

Should I port my code to a GPU?

Casey Battaglino - Aparna Chandramowliswaran - Jee Choi -
Kent Czechowski - M. Efe Guney (*UC Davis*) - Chris McClanahan -
Logan Moon - Dave Noble - Aashay Shringarpure (*Google*) -
Richard (Rich) Vuduc

Clusters, Clouds, and Grids for Scientific Computing

Flat Rock, North Carolina - September 9, 2010

**Georgia
Tech**



College of
Computing

Computational Science and Engineering



A: It depends. Opportunity cost?

- (1) Who are you?
- (2) What is your app?
- (3) What are your performance, productivity, and portability goals?

For most of us in this room, I'd say, "yes."

For the "average" apps developer, I'd say "not yet."

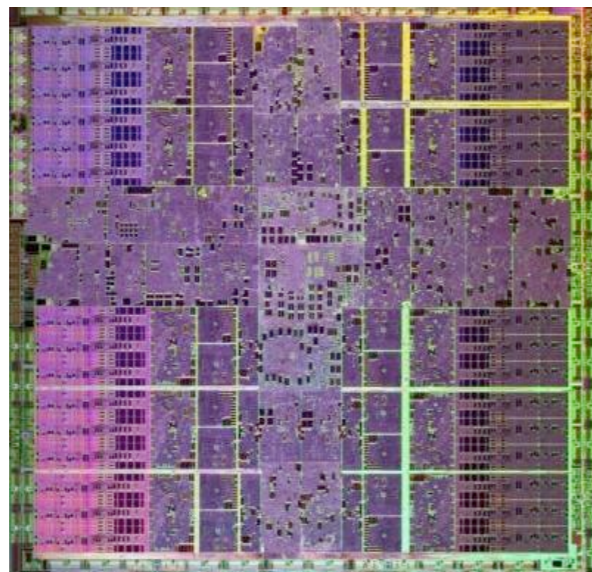


Q: Pay-off of a GPU port?

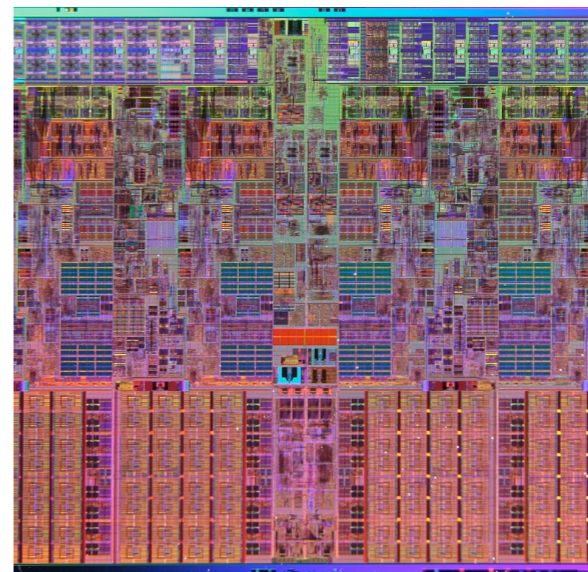
(Posed to me by Scott Klasky at ORNL)

- ▶ Meta-analysis, for **semi-irregular** sci. comp. + data analytics apps (sparse iterative + direct solvers; tree-based particle methods)
- ▶ **A: Given roughly same level of tuning & power*, ...**

GPU



IR



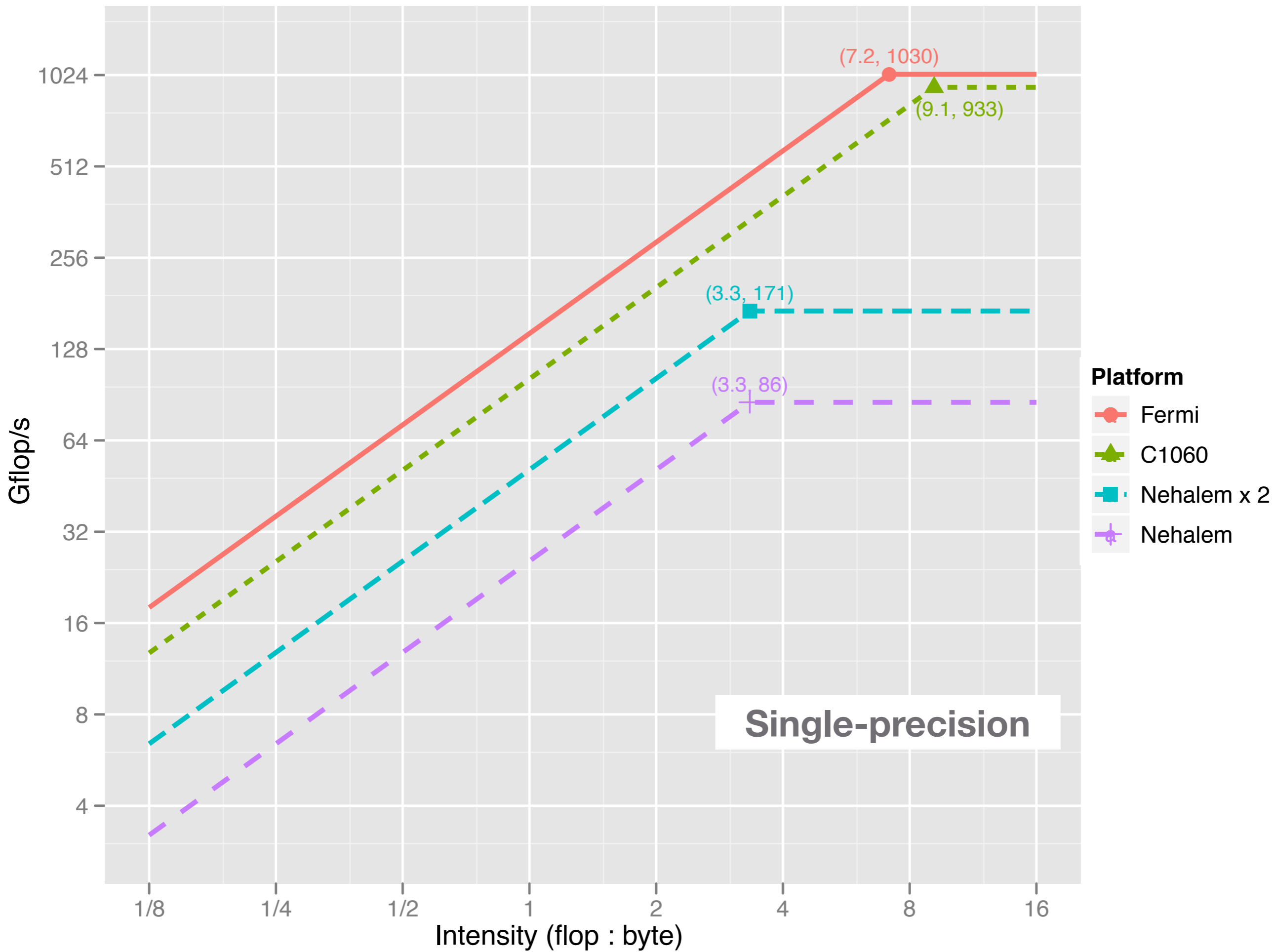
x 2 CPUs

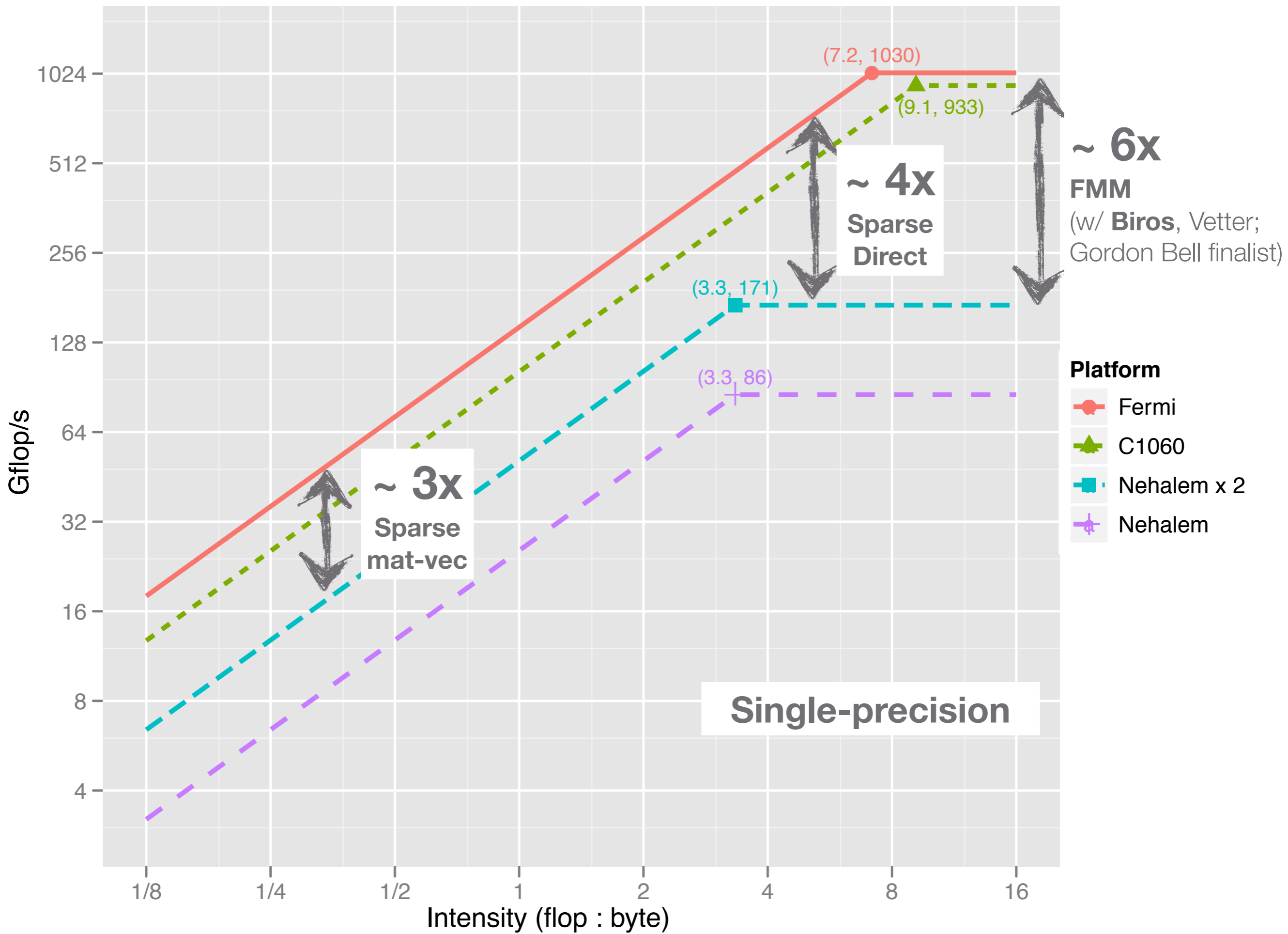
Reason 1: The potential is real, but might be less than you expect.

