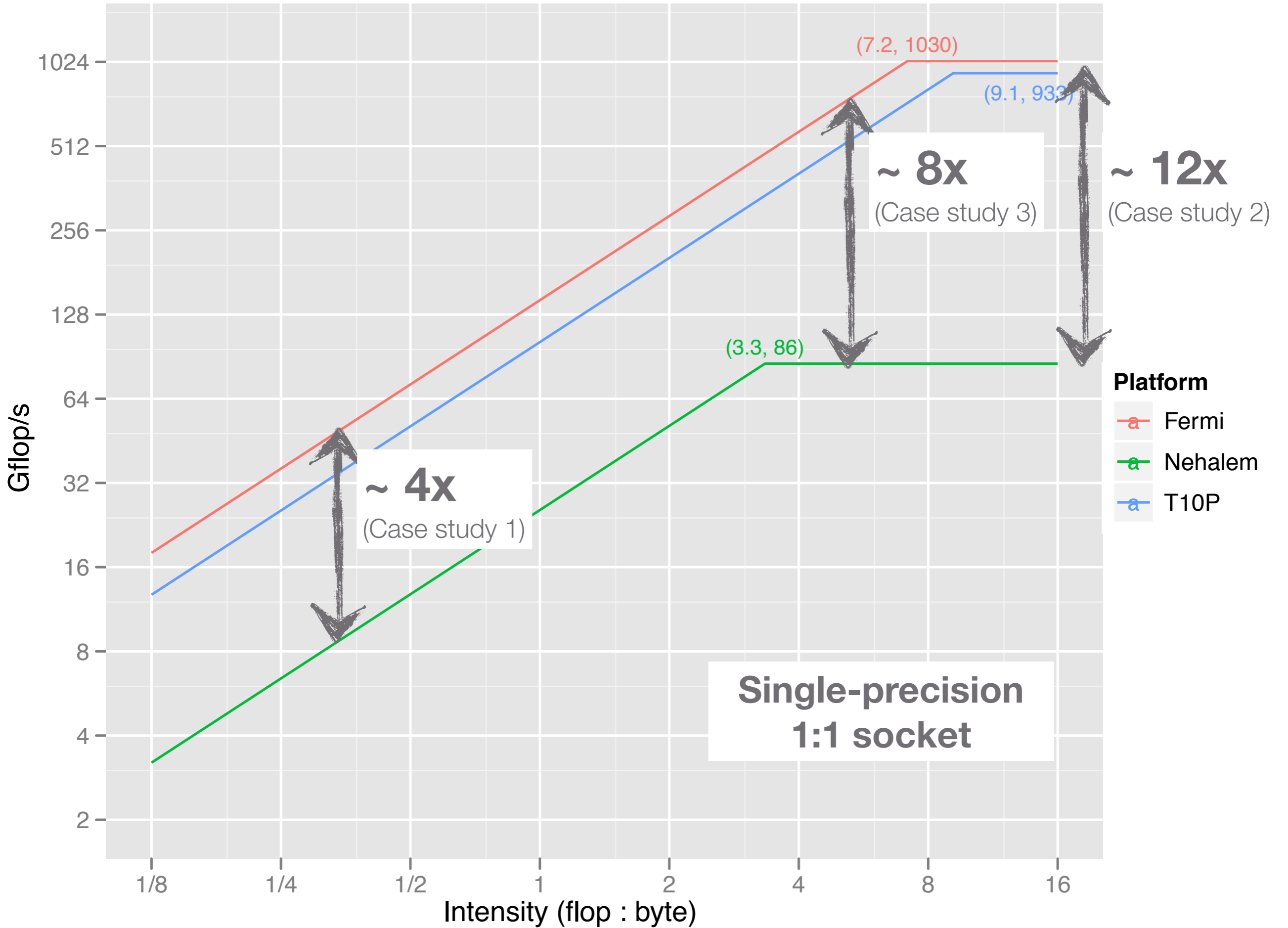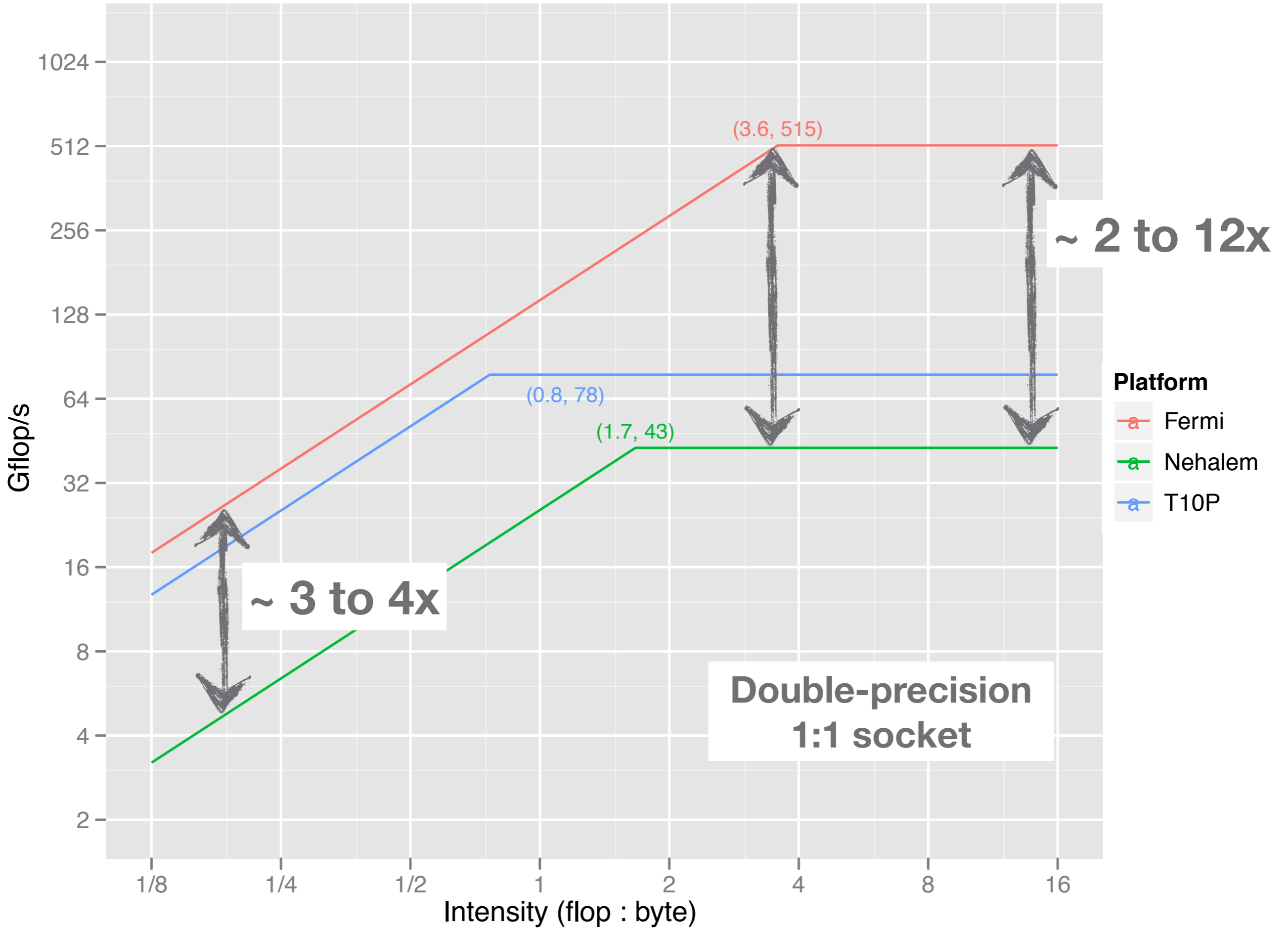
▶ Performance and productivity expectations

Thursday, July 1, 2010

Worth the effort?

