

Follow along:
hpcgarage.org/ppam13

How much (execution) time, energy, and power will my algorithm cost?

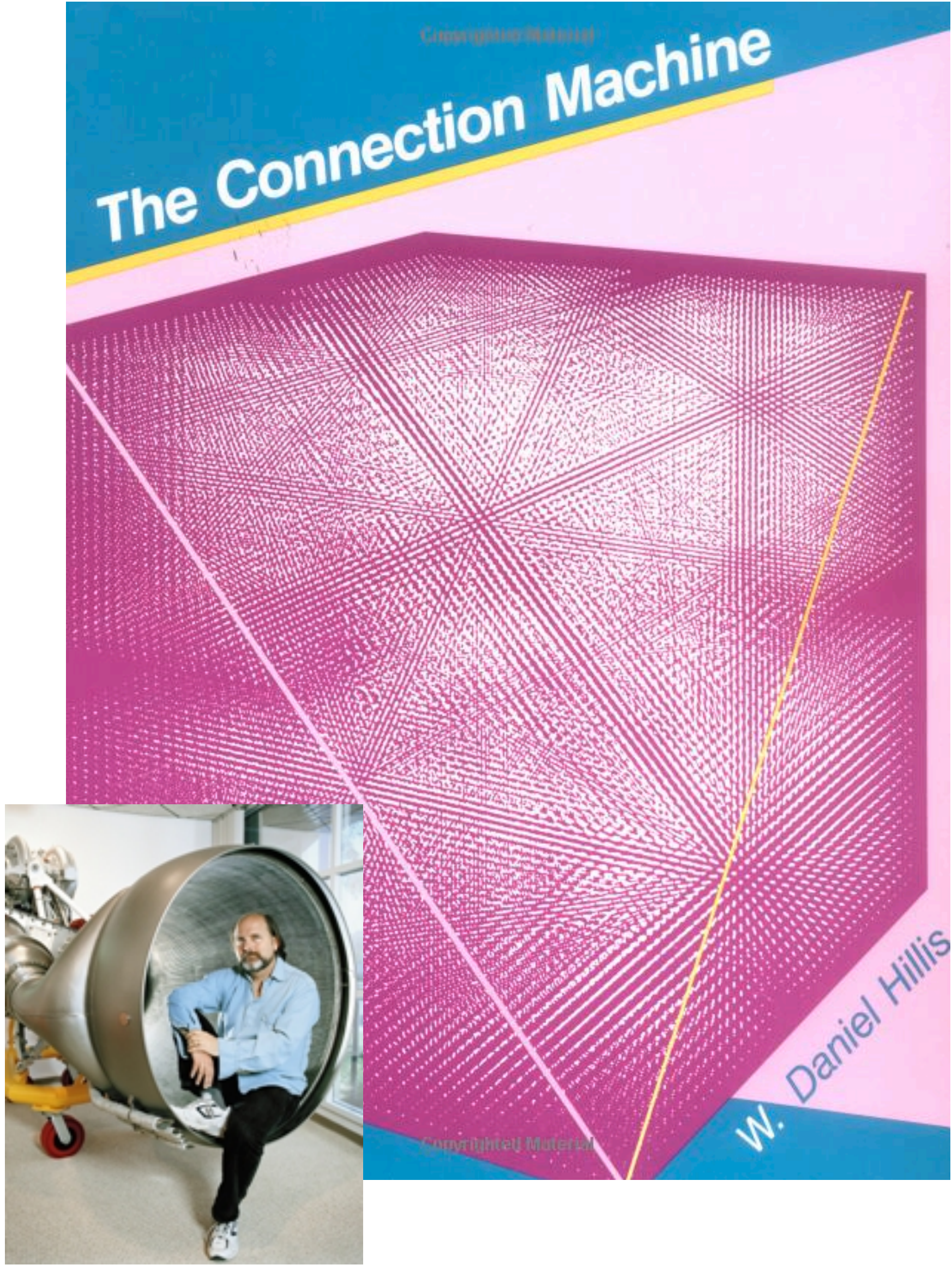
Aparna Chandramowliswaran [MIT] · **Jee Whan Choi** · **Kenneth (Kent) Czechowski** · Marat Dukhan · **Richard (Rich) Vuduc**
· Casey Battaglino · Danny Browne [GTRI] · Cong Hou [Google] · David (Dave) S. Noble, Jr. · Piyush Kumar Sao

September 11, 2013 – Int'l. Conf. Parallel Processing & Applied Mathematics – PPAM



ACM Doctoral Dissertation Award Winner (1985)

Follow along:
hpcgarage.org/ppam13



http://www.amazon.com/Connection-Machine-Artificial-Intelligence/dp/0262580977/ref=sr_1_1?ie=UTF8&qid=1319571691&sr=8-1
http://www.technologyreview.com/files/1158/1106_Q-A.tif.jpg