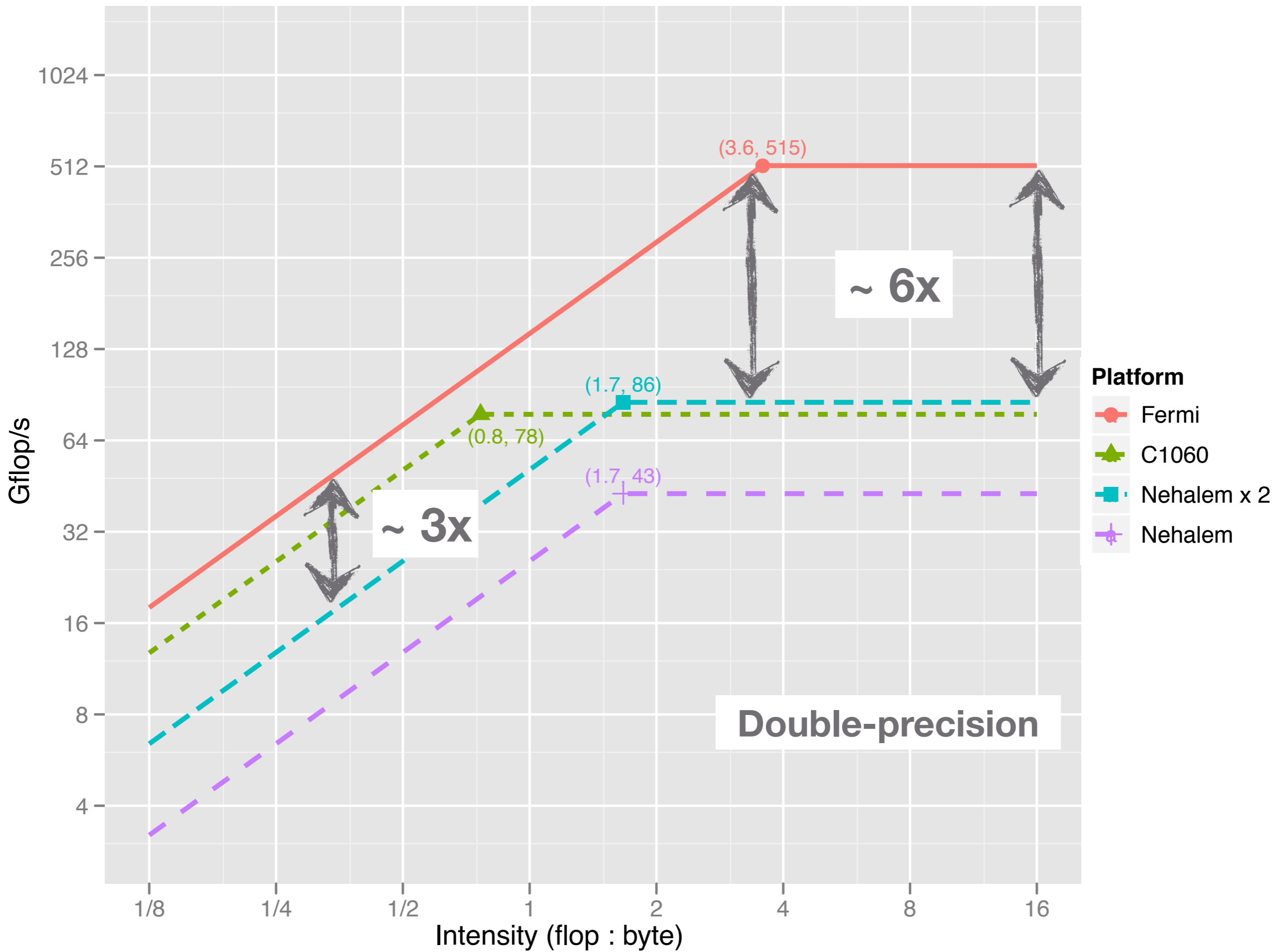


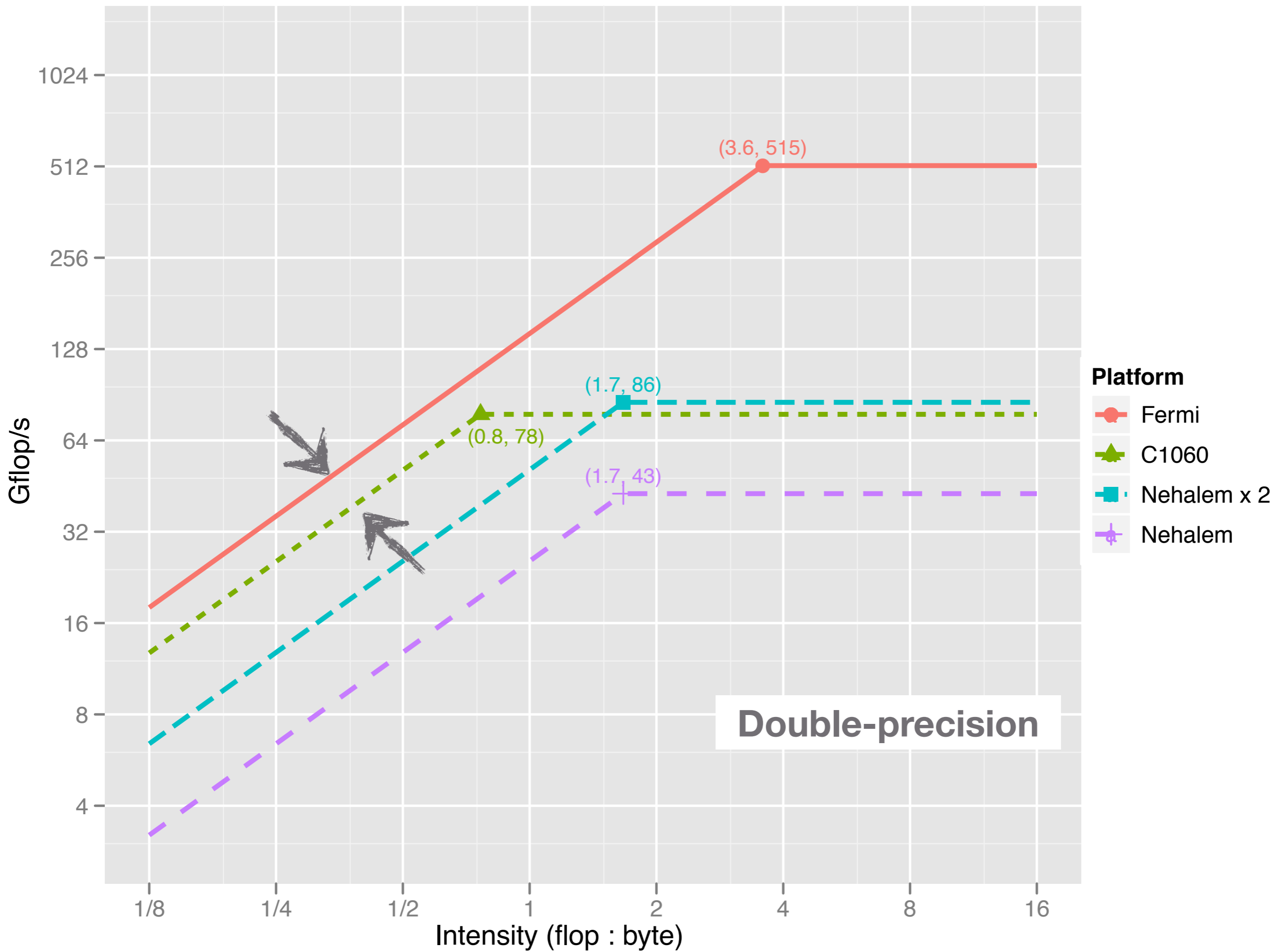
Architecture	Intel Nehalem X5550	NVIDIA T10P C1060	NVIDIA GT200 GTX 285	NVIDIA Fermi C2050
GHz	2.66	1.44	1.47	1.15
Sockets	2	1	1	1
Cores per socket	4	30	30	15*
Peak Gflop/s single (<i>double</i>)	170.6 (85.3)	933 (78)	1060 (88.4)	1030 (515)
Peak GB/s	51.2	102	159	144
Sys. Watts (sockets only)	375 (200)	200	204	247











Reason 2: Productivity: Though there is potential, there is also no free lunch.



Parallel Sorting (survey)

(Does **not** include Merrill & Grimshaw '10)



Parallel Sorting (survey)

(Does **not** include Merrill & Grimshaw '10)



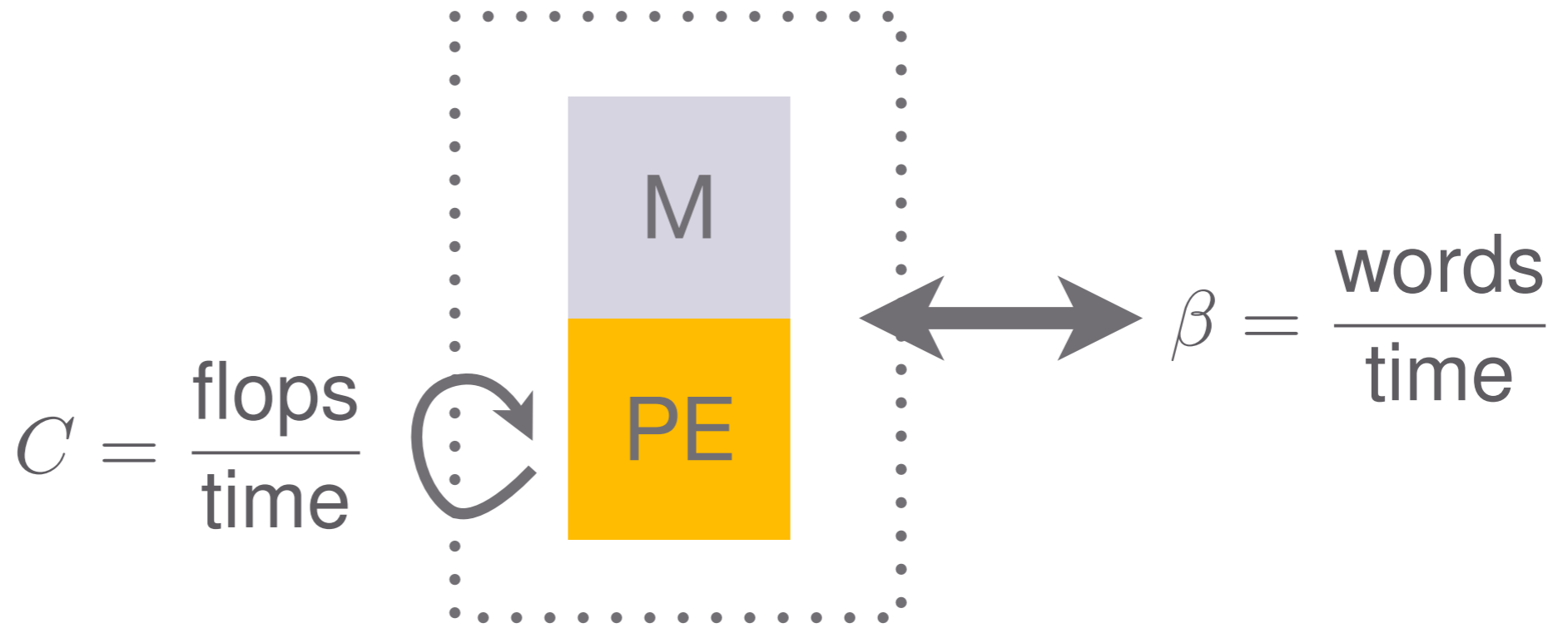
Parallel Sorting (survey)

(Does **not** include Merrill & Grimshaw '10)

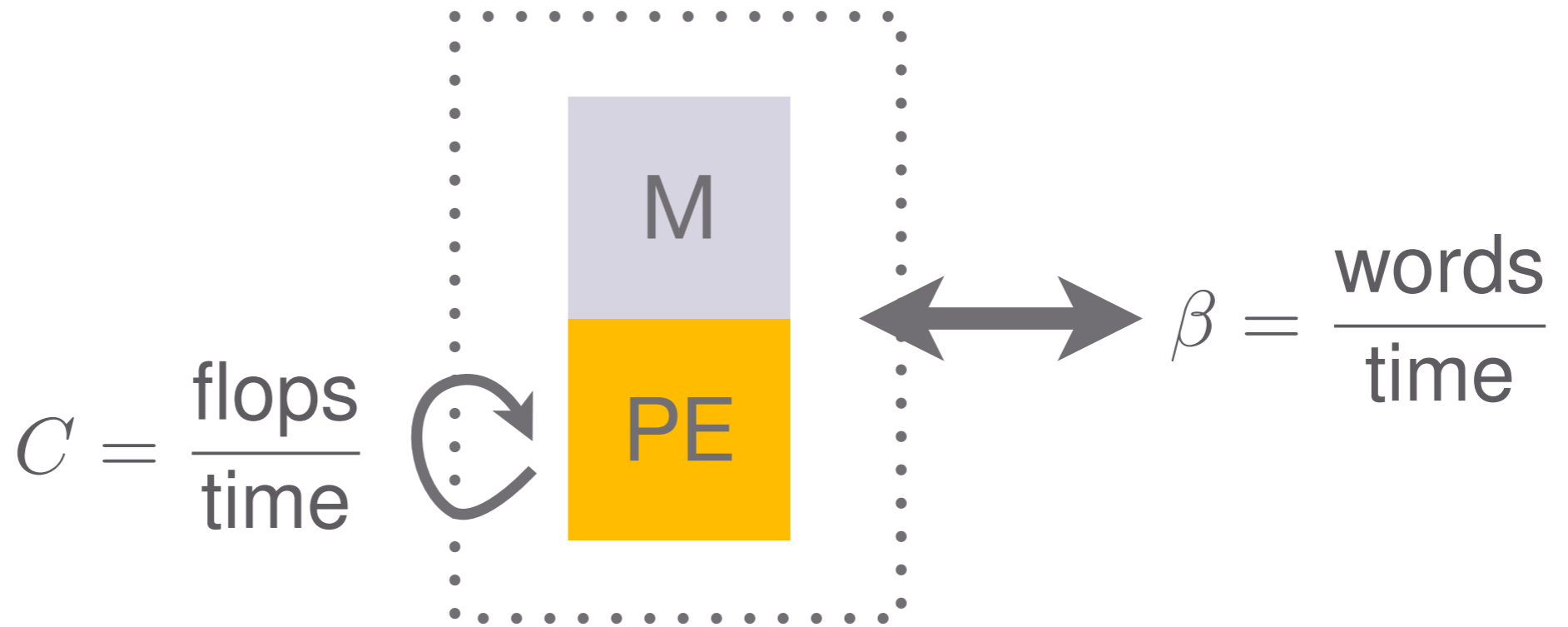


Reason 3: It's a moving target that might end up converging to what we had before.