Algorithm: Parallel sorting

Single-precision
1:1 socket

Thursday, July 1, 2010

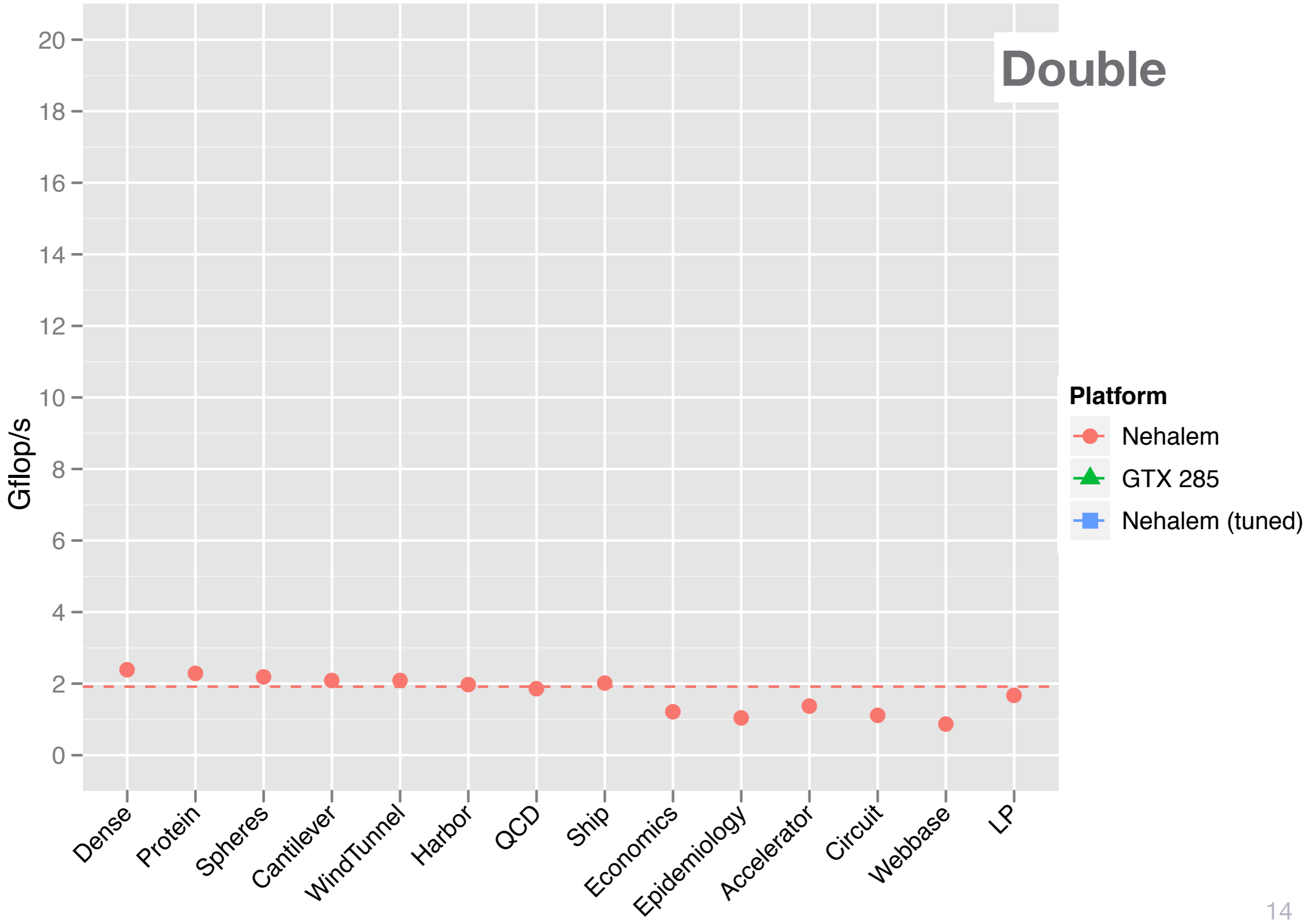▶ Case study 1: Sparse iterative solvers

S. Williams, N. Bell, J. Choi, M. Garland, L. Oliker, R. Vuduc. "SpMV on MC & Acc." In *BDK book* (2010).
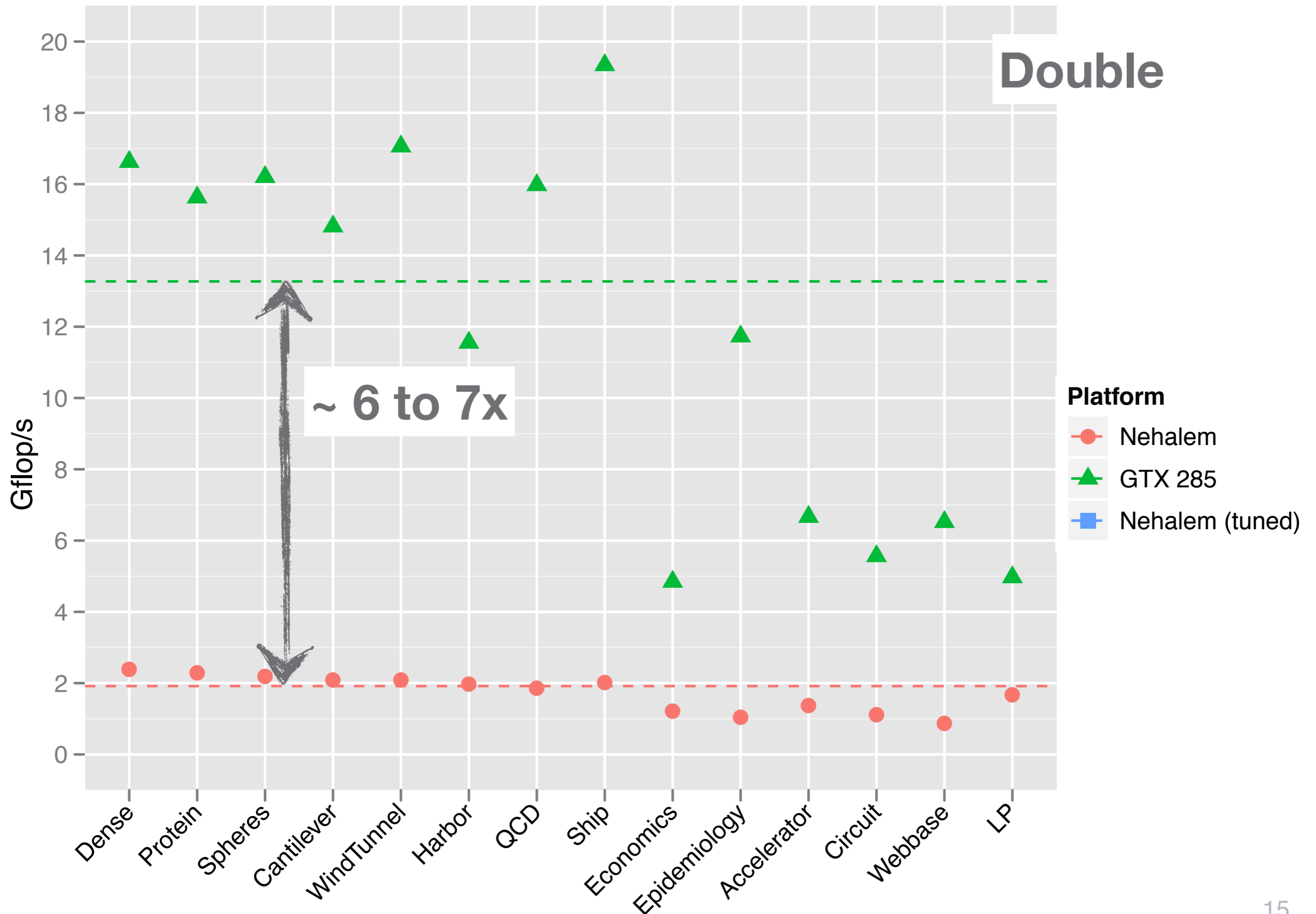
# Anatomy of a sparse iterative solver

```
do {
  …
  y ← A*x      "SpMV"
  …
} while (!converged)
```

▶ Bottleneck: Sparse matrix-vector multiply (SpMV)