6.10 Notes 135

7 New Computer Architectures and Their Relationship to Physics

or, Why Computer Science is No Good 137

7.1 Why Computer Science is No Good 137

7.2 Connection Machine Physics 139

7.3 New Hope for a Science of Computation 142

7.4 Notes 144

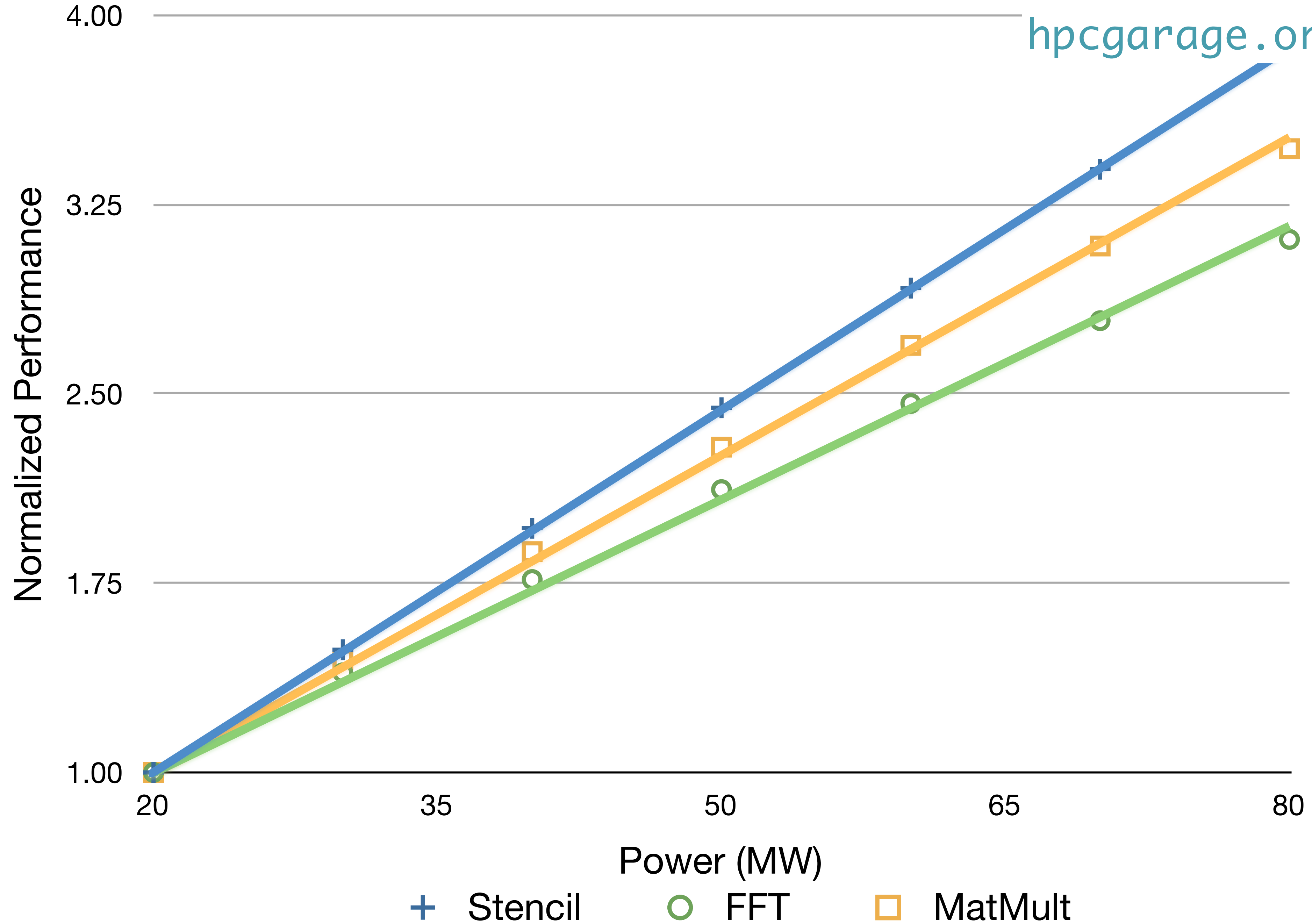
Annotated Bibliography 145

http://www.amazon.com/Connection-Machine-Artificial-Intelligence/dp/0262580977/ref=sr_1_1?ie=UTF8&qid=1319571691&sr=8-1
http://www.technologyreview.com/files/1158/1106_Q-A.tif.jpg

[Preview] Redressing Hillis? Algorithmic power scaling

Follow along:

hpcgarage.org/ppam13





Aparna Chandramowliswaran
 [Now R.S. @ MIT]
Fast multipole method



Jee Whan Choi
*Autotuning for
 power and energy*



Kent Czechowski
Co-design



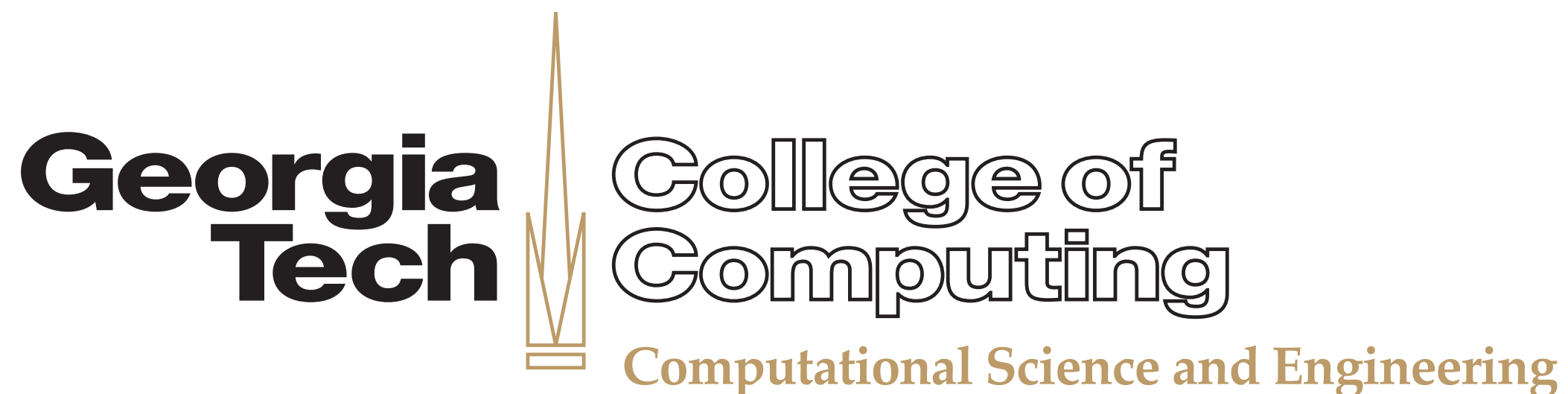
Marat Dukhan
*Math libraries
 & machine learning*