HPC 101

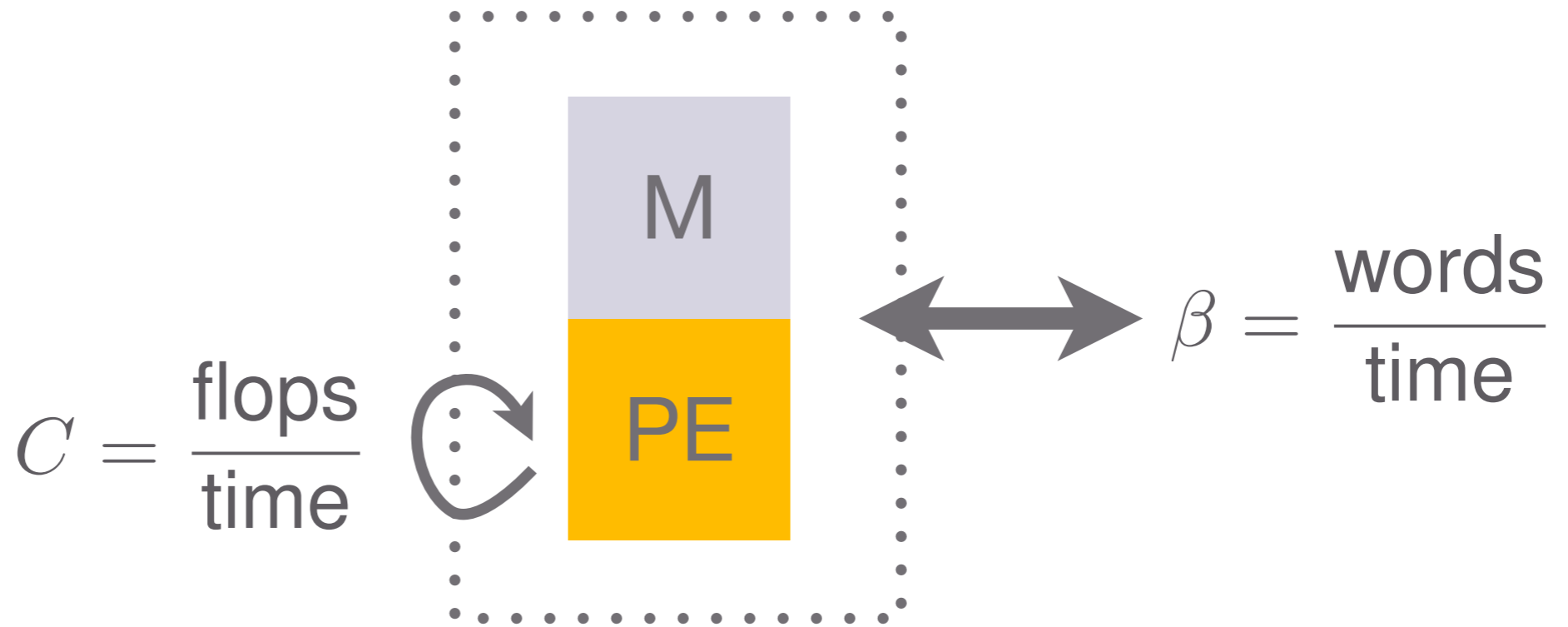


Balance equation for an I/O-optimal matrix multiply:

$$\frac{C}{\beta} = \Theta \left(\sqrt{M} \right)$$

HPC 101

See, for example, Kung (ISCA 1986).



Balance equation for an I/O-optimal matrix multiply:

$$\frac{C}{\beta} = \Theta(\sqrt{M})$$

For comparison-based sort:

$$\frac{C}{\beta} = \Theta(\log_2 M)$$

HPC 101

See, for example, Kung (ISCA 1986).

Having said all that, I am still **optimistic** about the role GPUs will play in current and future systems!

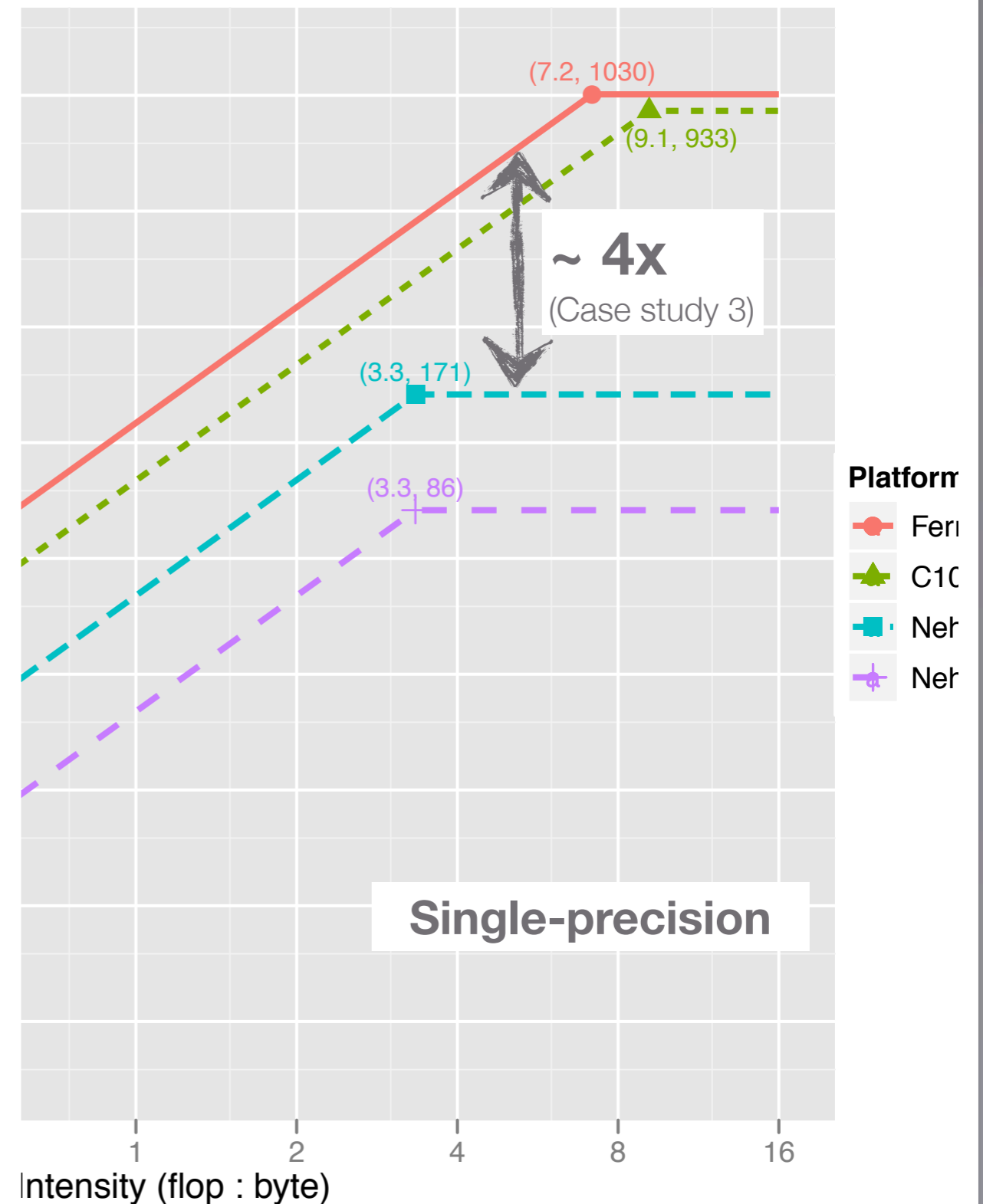


Sparse direct solvers

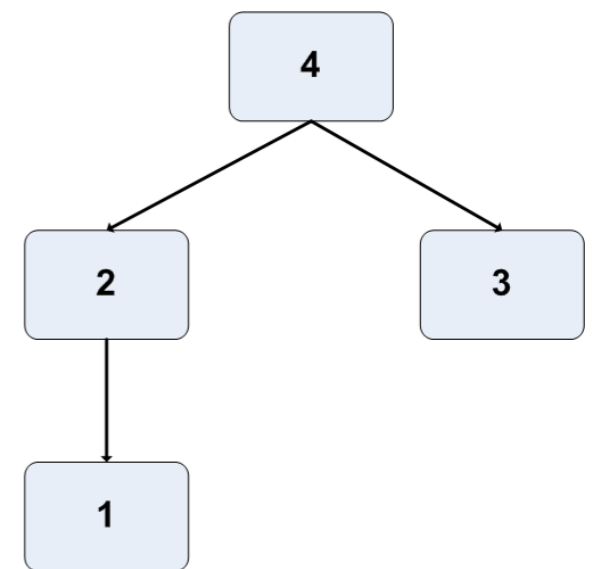
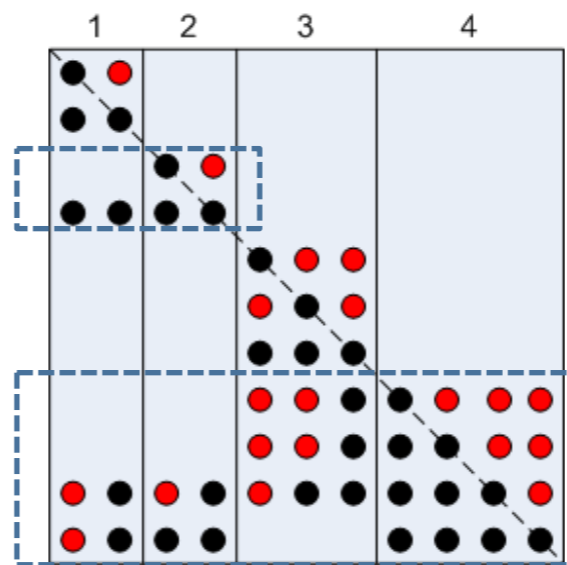
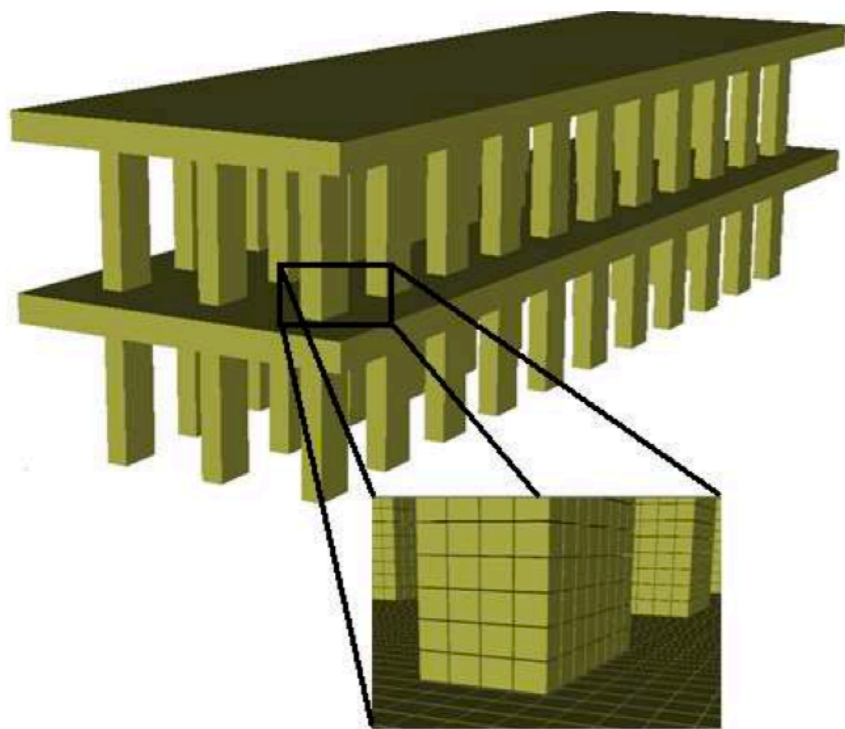
Kent Czechowski,
M. Efe Guney, R. Vuduc

(Work in progress)

M. Efe Guney. *High-performance direct solution of finite-element problems on multi-core processors*. Ph.D. Thesis, School of Civil Engineering, Georgia Tech, May 2010.