▶ Memory bandwidth-limited (stream A): GPU / CPU ~ 3x
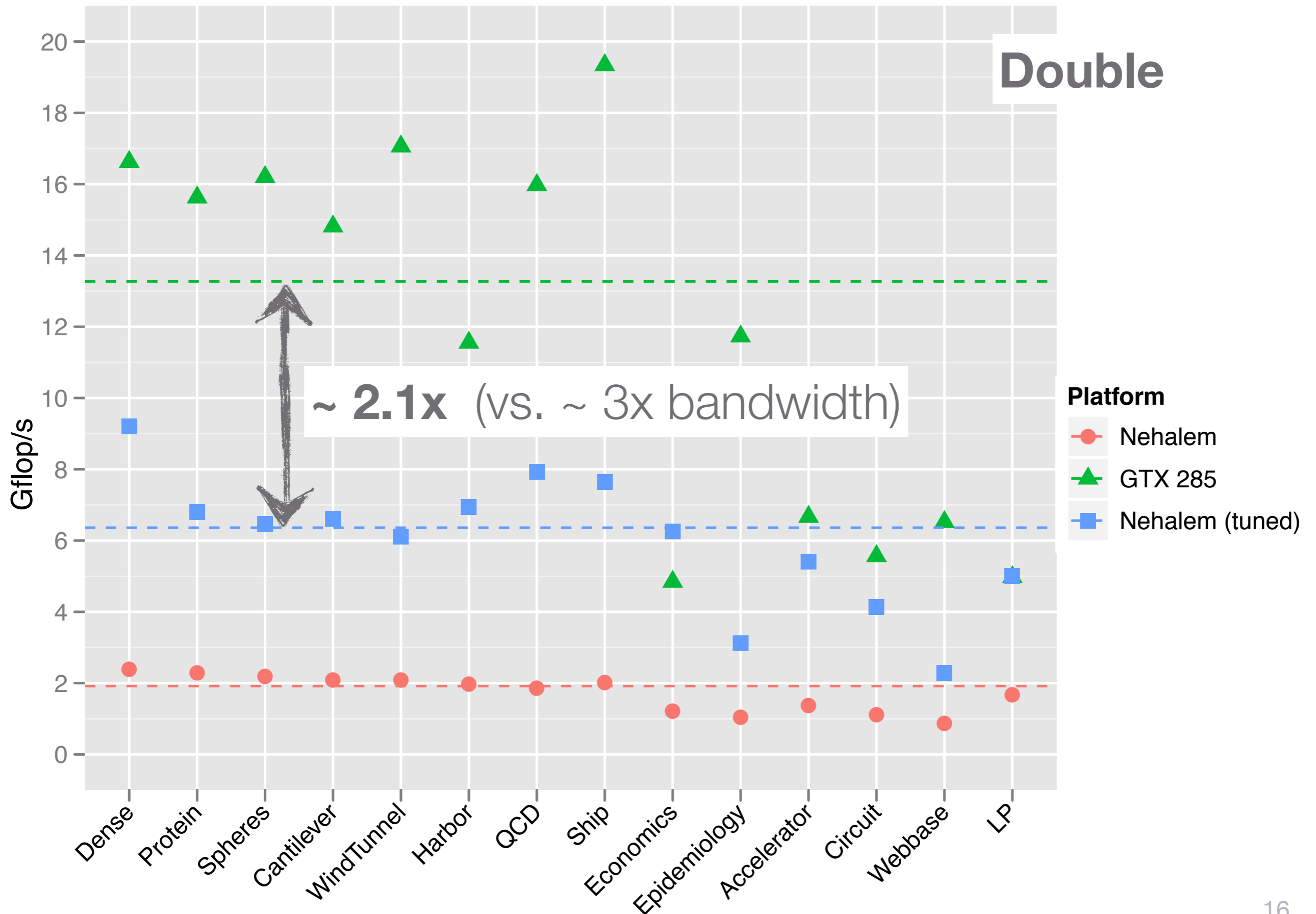
S. Williams, N. Bell, J. Choi, M. Garland, L. Oliker, R. Vuduc. "SpMV on MC & Acc." In *BDK book* (2010).

Thursday, July 1, 2010

**Double**

S. Williams, N. Bell, J. Choi, M. Garland, L. Oliker, R. Vuduc. "SpMV on MC & Acc." In *BDK book* (2010).

**Double**

~ **2.1x**  (vs. ~ 3x bandwidth)

Gflop/s

Platform
- Nehalem
- GTX 285
- Nehalem (tuned)

Dense · Protein · Spheres · Cantilever · WindTunnel · Harbor · QCD · Ship · Economics · Epidemiology · Accelerator · Circuit · Webbase · LP

S. Williams, N. Bell, J. Choi, M. Garland, L. Oliker, R. Vuduc. "SpMV on MC & Acc." In *BDK book* (2010).

16

# Beyond the kernel…

▶ Optimal data structures differ between CPU & GPU
  ⇒ Startup, transfer cost

▶ In distributed memory, need to transfer vectors
  ⇒ PCIe limits

▶ Need ~ 100 iterations to break even, ~ 840 to get 2x on actual solver
  ⇒ Big or not? Depends on app & input

▶ Case study 2: The Fast Multipole Method

(1) Lashuk, et al., SC'09.  (2) Chandramowlishwaran, Williams, et al., IPDPS'10.

# Blood cell simulation

(movie) In 2-D with deformations

19

# Anatomy of an FMM



**Want**:
All pair interactions among green dots
$\rightarrow$ O($n^2$)

**FMM idea**:
Build tree, traverse & prune (approx.)
$\rightarrow$ O($n$ log $n$), with accuracy guarantee

▶ Bottleneck: "Little" all-pair interactions ($b^2$) among leaves
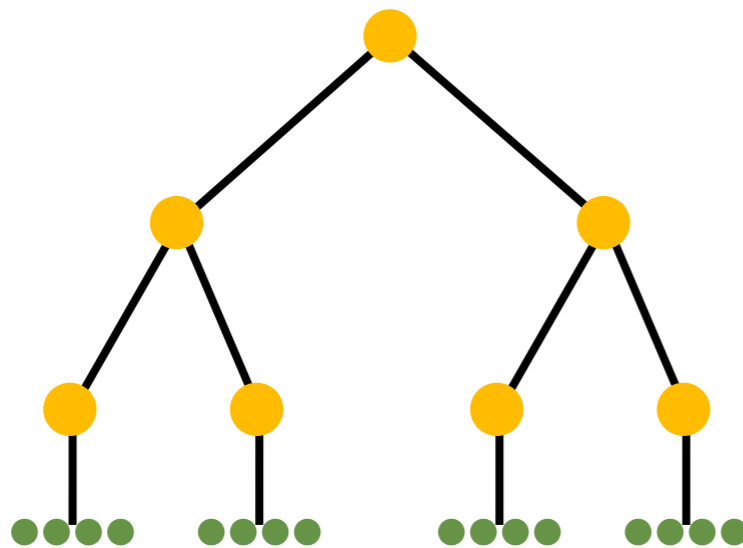
▶ High compute intensity: 12x possible?

Thursday, July 1, 2010

# The story

- ▶ Baseline: Lashuk, et al., SC'09

  - ▶ FMM that scales to 100k cores of "Kraken" machine @ UTK

  - ▶ Good overall scalability, but low within-node performance

- ▶ Try GPUs?

  - ▶ Gumerov & Duraiswami (JCP'08) suggest 30—60x speedups on GPUs

  - ▶ We successfully replicate on 256 GPU system at UIUC (1 MPI task / GPU)