2010 Gordon Bell Prize (+ G. Biros)
 2010 IPDPS Best Paper (+ K. Knobe, Intel CnC lead)
 2012 SIAM Data Mining Best Paper (D. Lee [GE Research] + A. Gray)

See our recent 2013 IPDPS papers, posted at: hpcgarage.org/ppam13

* Jee Choi — “A roofline model of energy.”

* Kent Czechowski — “A theoretical framework for algorithm-architecture co-design.”



Rebutting Hillis, Part 1:

First principles:

Rooflines in **time**,
arch lines in **energy**,
and “**power** lines”



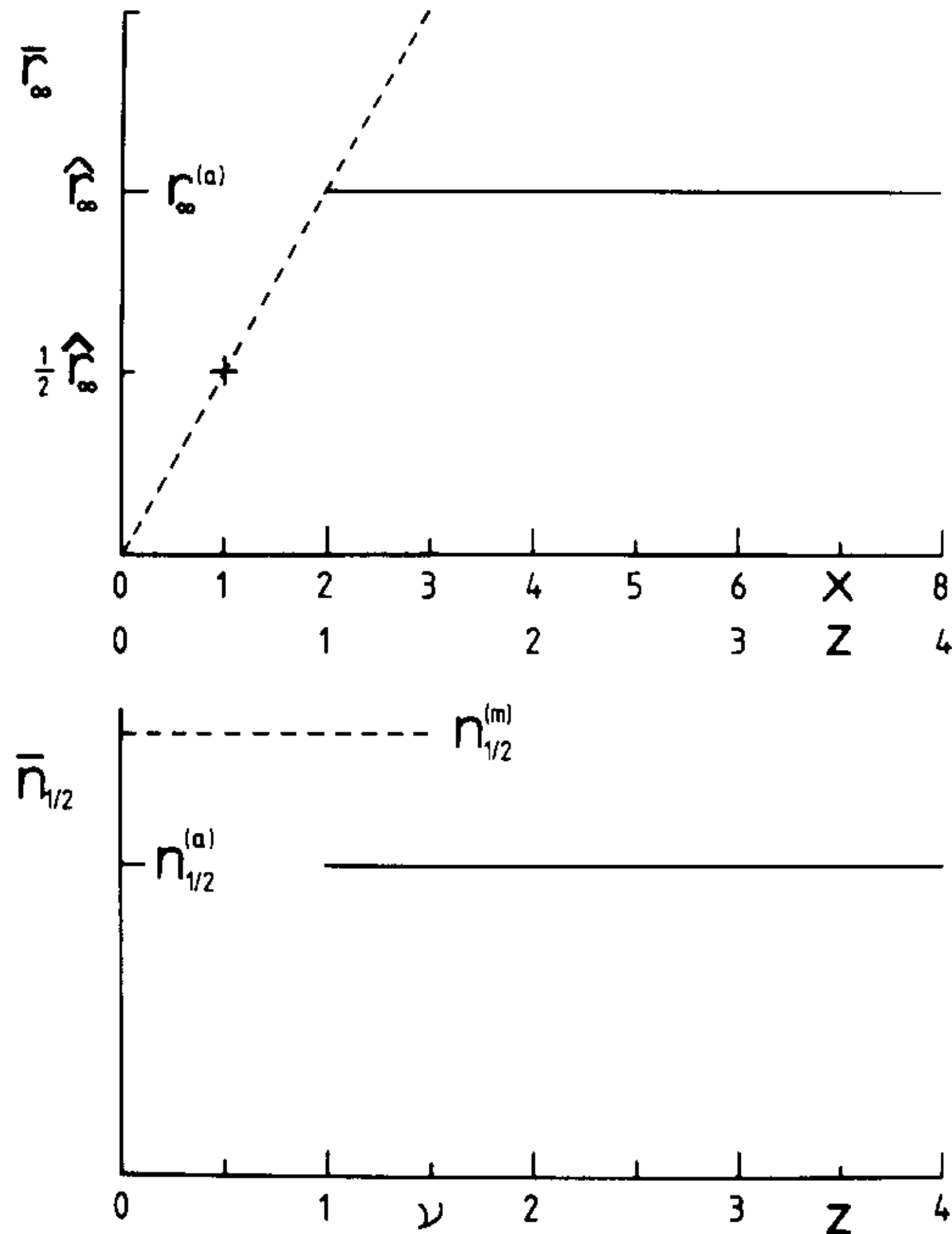
Jee



Marat

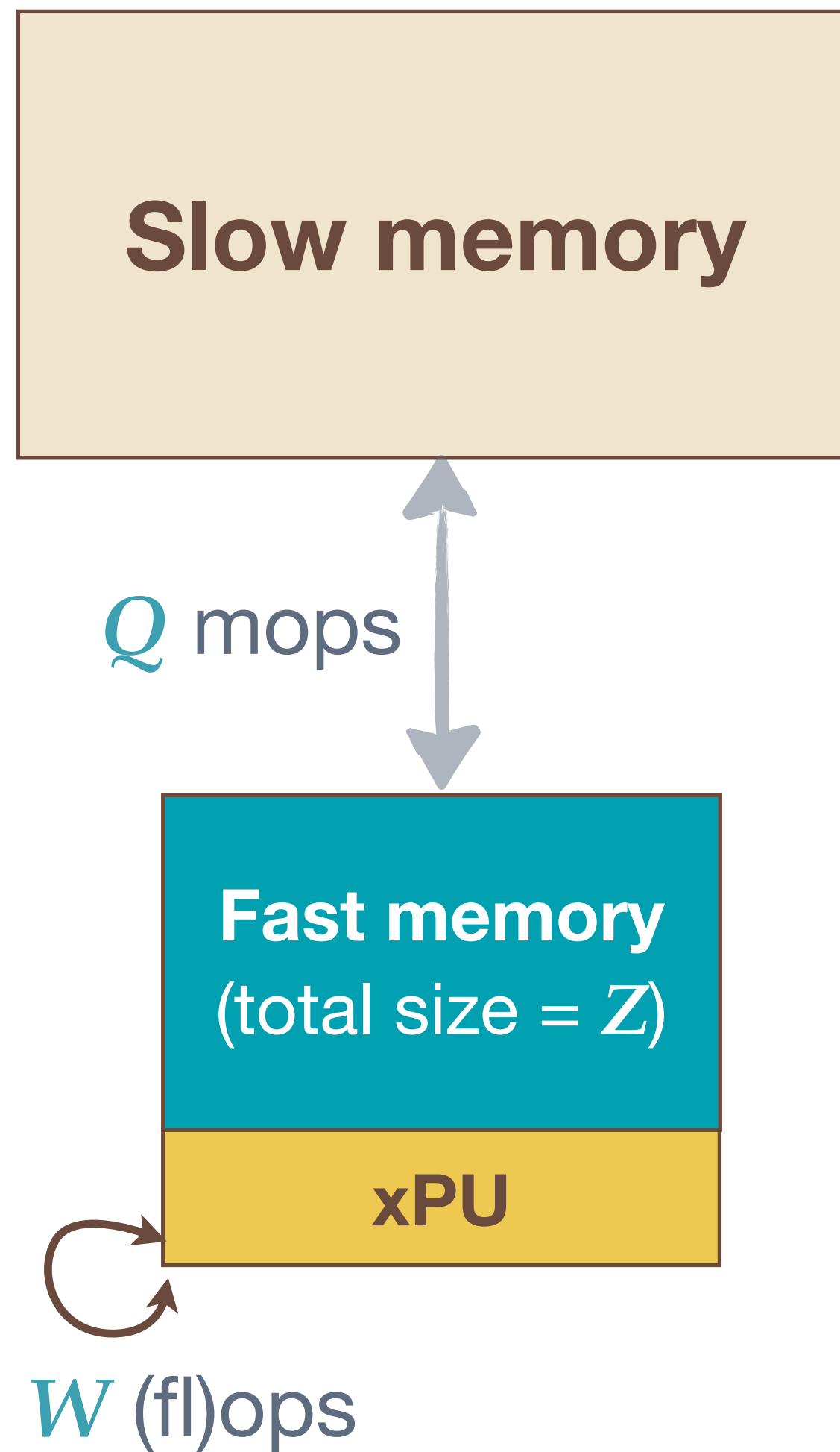
Energy and power analogues of the time-based “roofline model” of Hockney & Curington (1989) and Williams et al. (2009)

R.W. Hockney, I.J. Curington / $f_{1/2}$: A parameter to characterize bottlenecks



R.W. Hockney and I.J. Curington (1989).
 “ $f_{1/2}$: A parameter to characterize memory and communication bottlenecks.”
Parallel Computing, **10**(3), 277–286.
 doi: [10.1016/0167-8191\(89\)90100-2](https://doi.org/10.1016/0167-8191(89)90100-2)

Fig. 2. The variation of $(\bar{r}_\infty, \bar{n}_{1/2})$ with f for the case of a combined memory I/O and arithmetic pipeline, when the I/O and arithmetic can be overlapped. Full lines: parameters when arithmetic dominates, equations (14b); dotted lines: parameters when I/O dominates, equations (13b). Notation as Fig. 1, and $\nu = 1.5$.