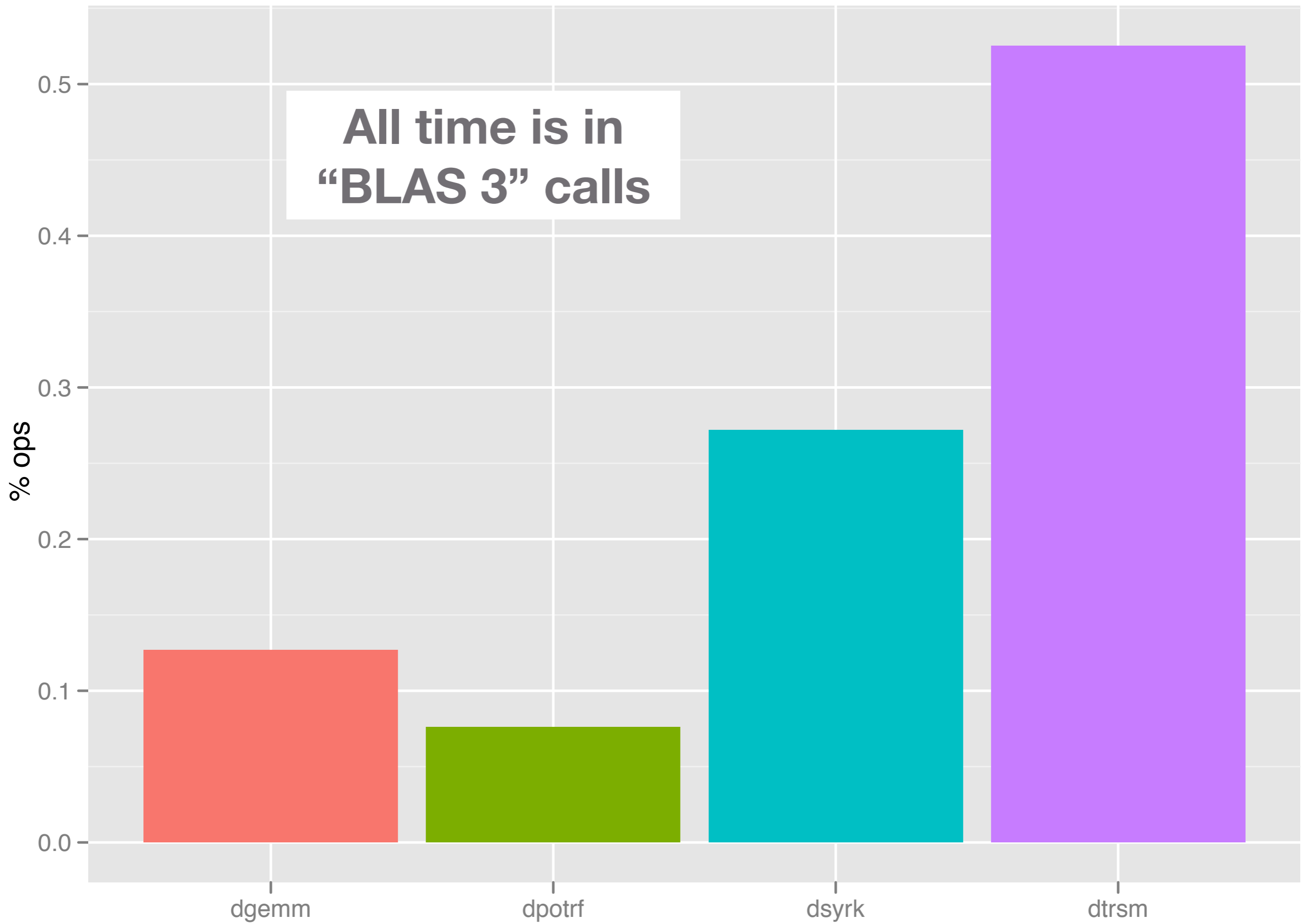


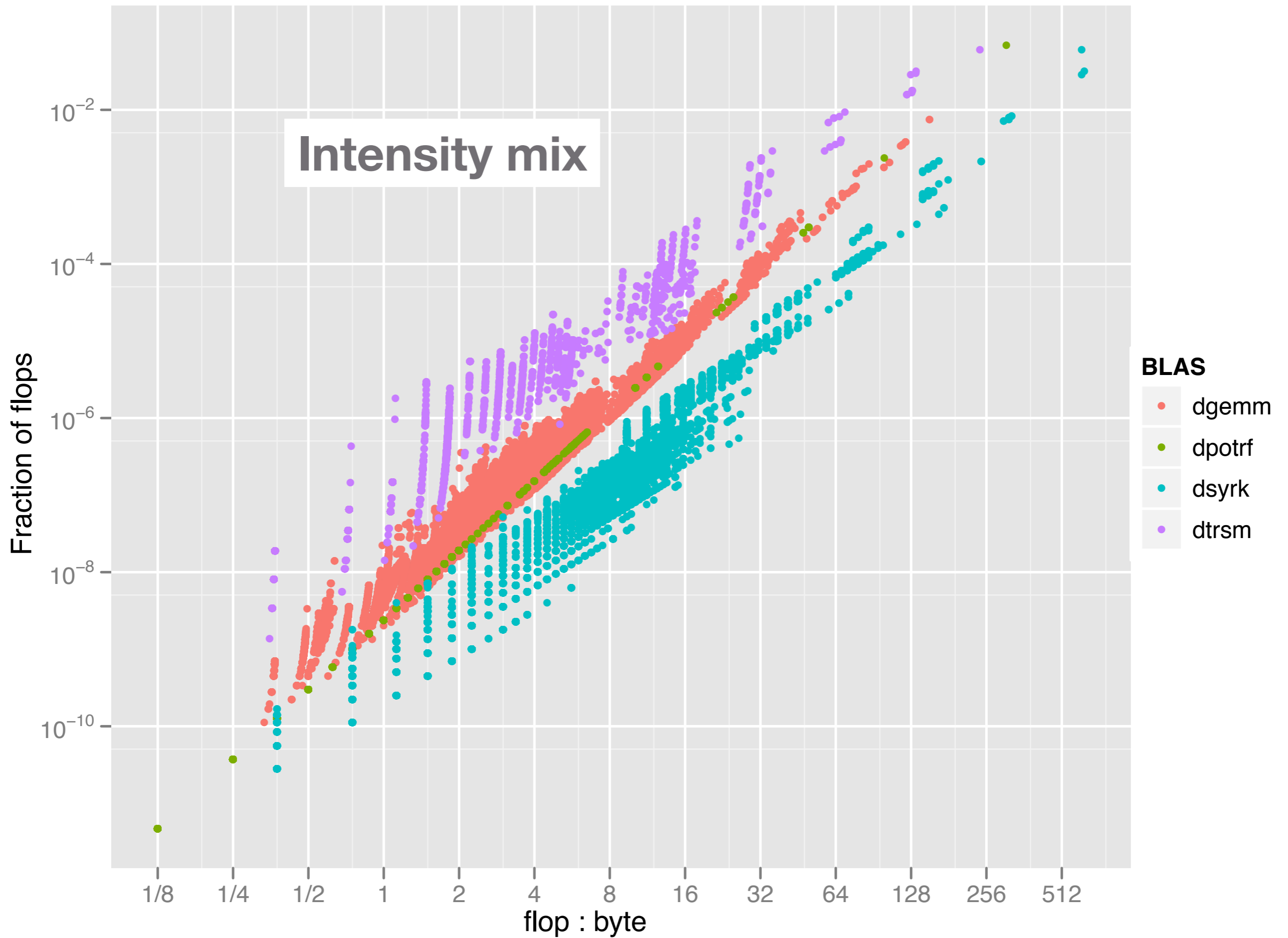
Anatomy of a sparse direct solver

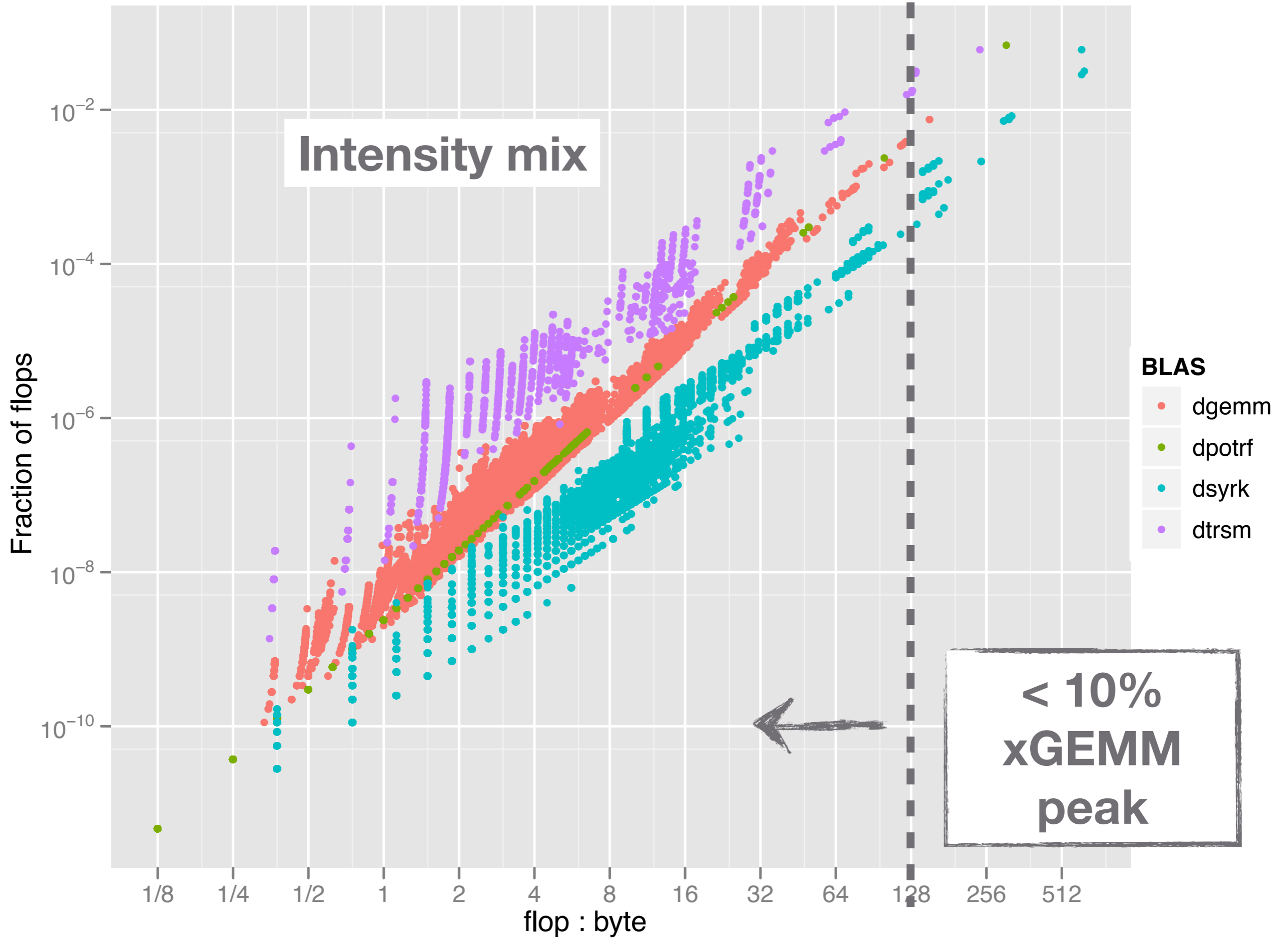


- ▶ Sparse Cholesky factorization, $A = L \cdot L^T$, where A & L are sparse
- ▶ Mixed compute intensity, average of ~ 4 flops : byte for sample problem