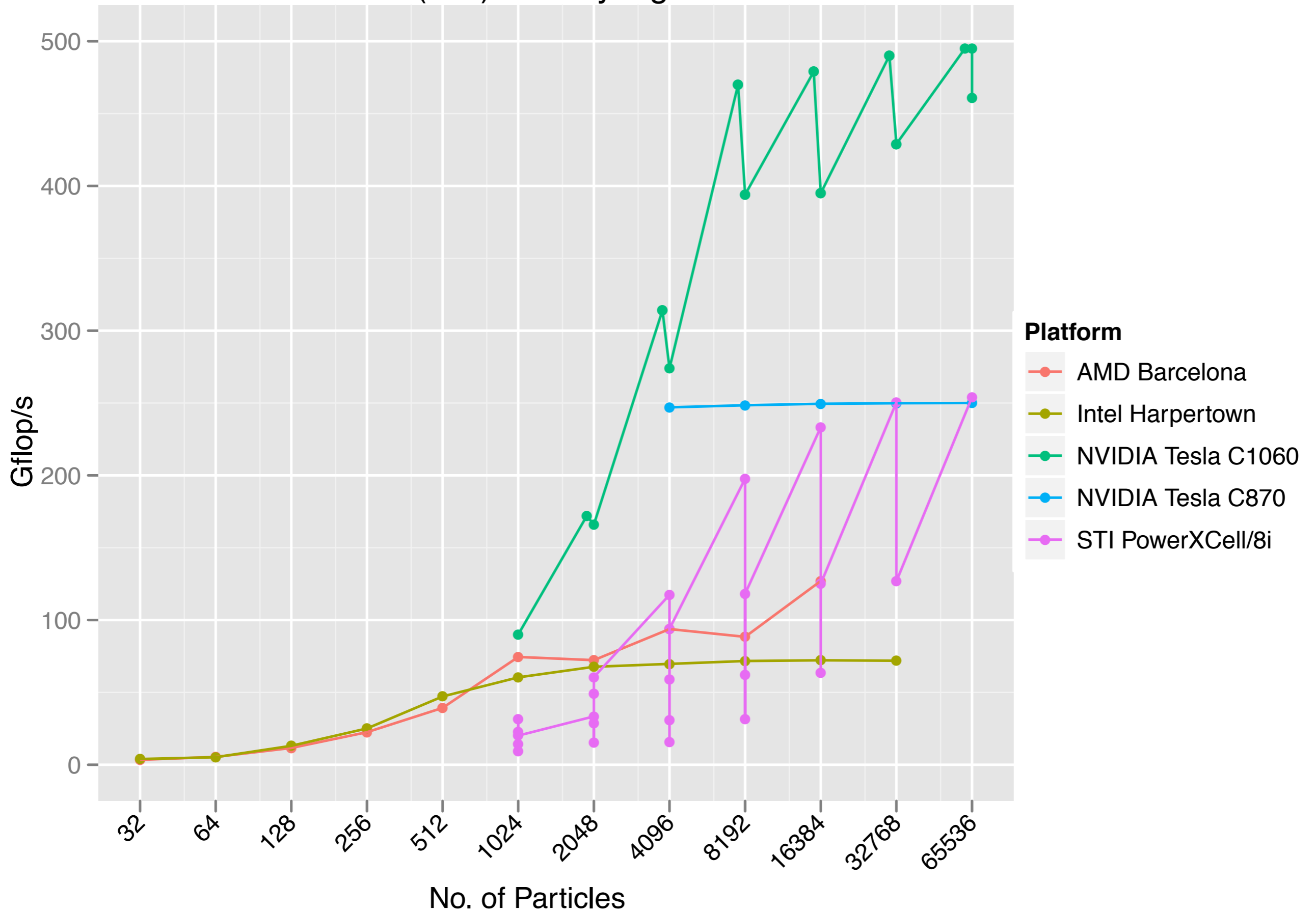- ▶ But, stubborn student (Aparna C.) is skeptical of speedup

Thursday, July 1, 2010

# Recall: Anatomy of an FMM



**Want**:
All pair interactions among green dots
$\rightarrow$ O($n^2$)

**FMM idea**:
Build tree, traverse & prune (approx.)
$\rightarrow$ O($n \log n$), with accuracy guarantee

► Bottleneck: **"Little" all-pair interactions ($b^2$) among nodes**

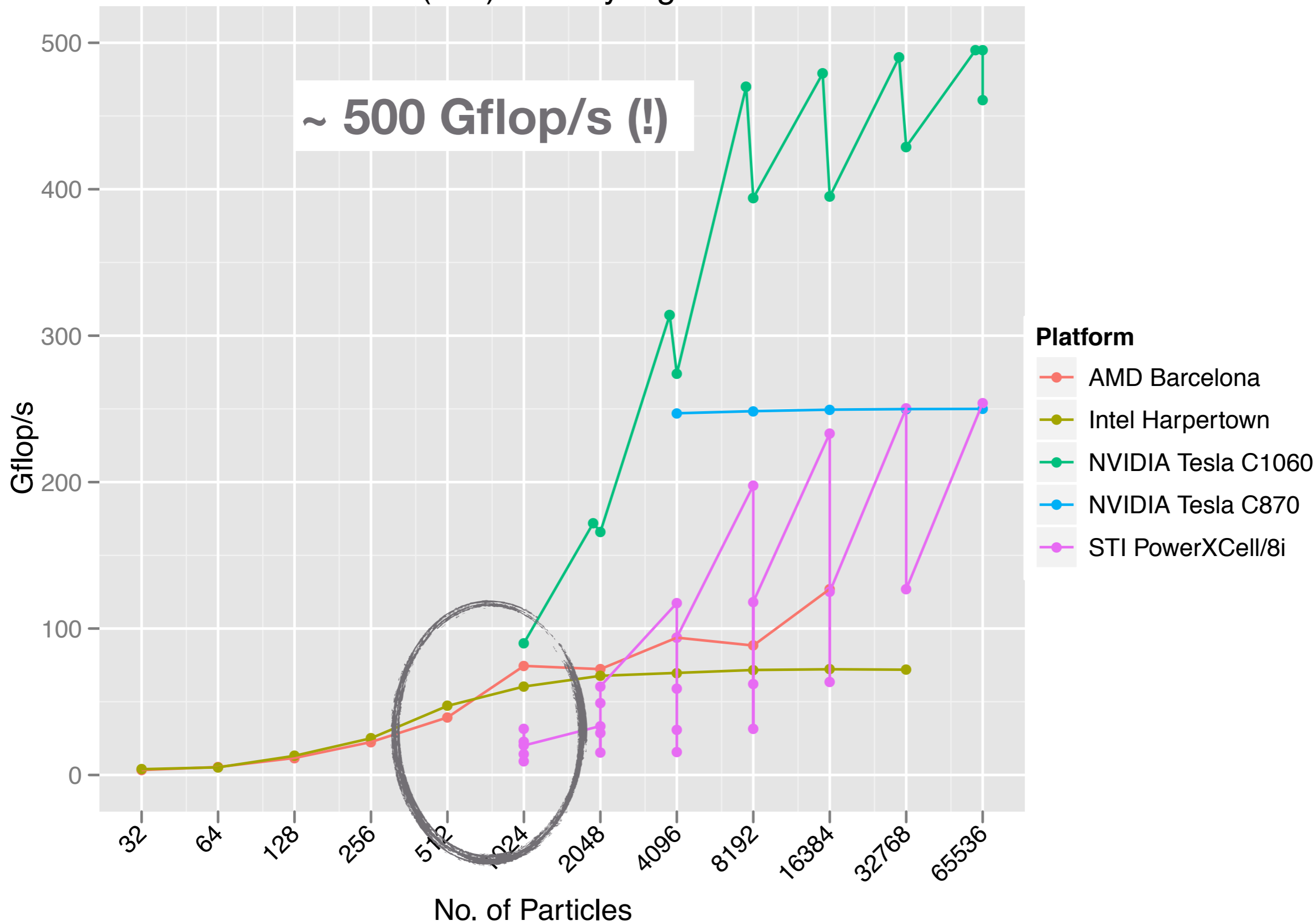► High compute intensity: 12x possible?

# Direct O(n^2) n-body algorithm



Thursday, July 1, 2010

# Direct O(n^2) n–body algorithm



~ 500 Gflop/s (!)