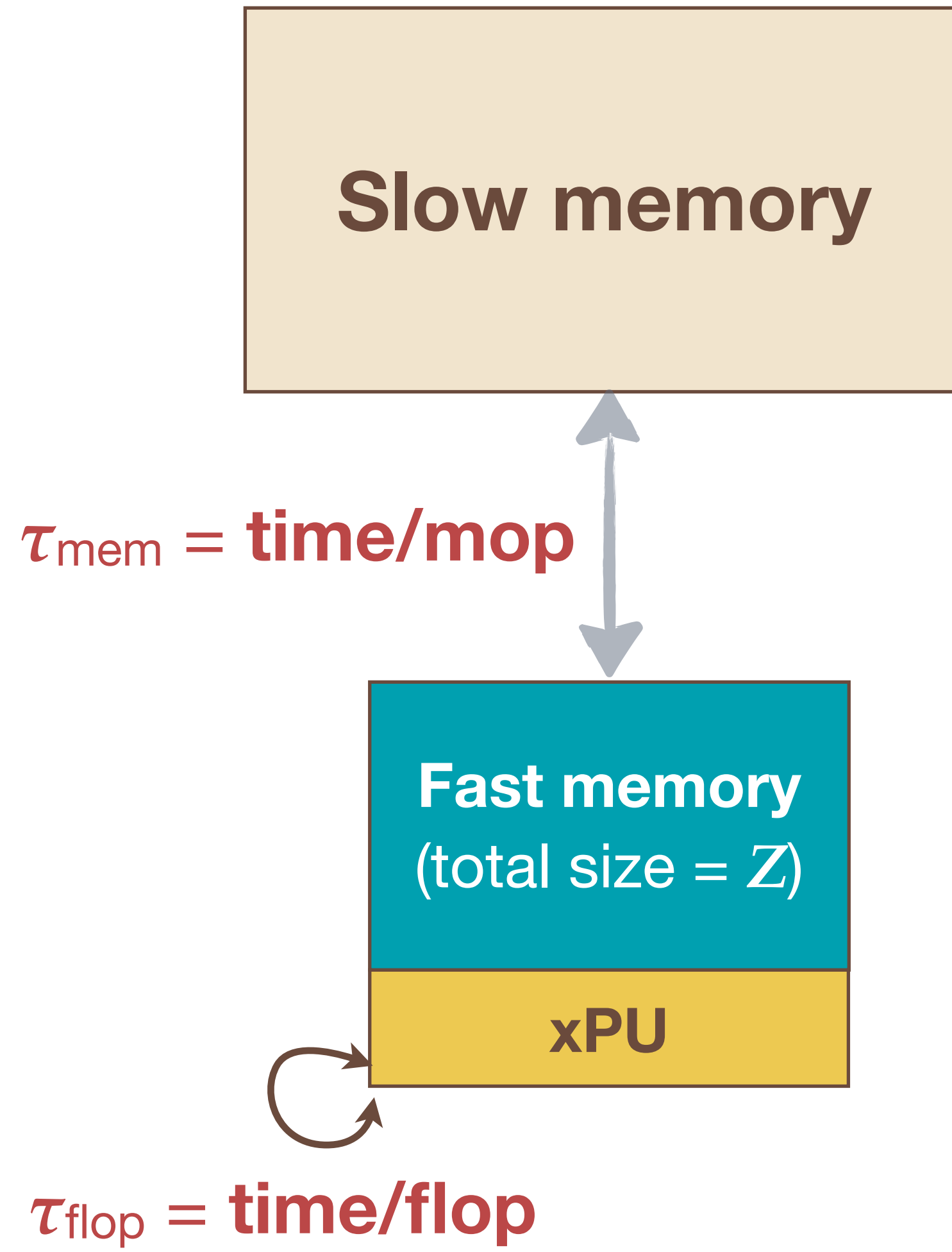


W \equiv # (fl)ops

Q \equiv # mem. ops (mops)