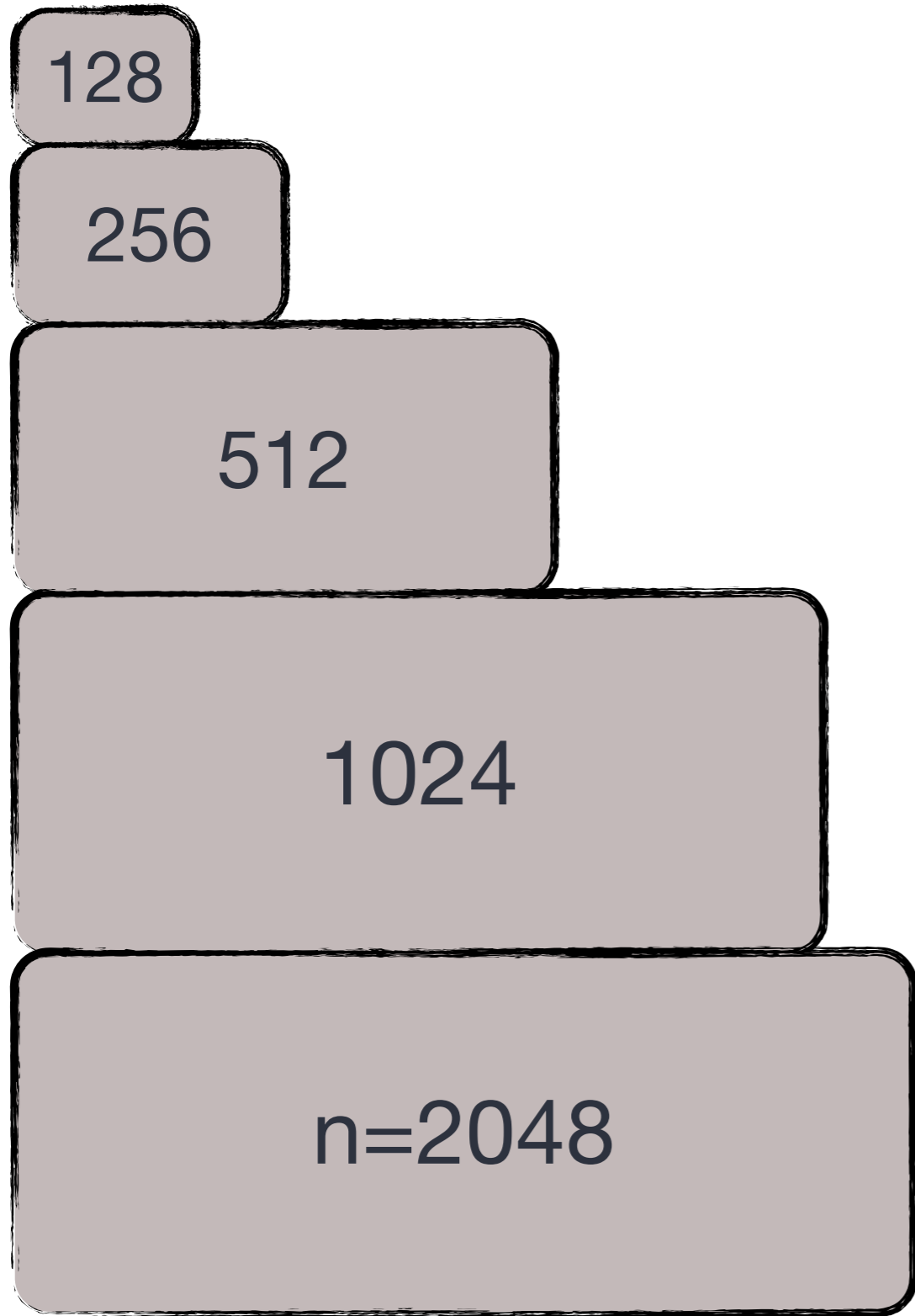
**All time is in
“BLAS 3” calls**





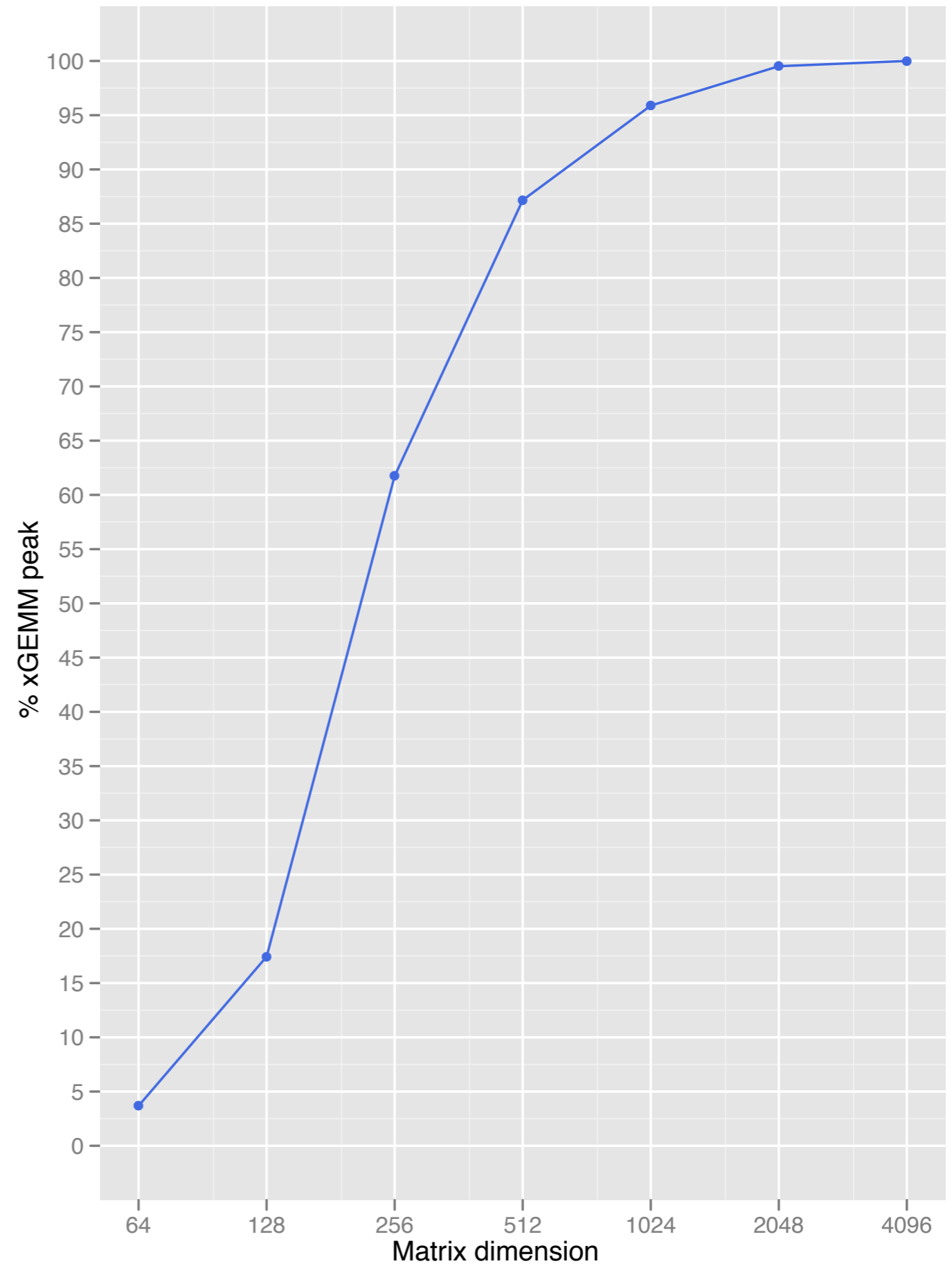
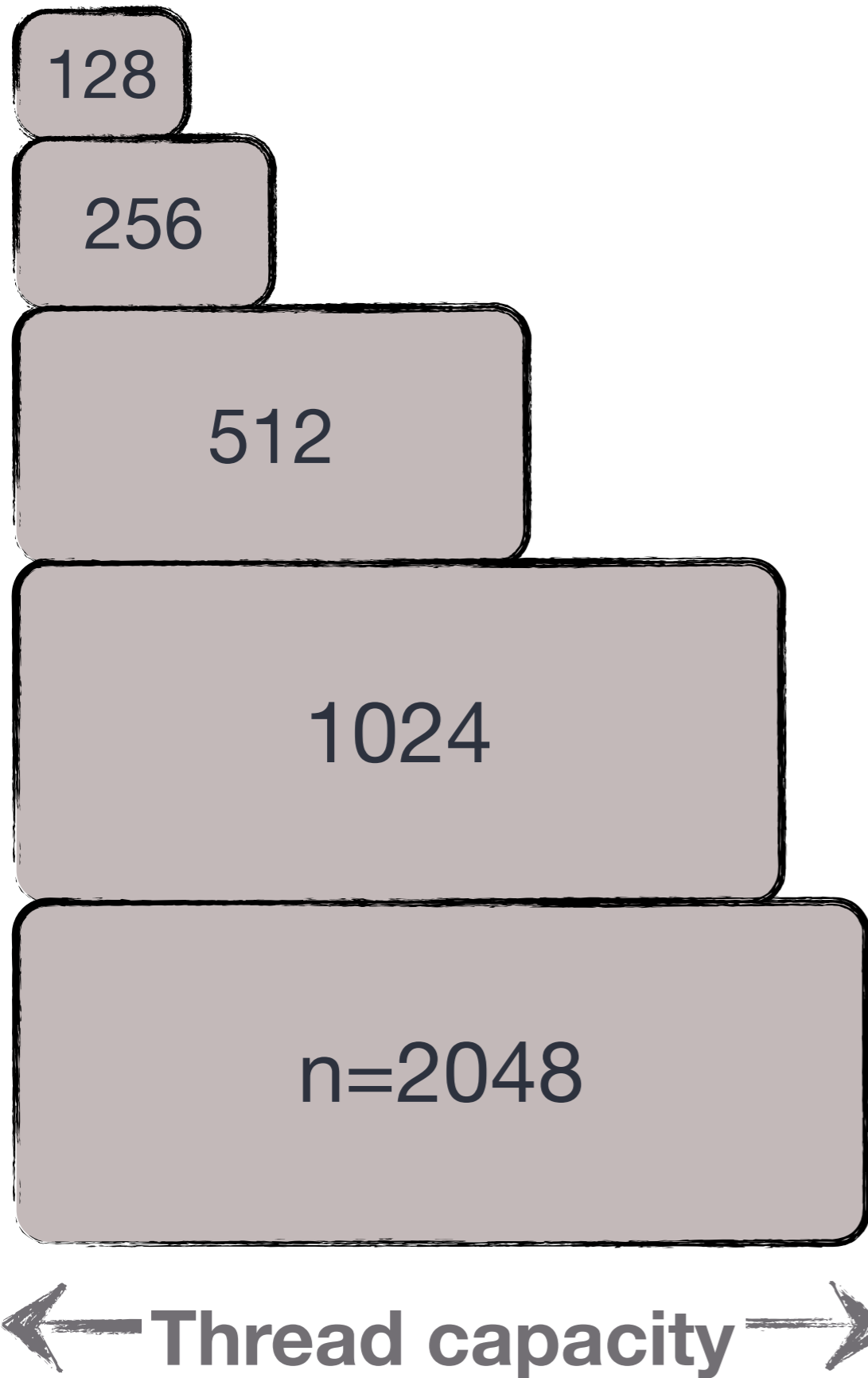


Time

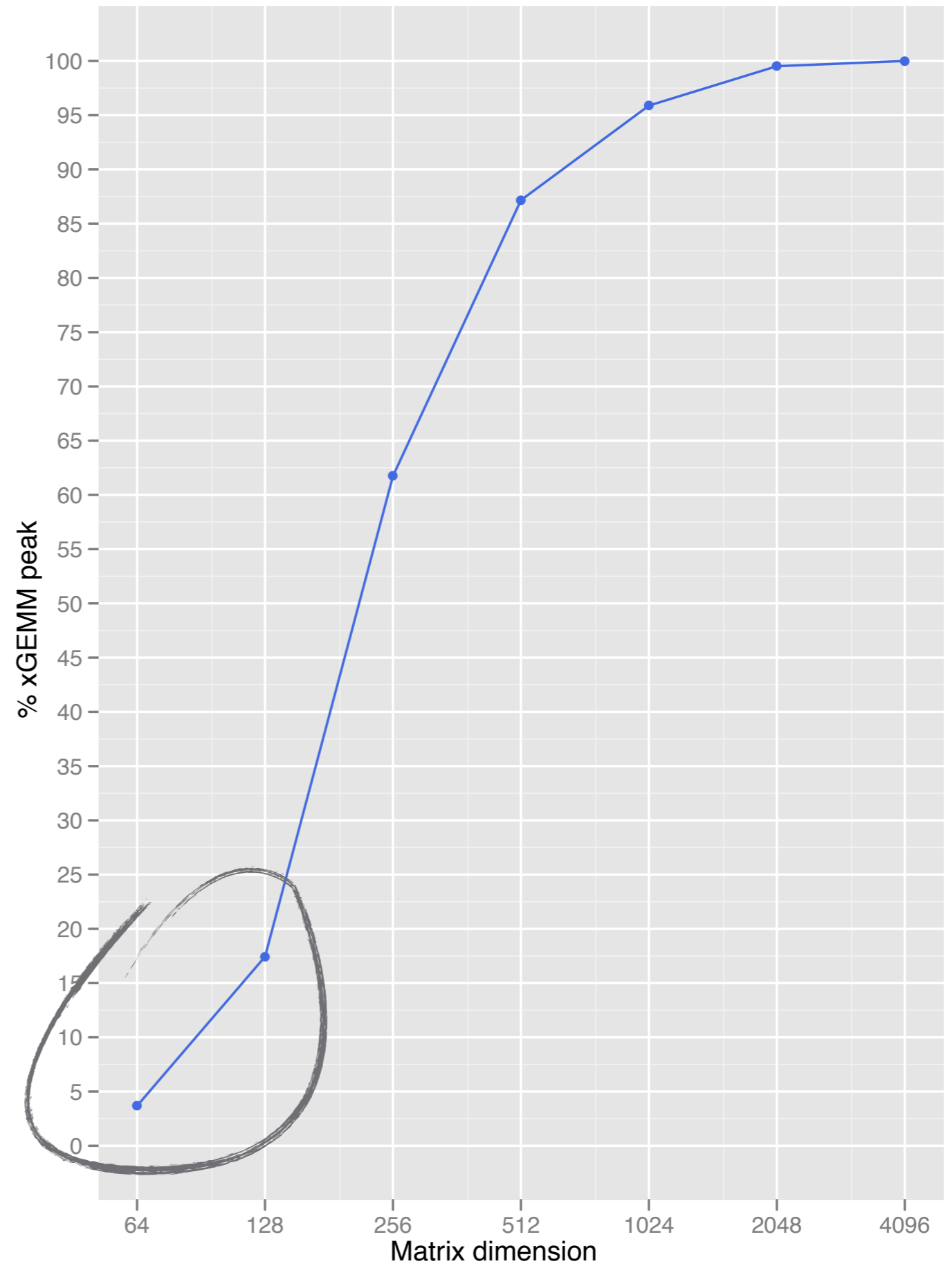
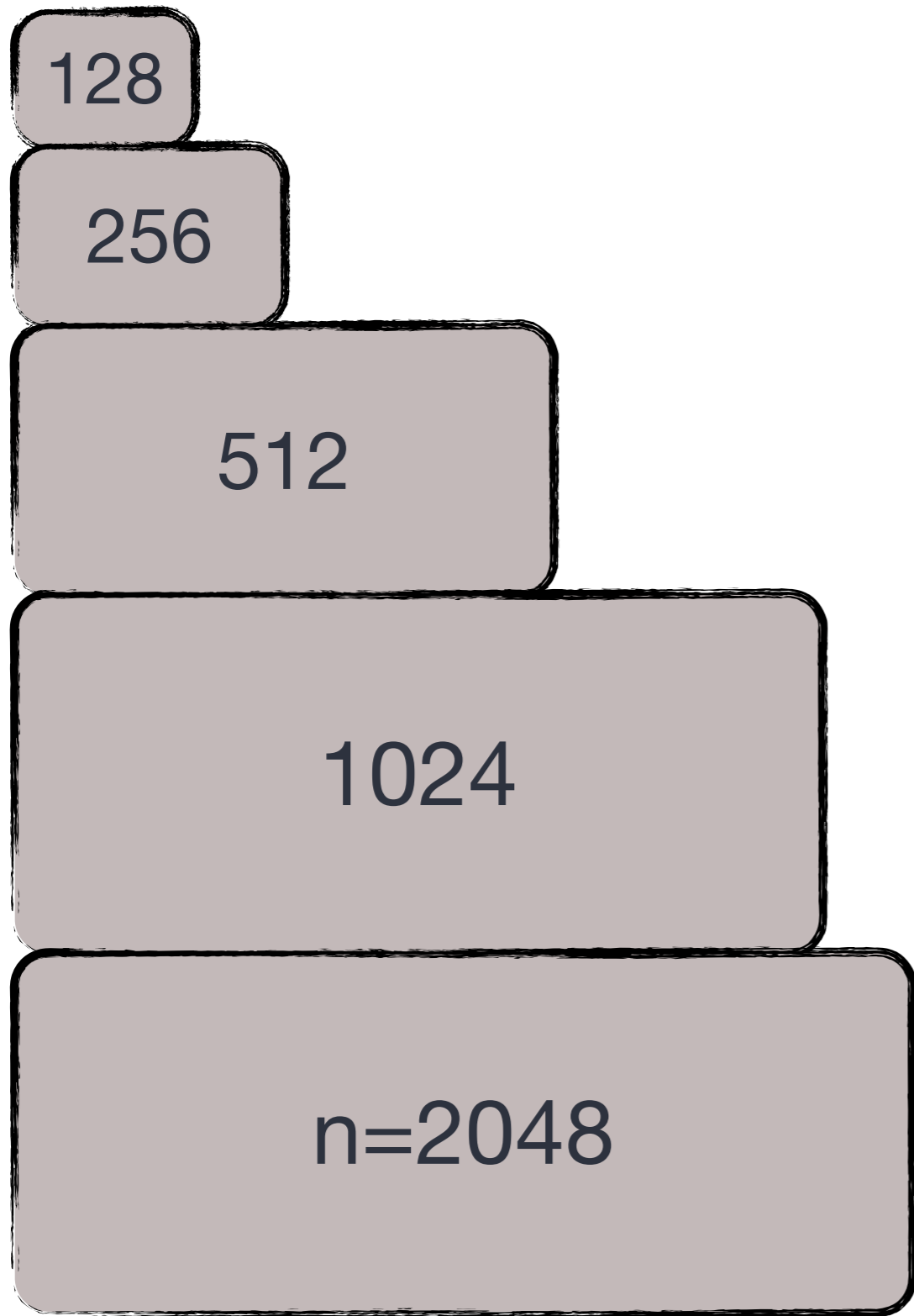