Gflop/s
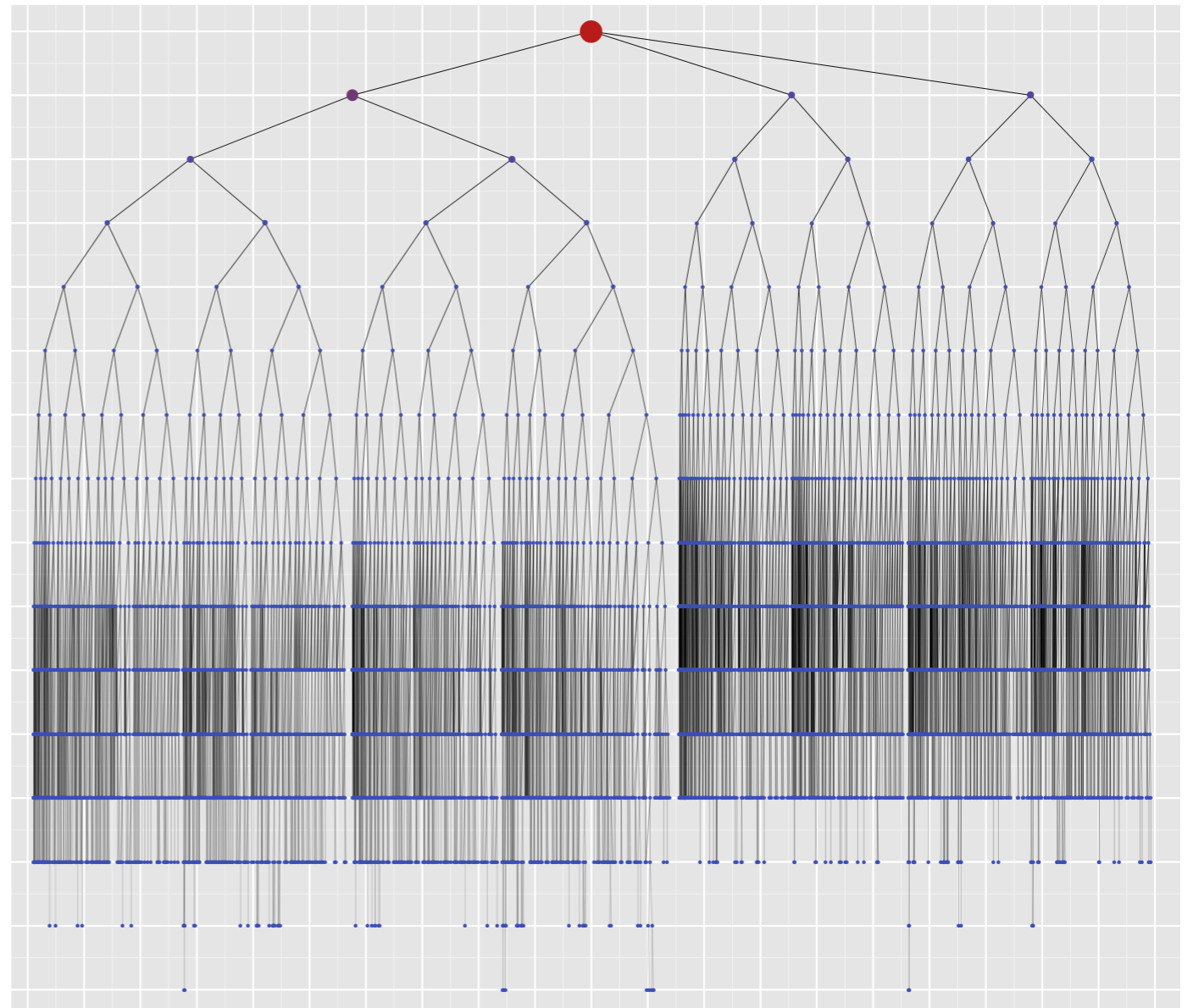
500
400
300
200
100
0

No. of Particles

32  64  128  256  512  1024  2048  4096  8192  16384  32768  65536

**Platform**

- AMD Barcelona
- Intel Harpertown
- NVIDIA Tesla C1060
- NVIDIA Tesla C870
- STI PowerXCell/8i

Thursday, July 1, 2010

# Direct O(n^2) n-body algorithm

~ 500 Gflop/s (!)

Gflop/s

No. of Particles

**Platform**

- AMD Barcelona