$= Q(Z)$

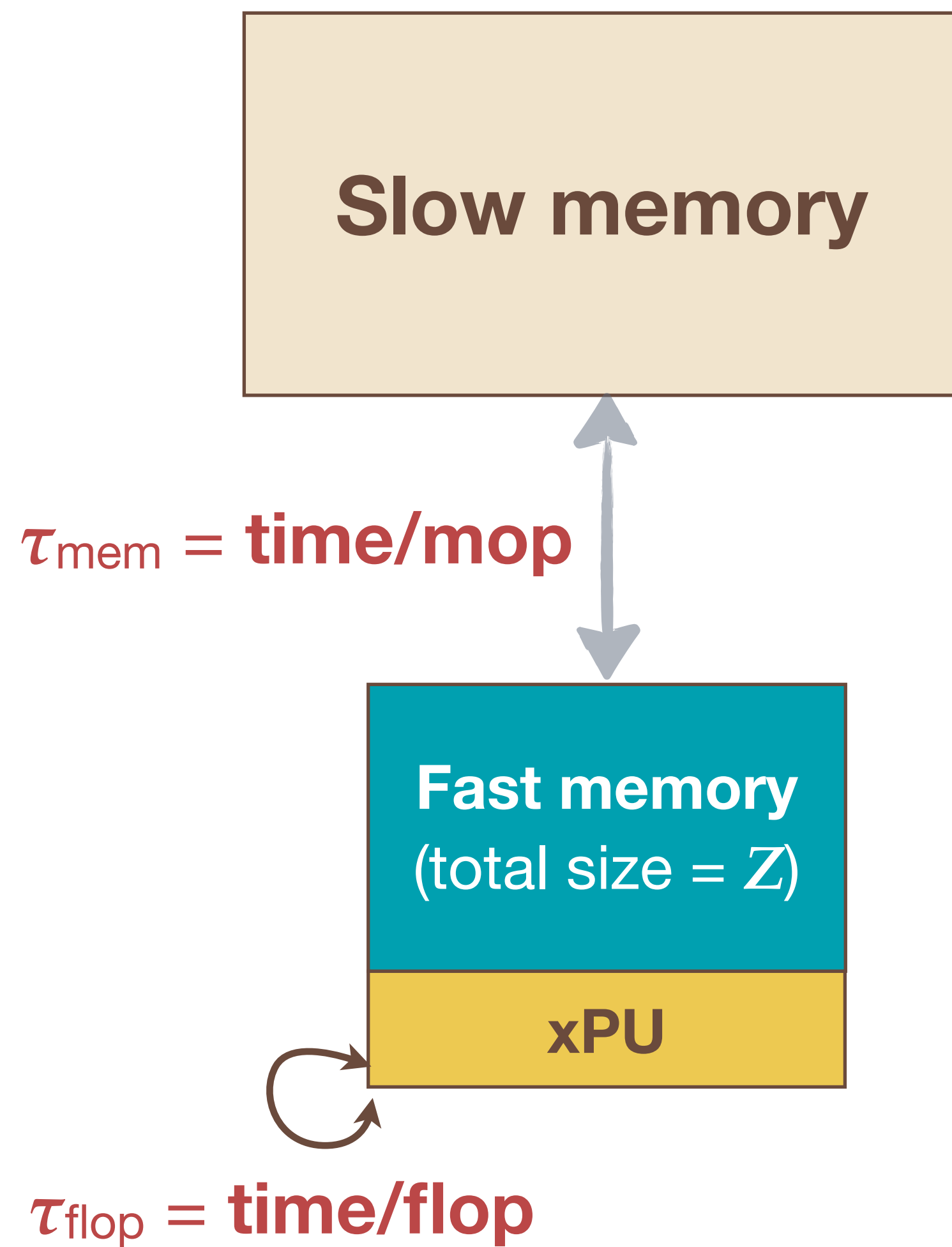
von Neumann bottleneck



$$T = \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}})$$

von Neumann bottleneck

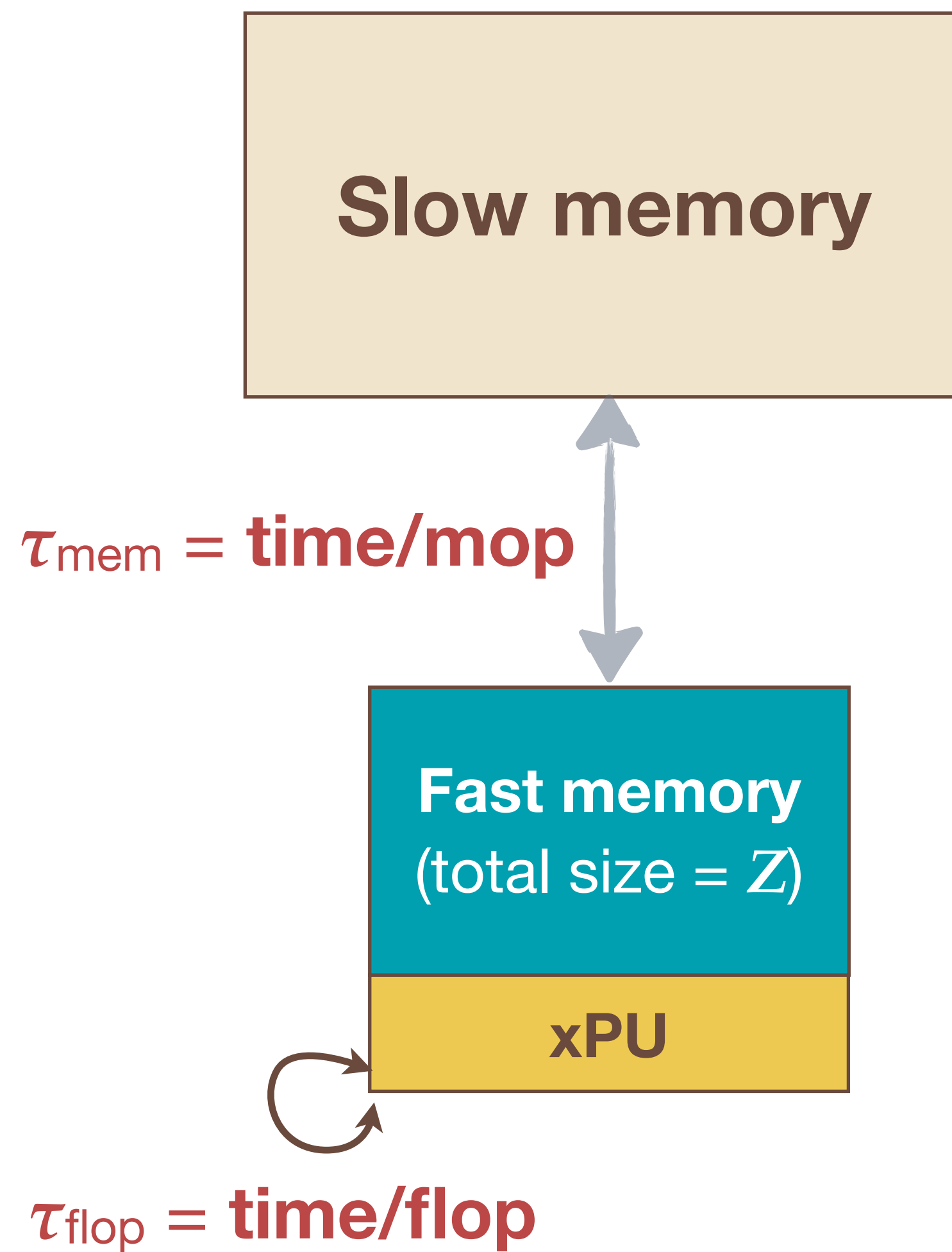
Balance analysis — Kung (1986); Hockney & Curington (1989); Bleloch (1994); McCalpin (1995); Williams et al. (2009); Czechowski et al. (2011); ...