← Thread capacity →

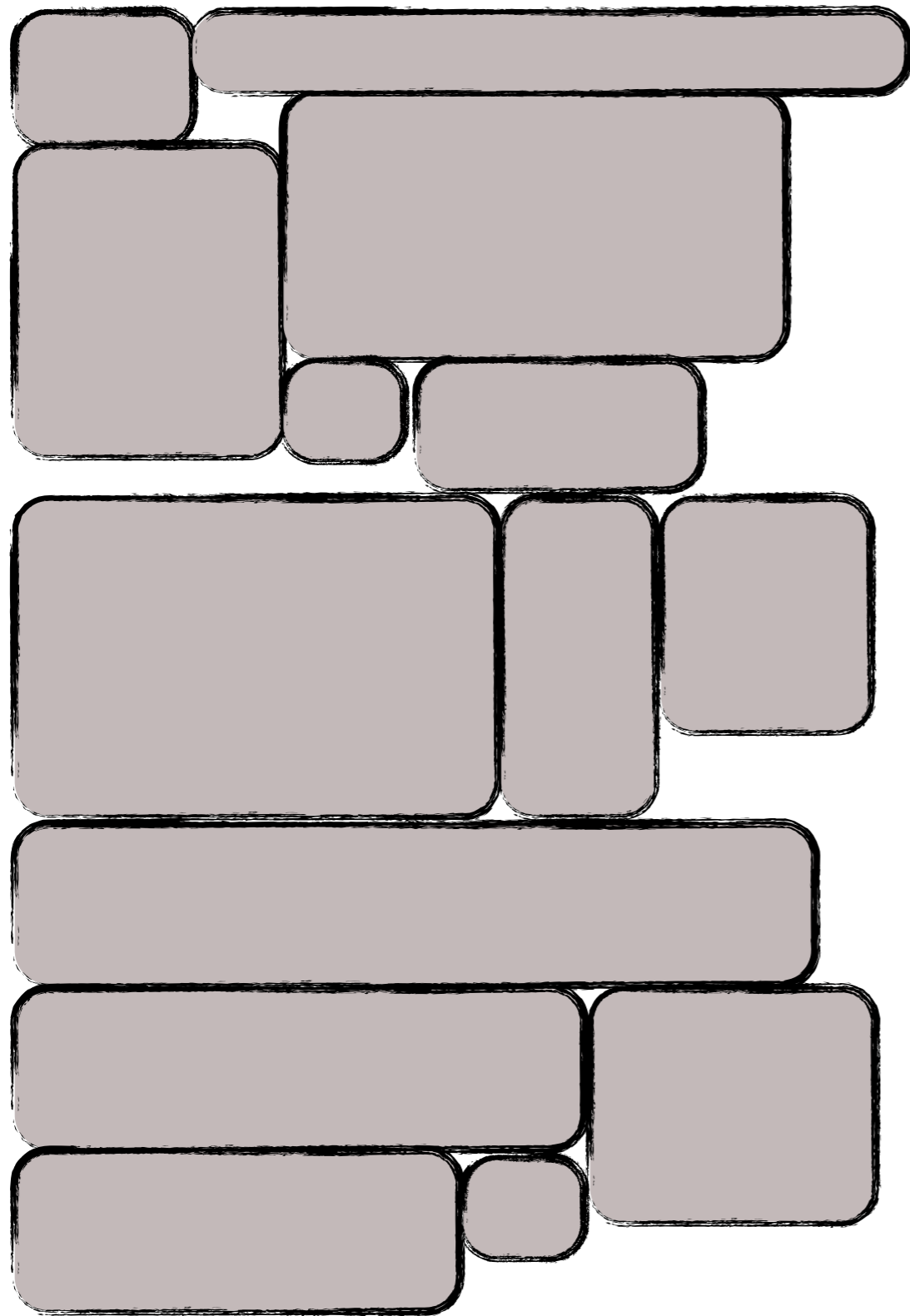
Time



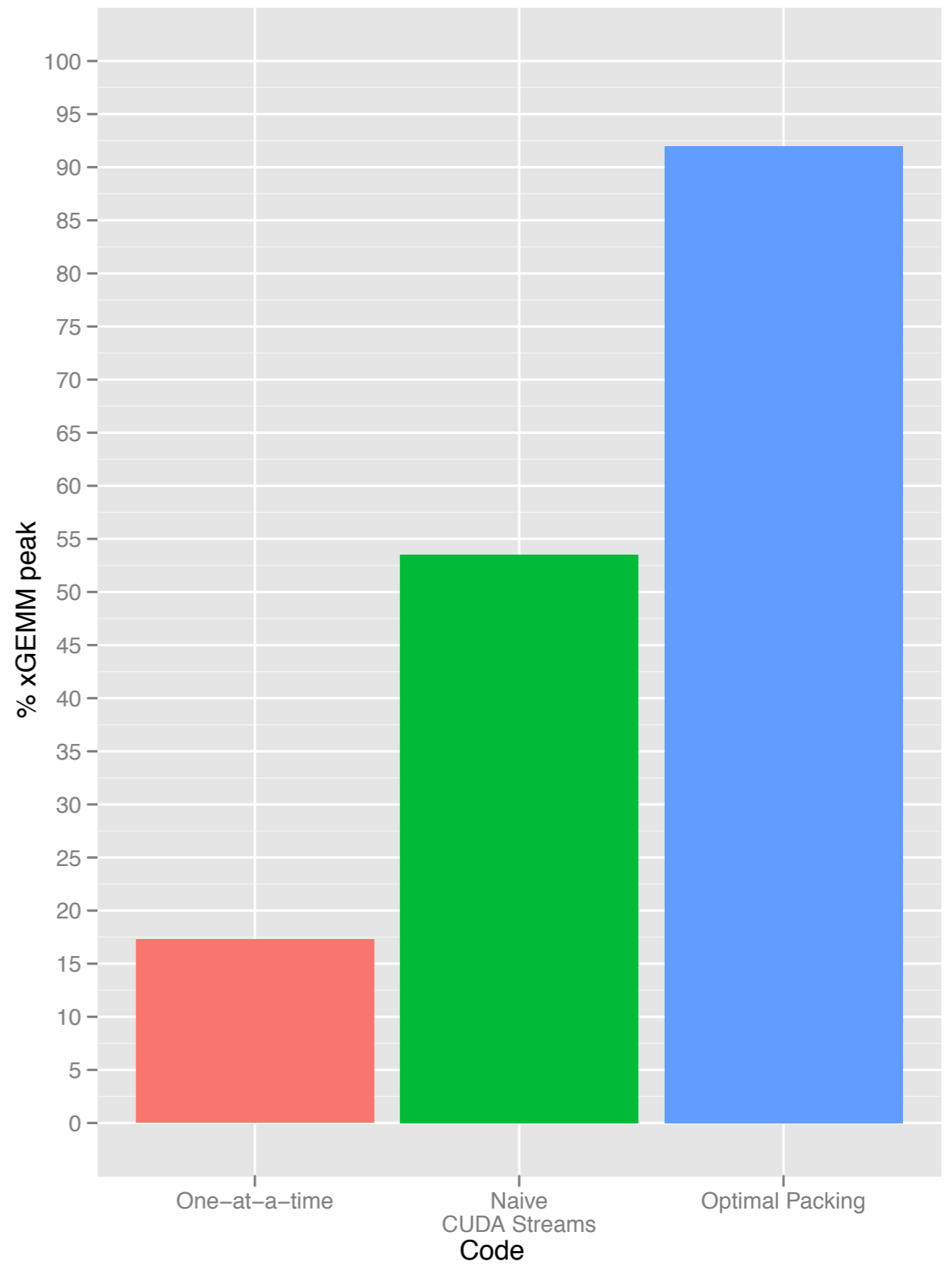
Time



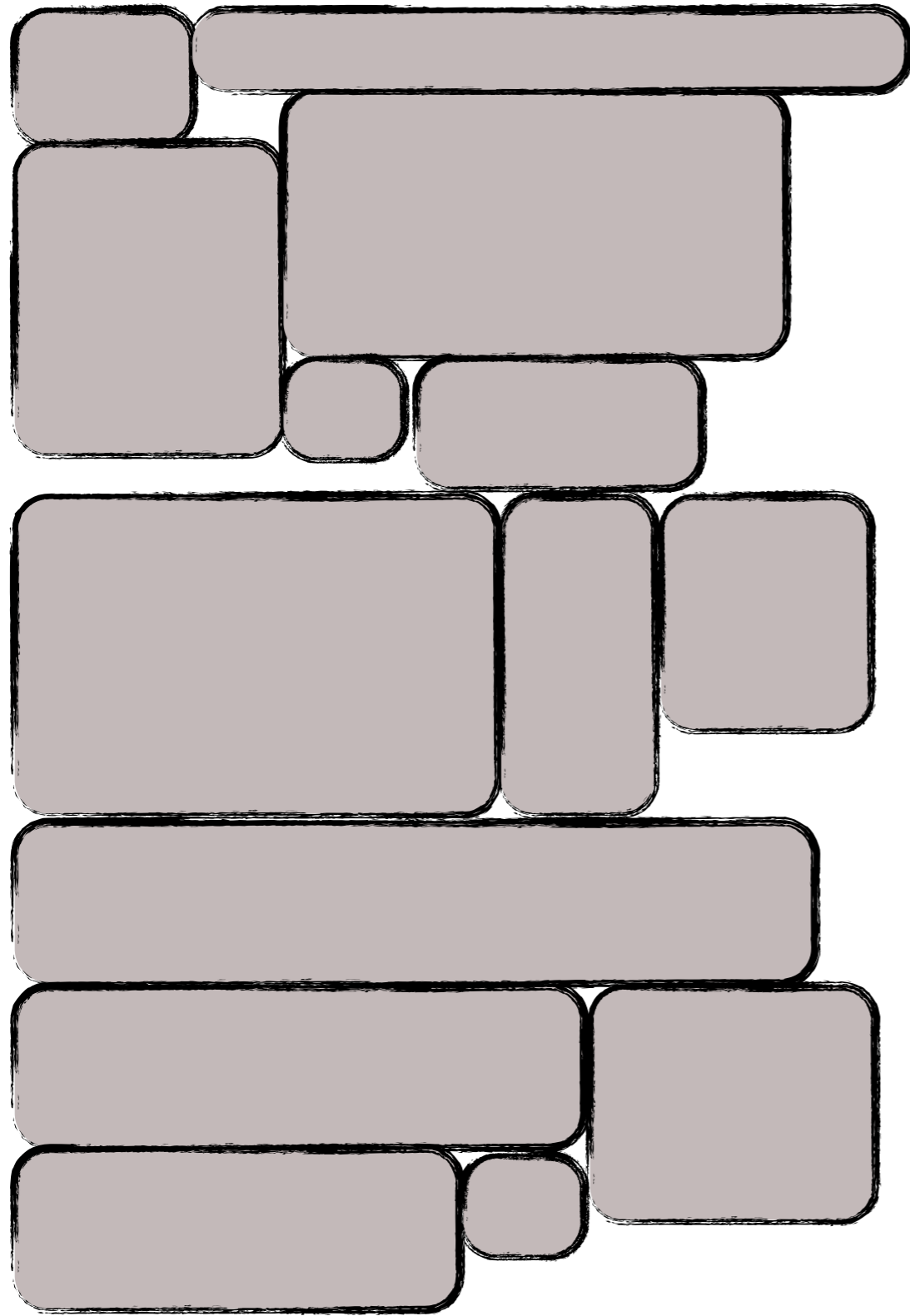
Time



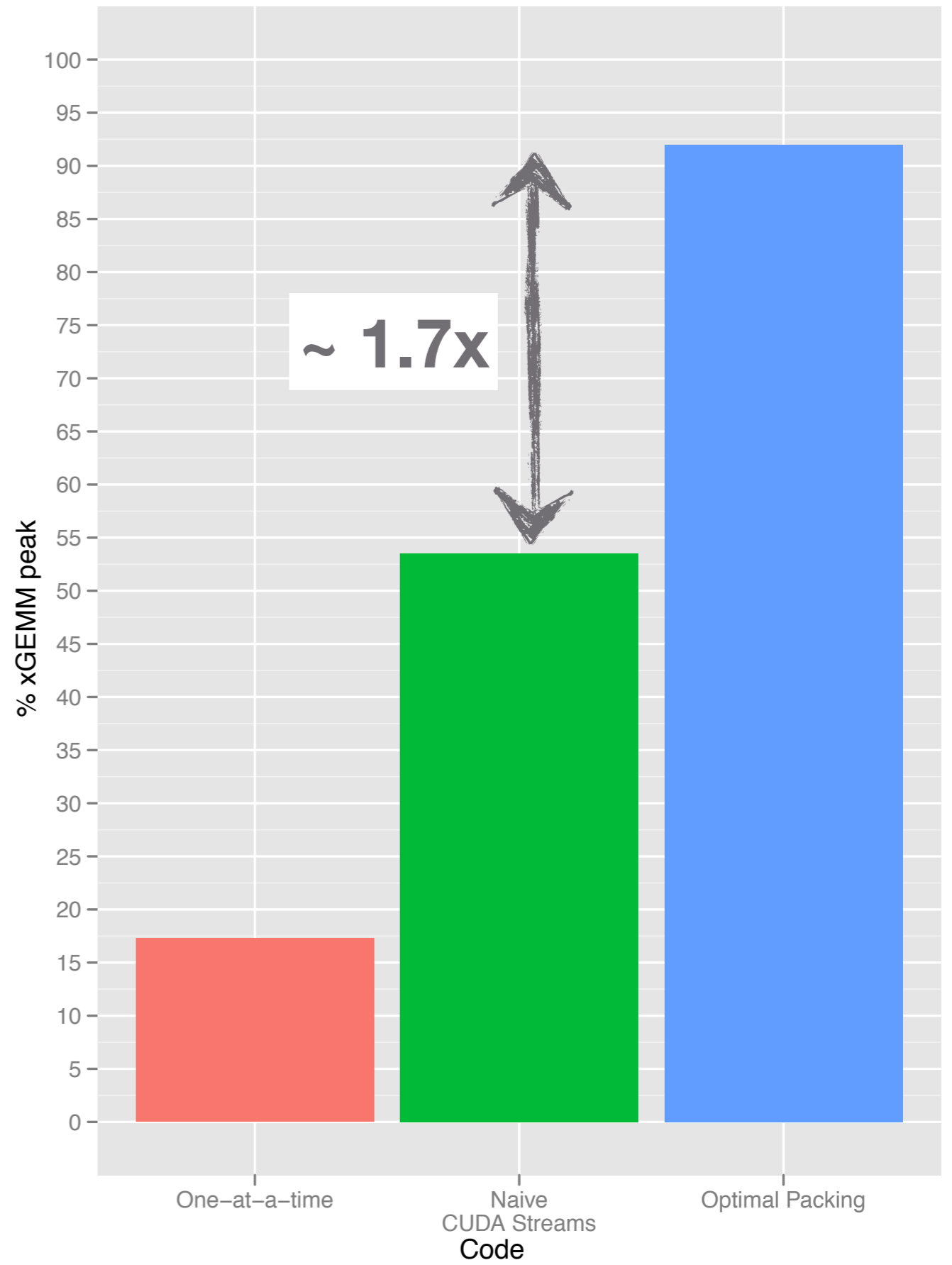