- Intel Harpertown
- NVIDIA Tesla C1060
- NVIDIA Tesla C870
- STI PowerXCell/8i

Thursday, July 1, 2010

▶ Case study 3: Sparse direct solvers

(1) M. Efe Guney, Ph.D. Thesis, Georgia Tech, May '10.  (2) Guney, Czechowski, et al. (*in progress*)
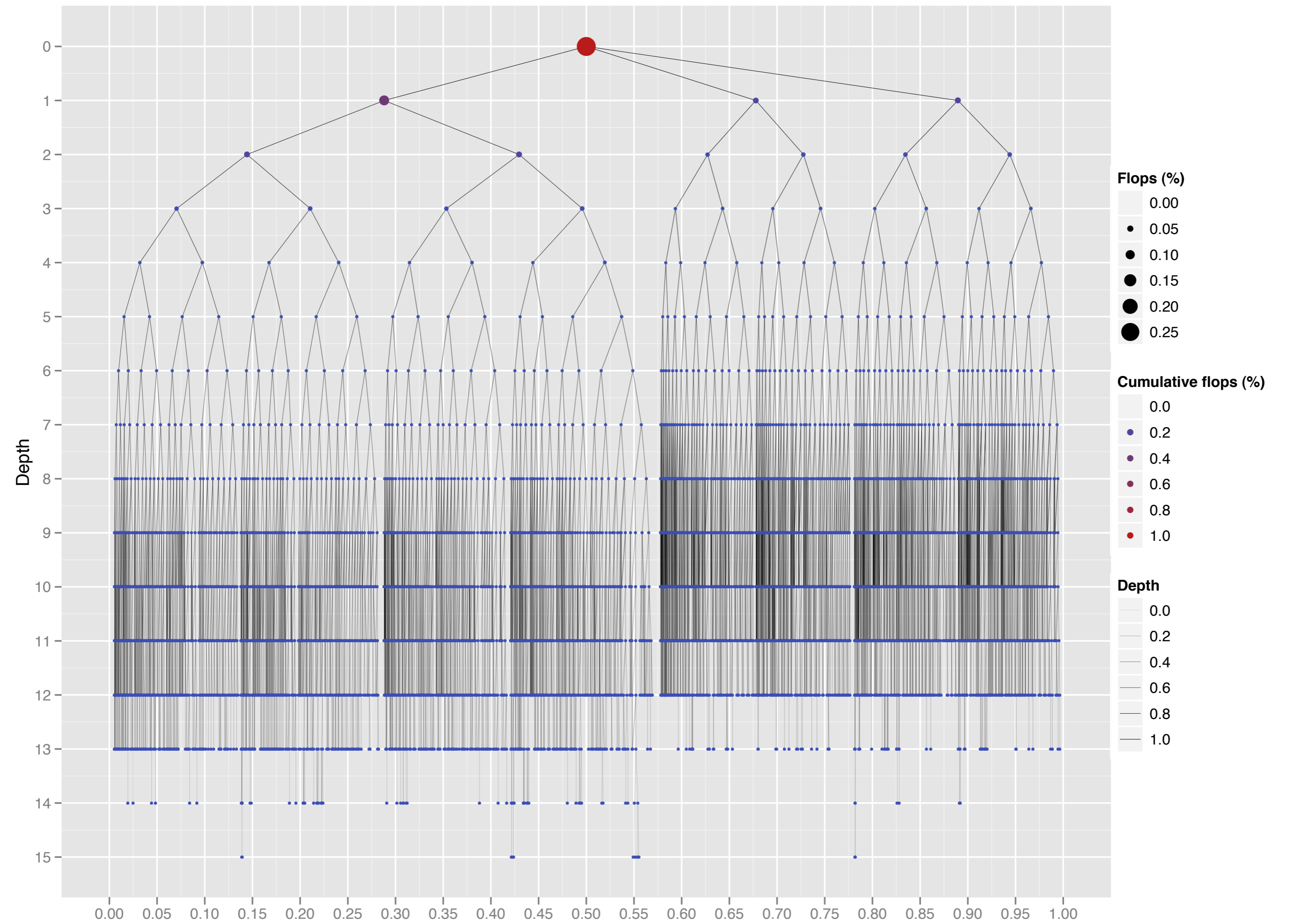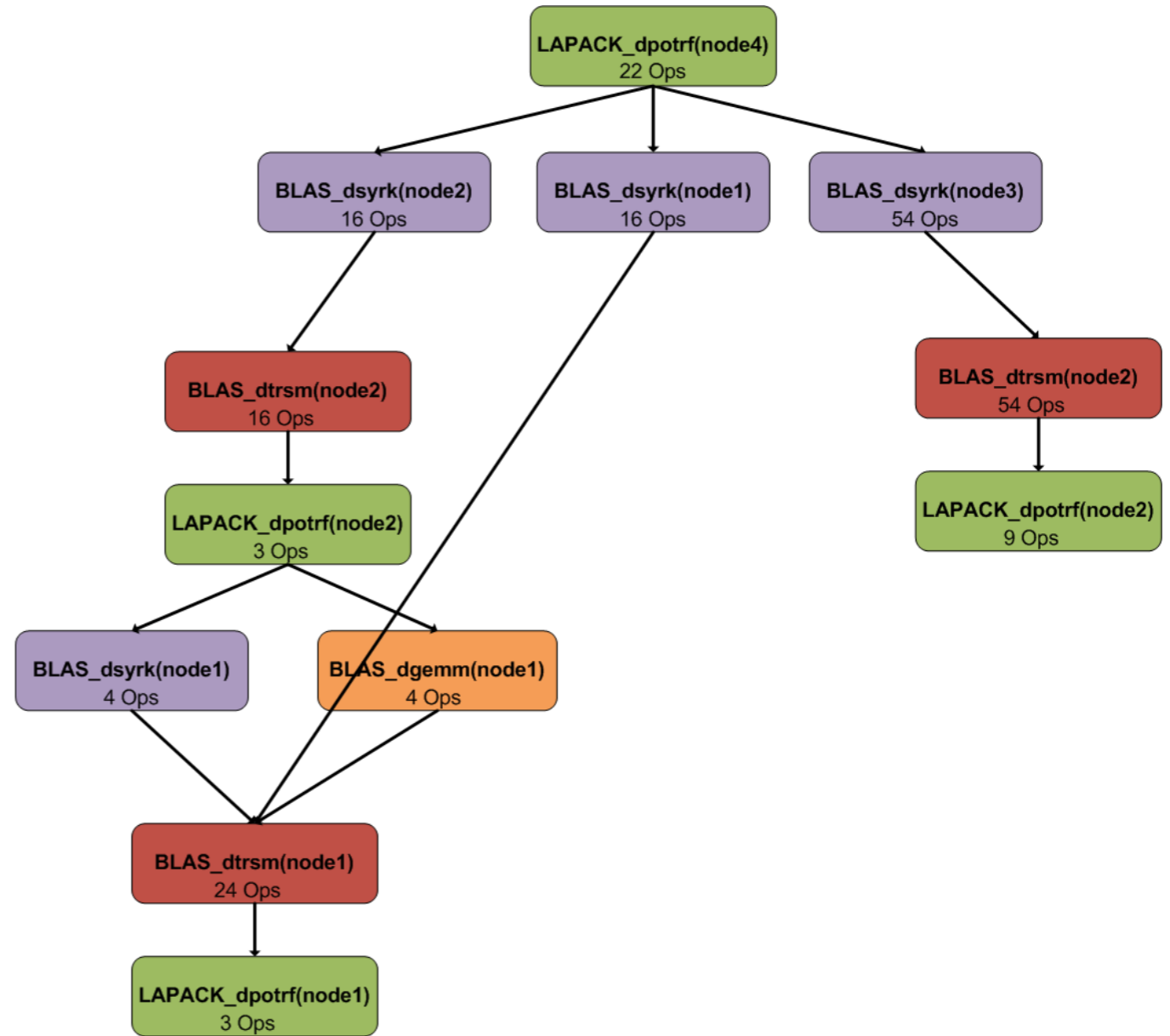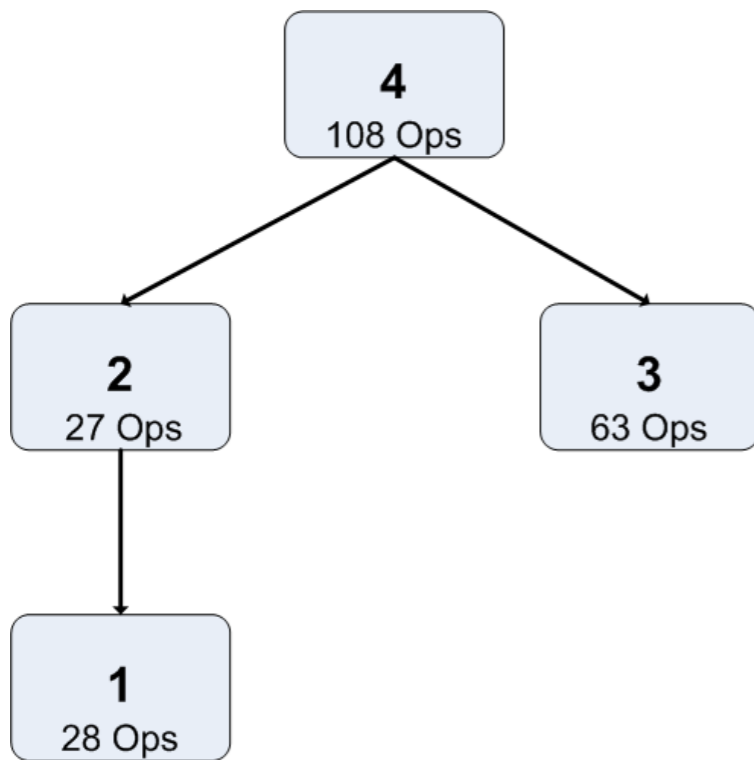
# Anatomy of a sparse direct solver