$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right)
 \end{aligned}$$

von Neumann bottleneck

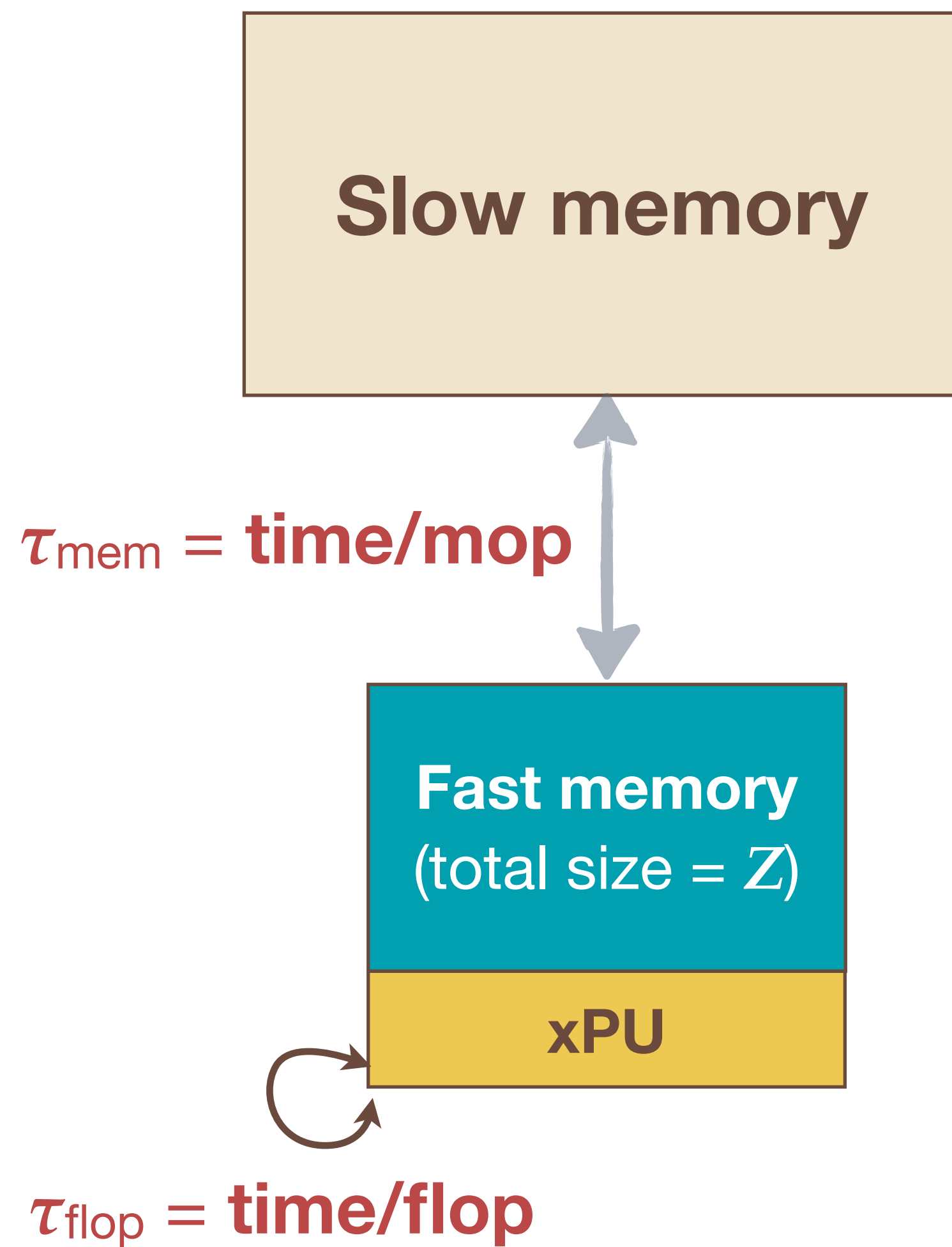
Balance analysis — Kung (1986); Hockney & Curington (1989); Bleloch (1994); McCalpin (1995); Williams et al. (2009); Czechowski et al. (2011); ...



$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

von Neumann bottleneck

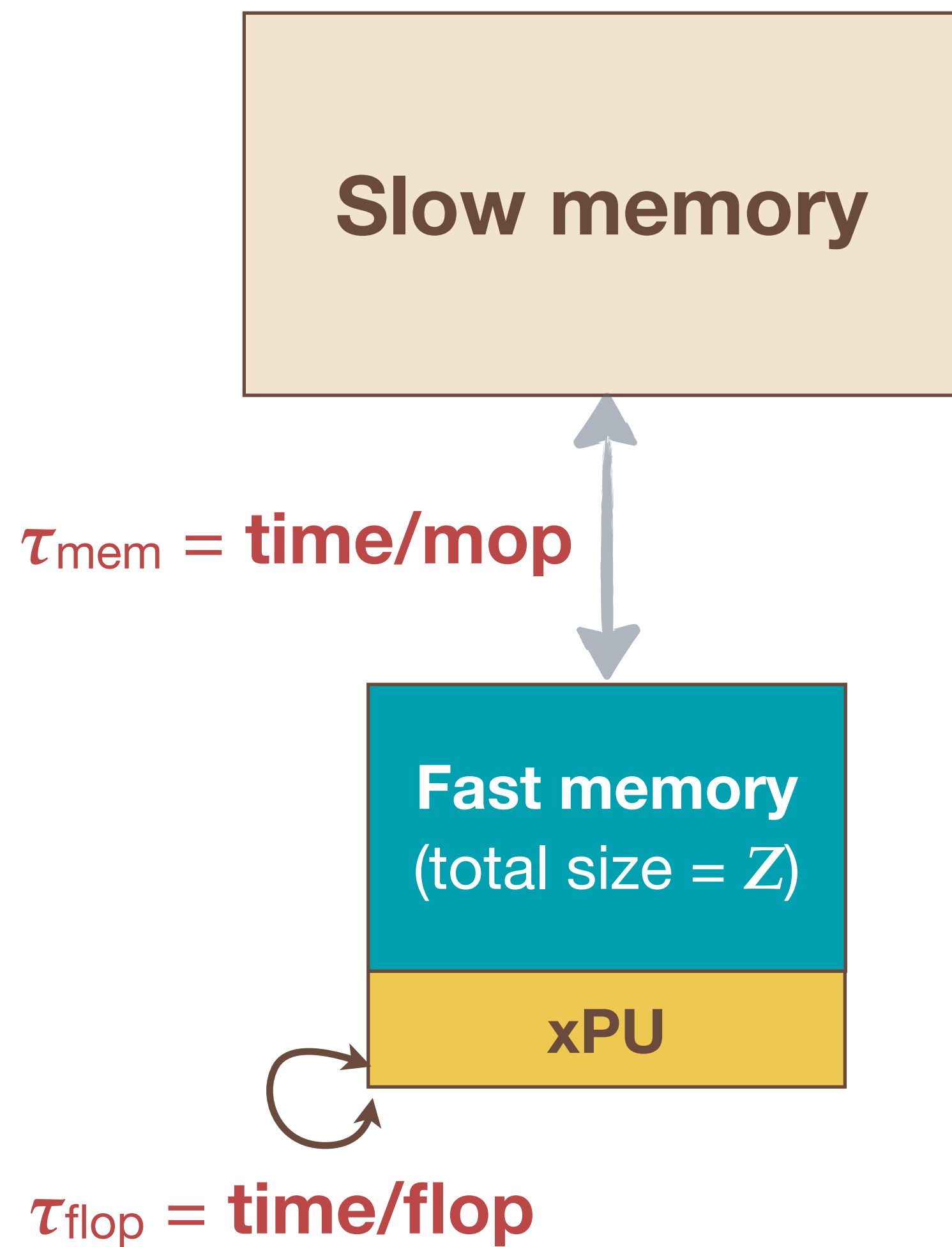
Balance analysis — Kung (1986); Hockney & Curington (1989); Blelloch (1994); McCalpin (1995); Williams et al. (2009); Czechowski et al. (2011); ...