← **Thread capacity** →



Time

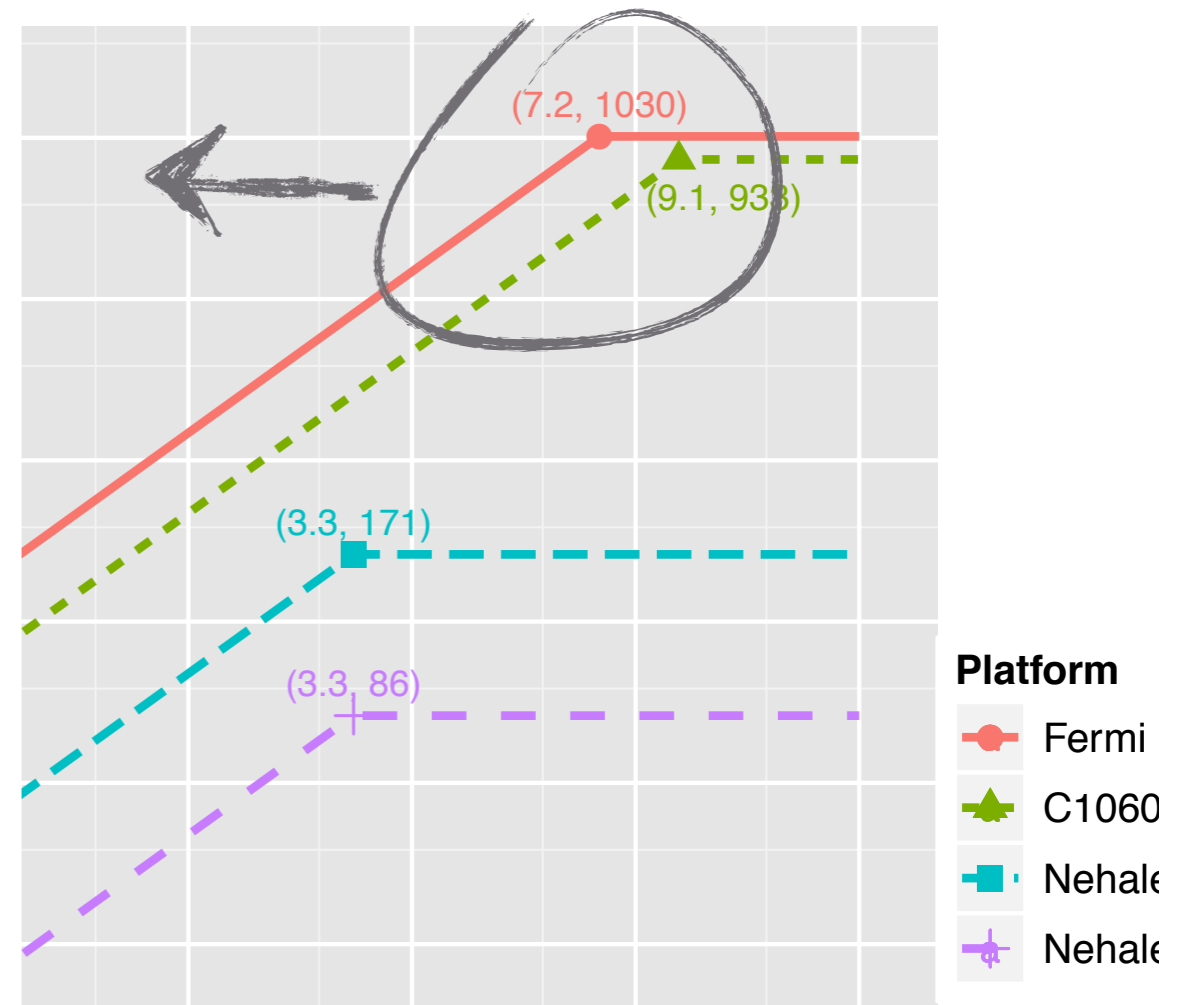


← Thread capacity →

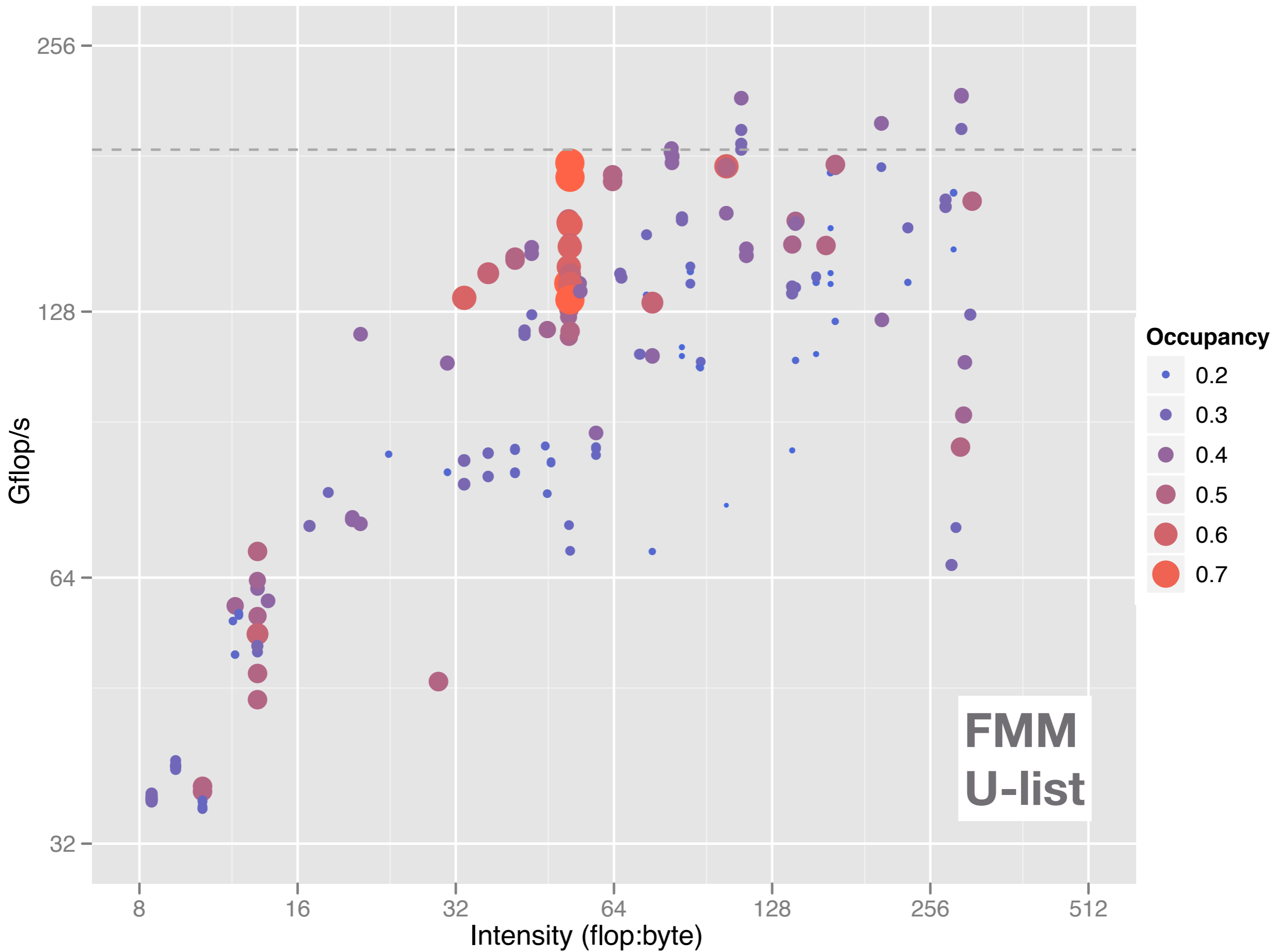


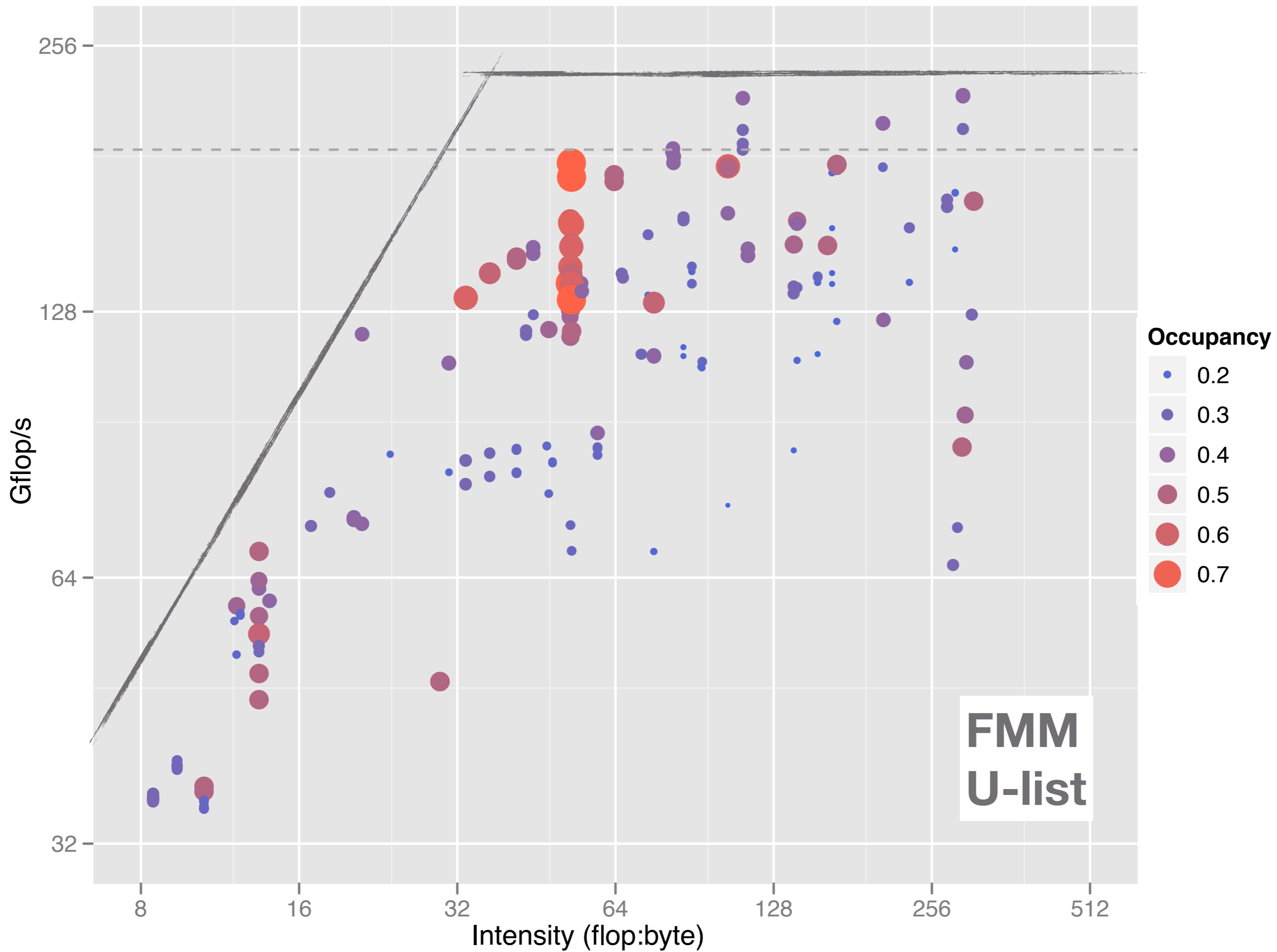
Some concluding questions...