▶ Sparse Cholesky factorization, $A = L \cdot L^T$, where $A$ & $L$ are sparse

▶ Mixed compute intensity, average of ~ 4 flops : byte for sample problem

# Real elimination tree example

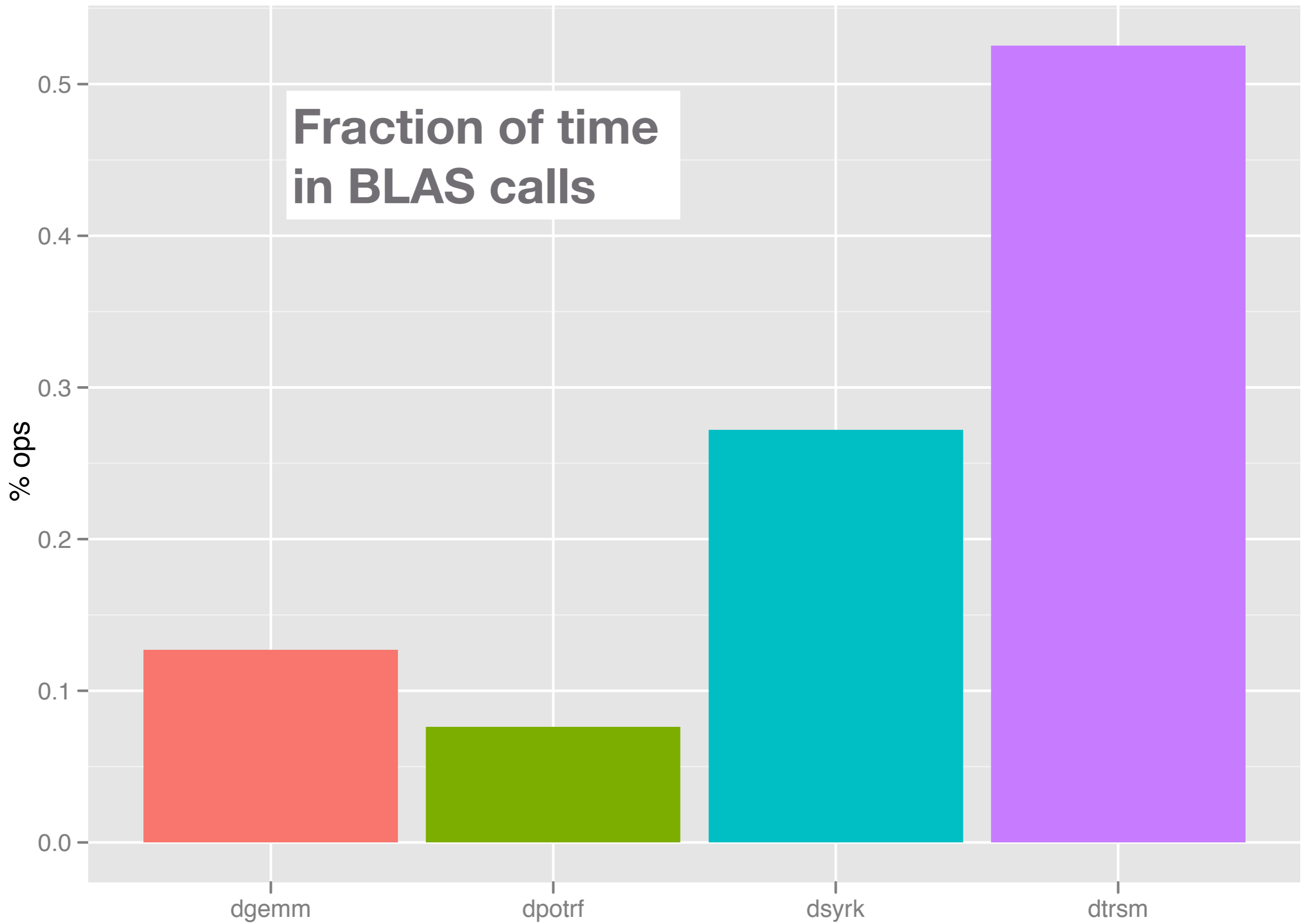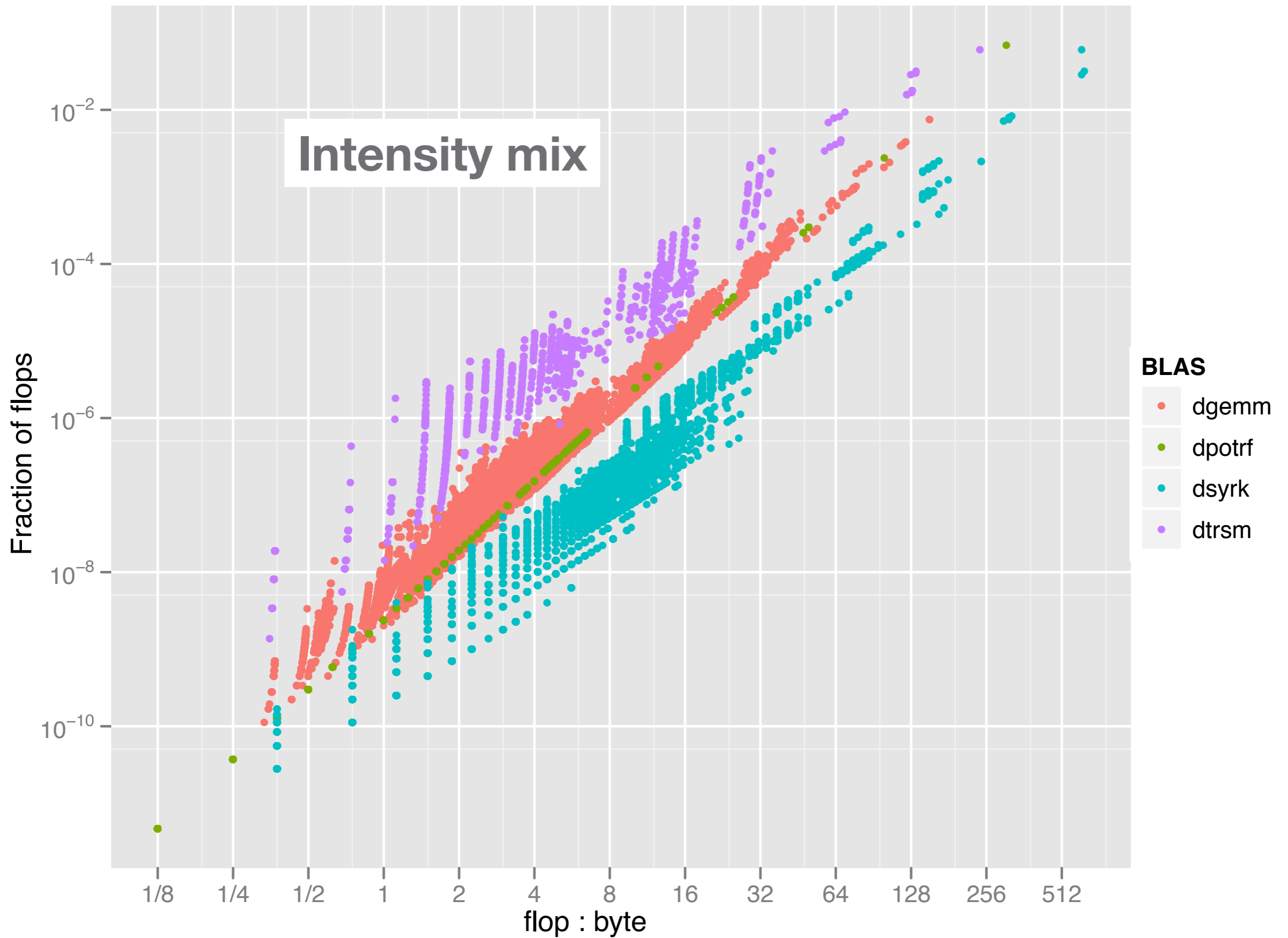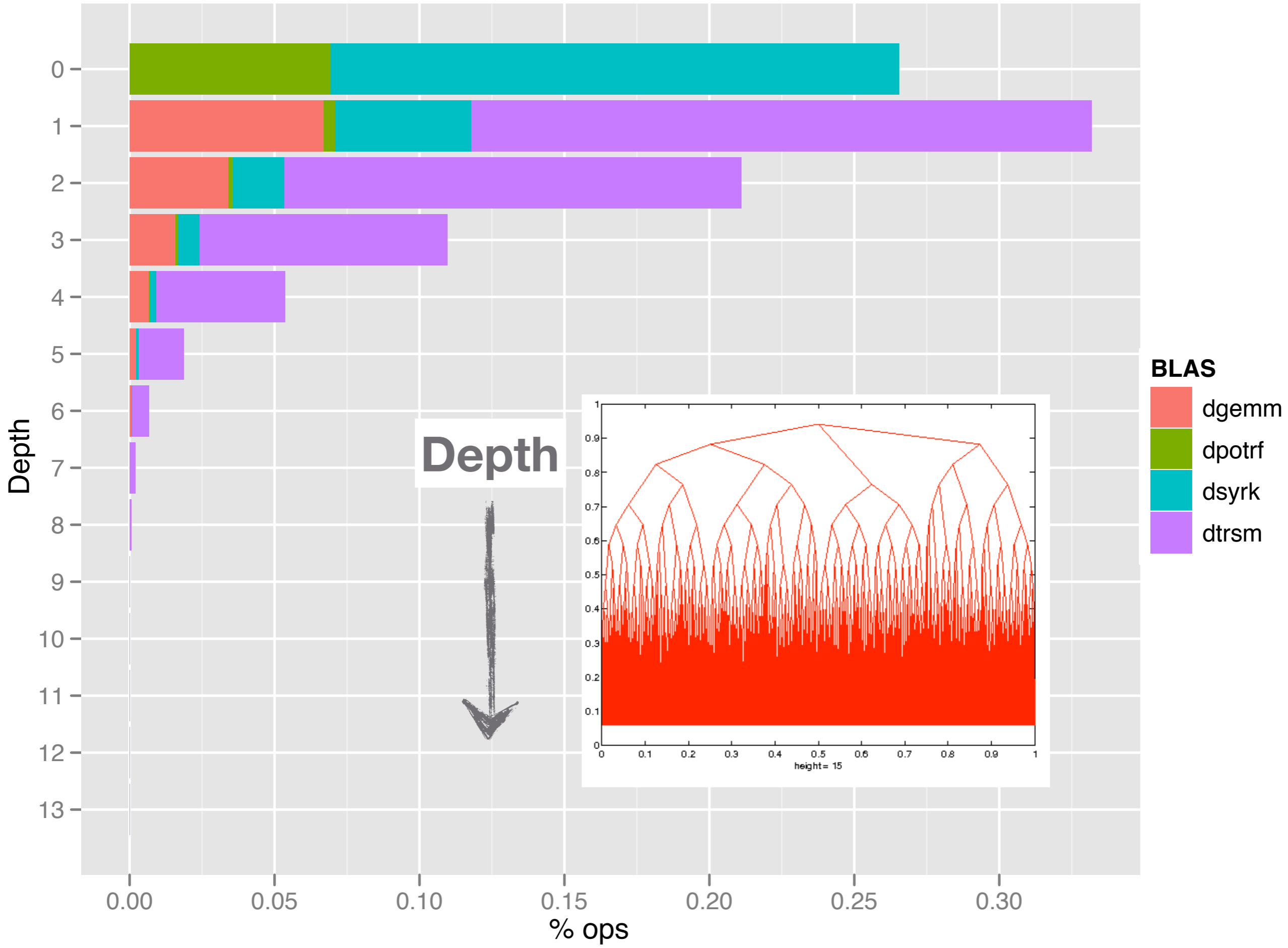Independent subtrees may be processed in parallel.

Thursday, July 1, 2010

# Finer-grained dependencies

Colored circles on the right are BLAS calls on operands of varying size.