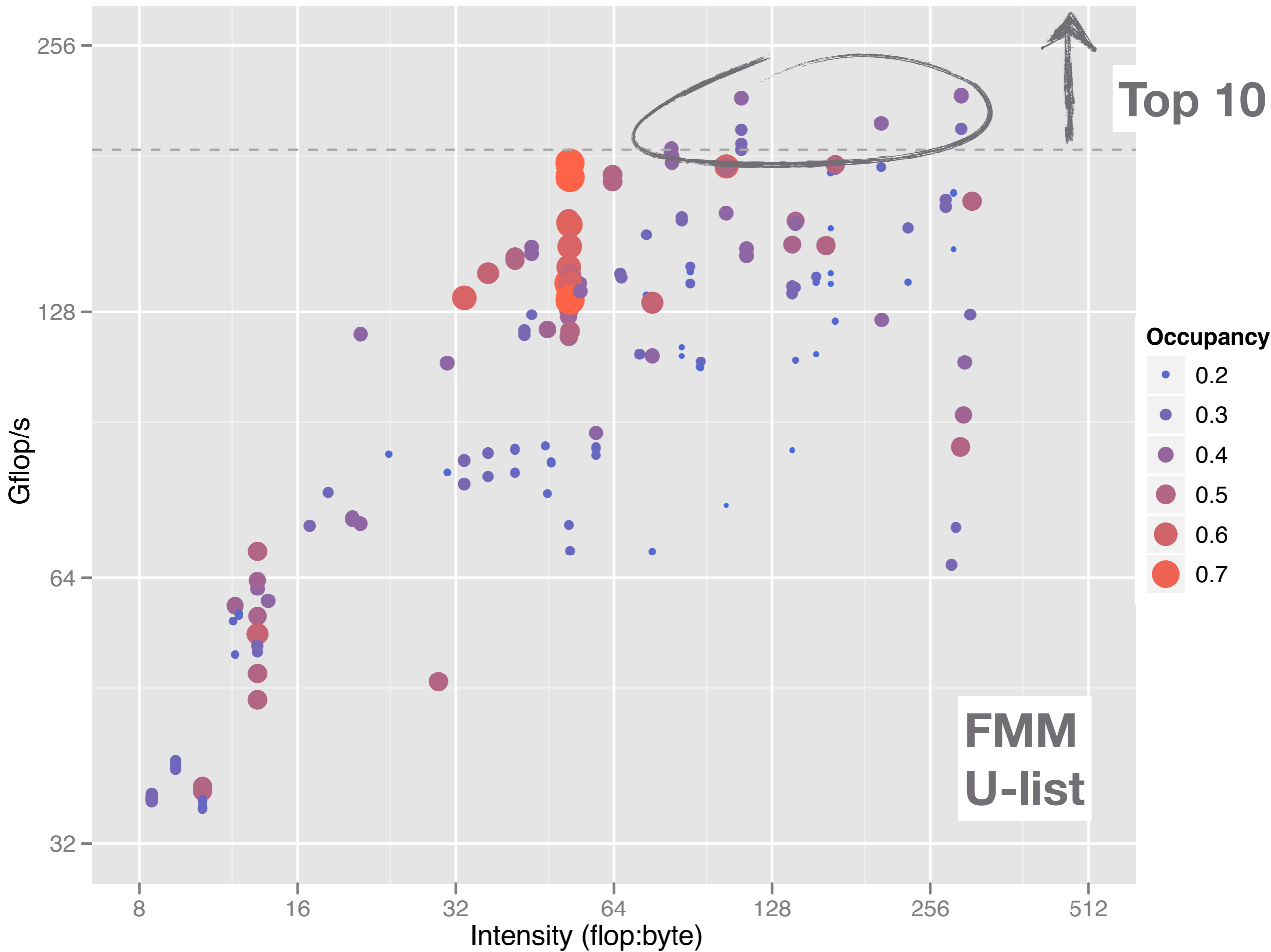
- * What is the right way to think about **opportunity cost**?
- * What are the **performance principles** for reducing tuning effort?
- * What **applications** will lead to better designs?

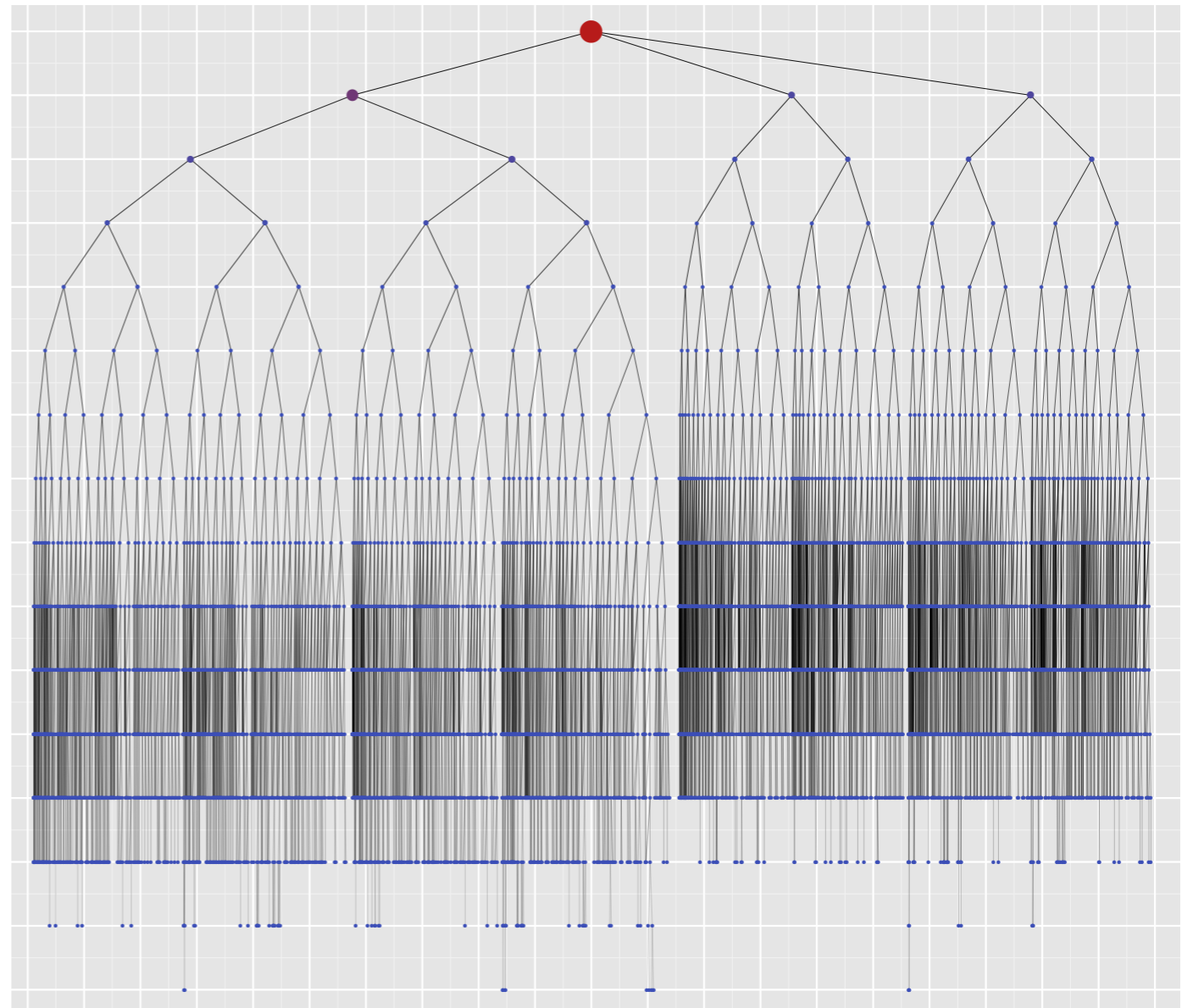
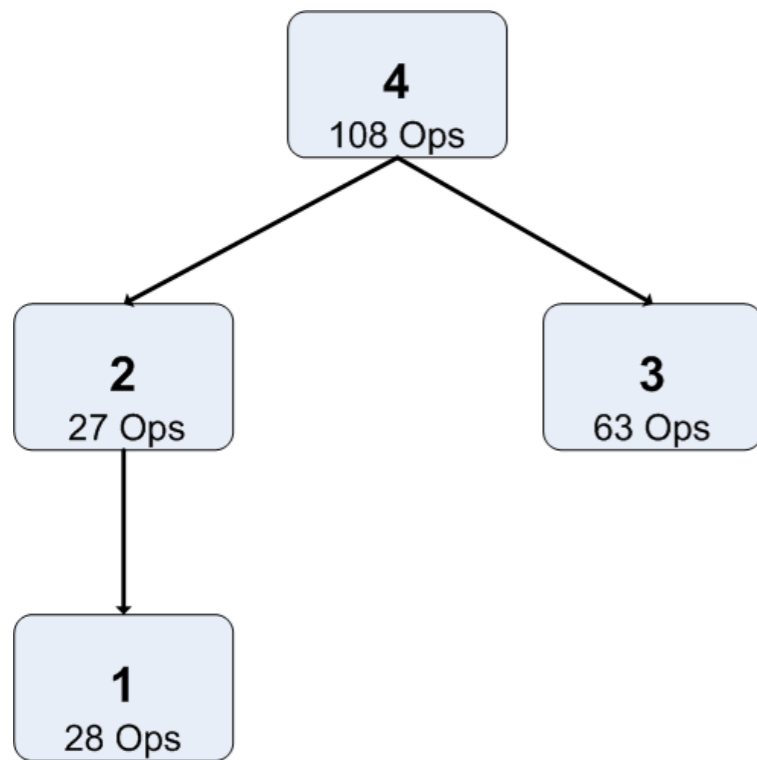


▶ Backup slides



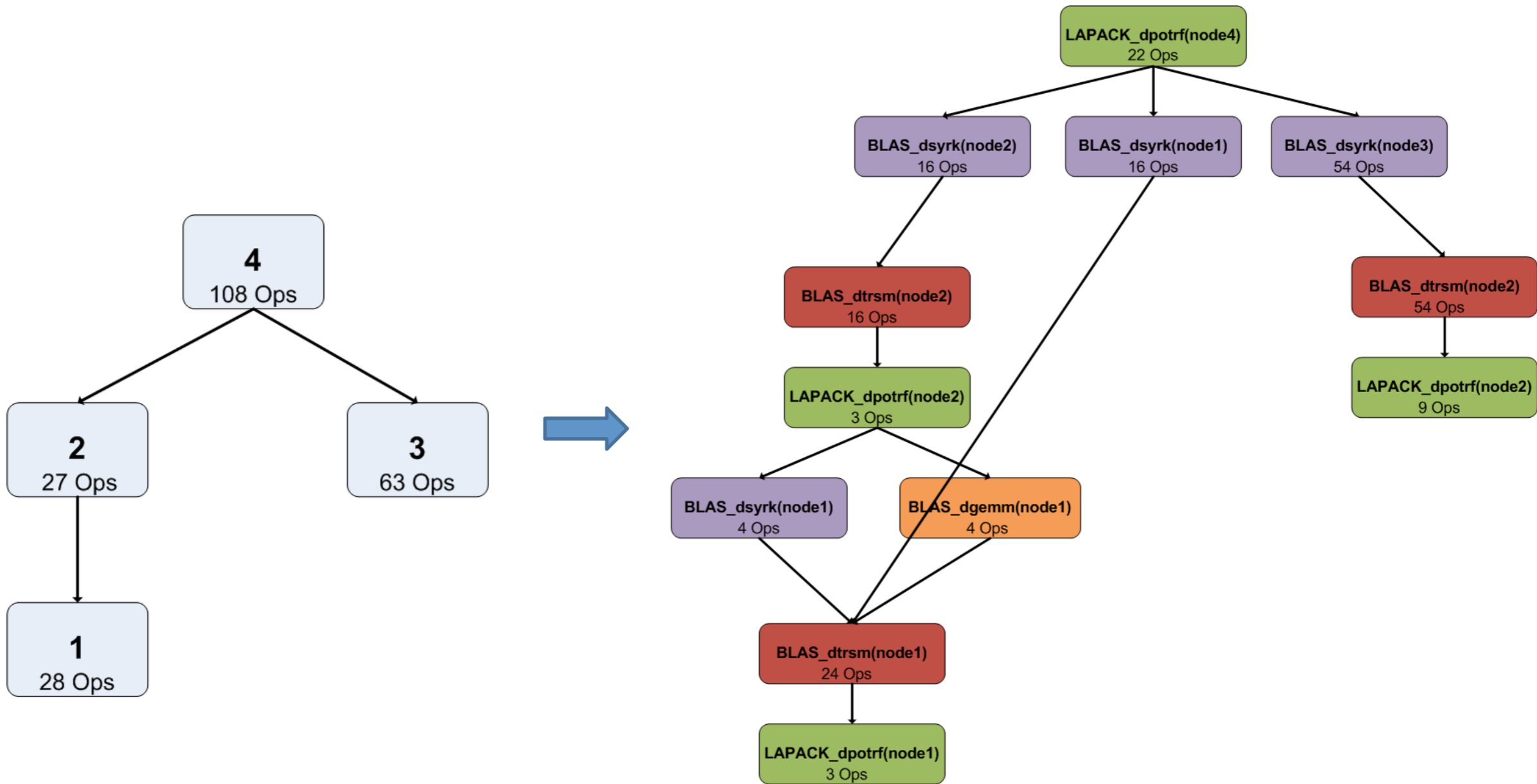






Parallelism = Elimination tree

Independent subtrees may be processed in parallel.



Finer-grained dependencies

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