Fraction of time
in BLAS calls

% ops

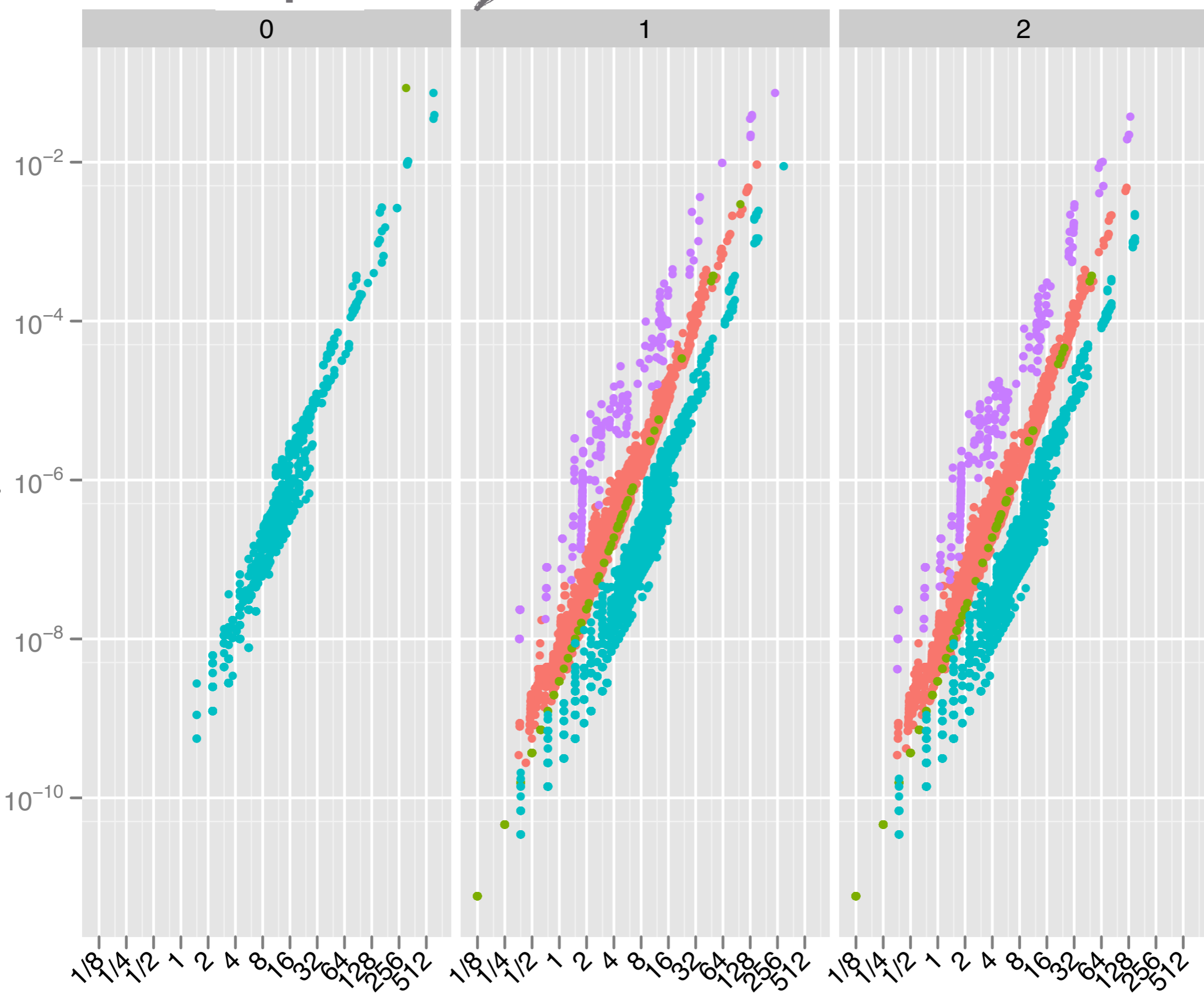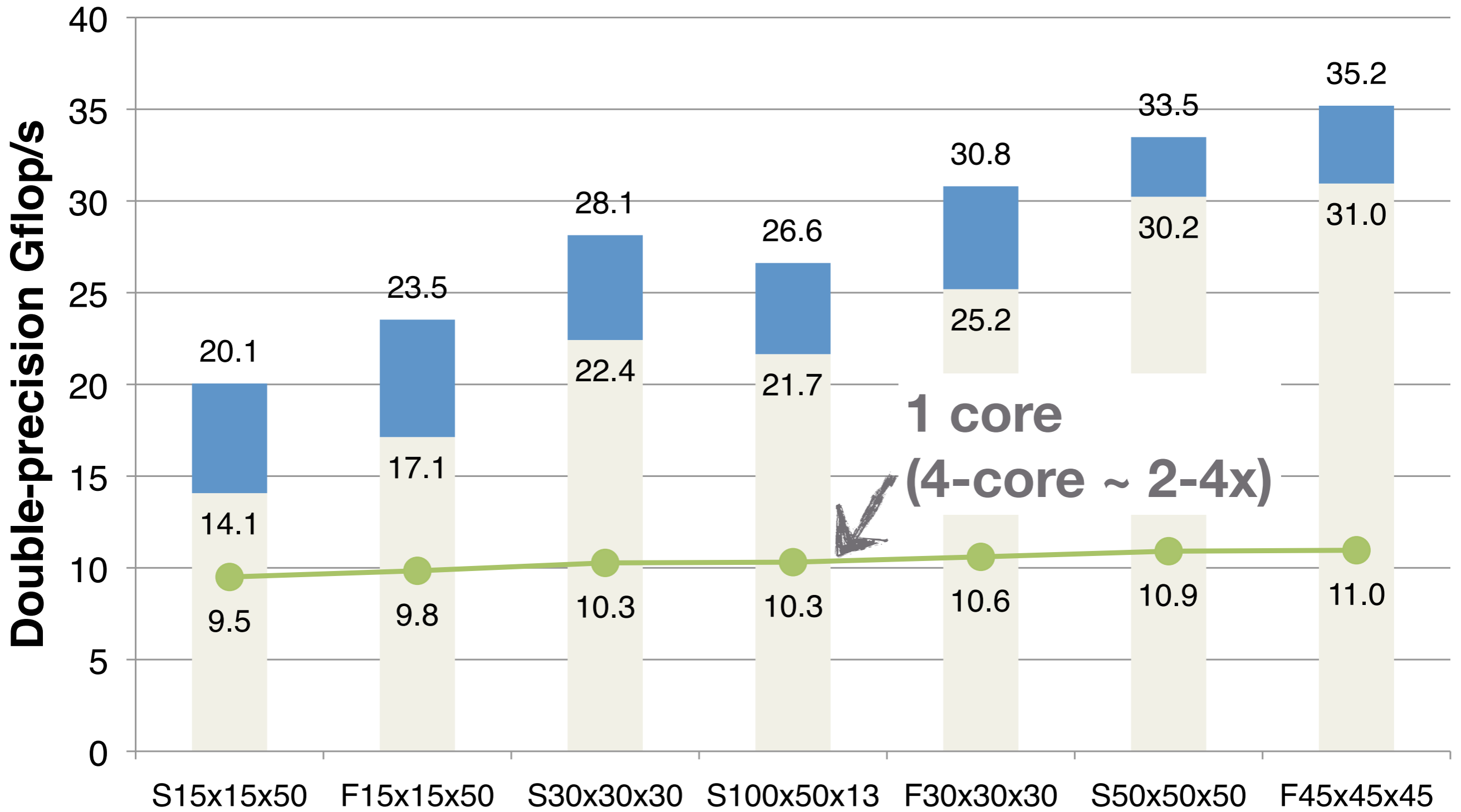dgemm    dpotrf    dsyrk    dtrsm

Thursday, July 1, 2010

Intensity mix

Fraction of flops vs flop : byte

BLAS
- dgemm
- dpotrf
- dsyrk
- dtrsm

Thursday, July 1, 2010

# Summary:
# Definite potential, but no "magic"

▶ For a set of semi-irregular computations, **1 GPU** ≅ **1 to 2 CPUs**.
  ⇒ Study heterogeneous computing, but scale-back expectations.

▶ Barriers?

  ▶ Bandwidth-bound: Aggregate GPU : CPU bandwidth < 3x

  ▶ PCIe: Will we really replicate GPU memory system in on-die CPU/GPU?

  ▶ Compute-bound: 5 to 10x potential, but how? E.g., FMM granularity mismatch

  ▶ Productivity: To 0th order, tuning required on all platforms (Bailey's Rule #6)