$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

Minimum time

von Neumann bottleneck

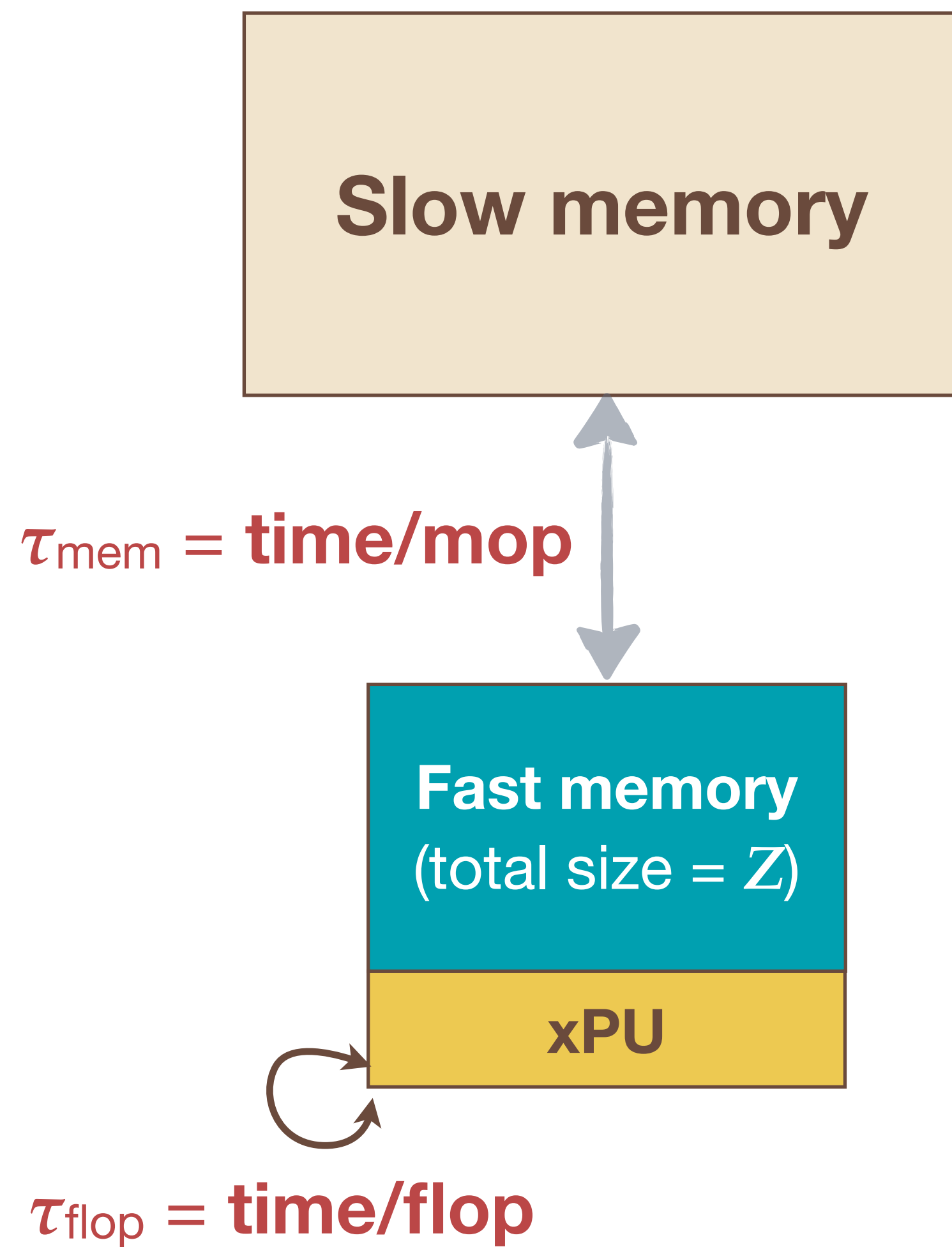


$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

Minimum time

Intensity
(flop : mop)

von Neumann bottleneck



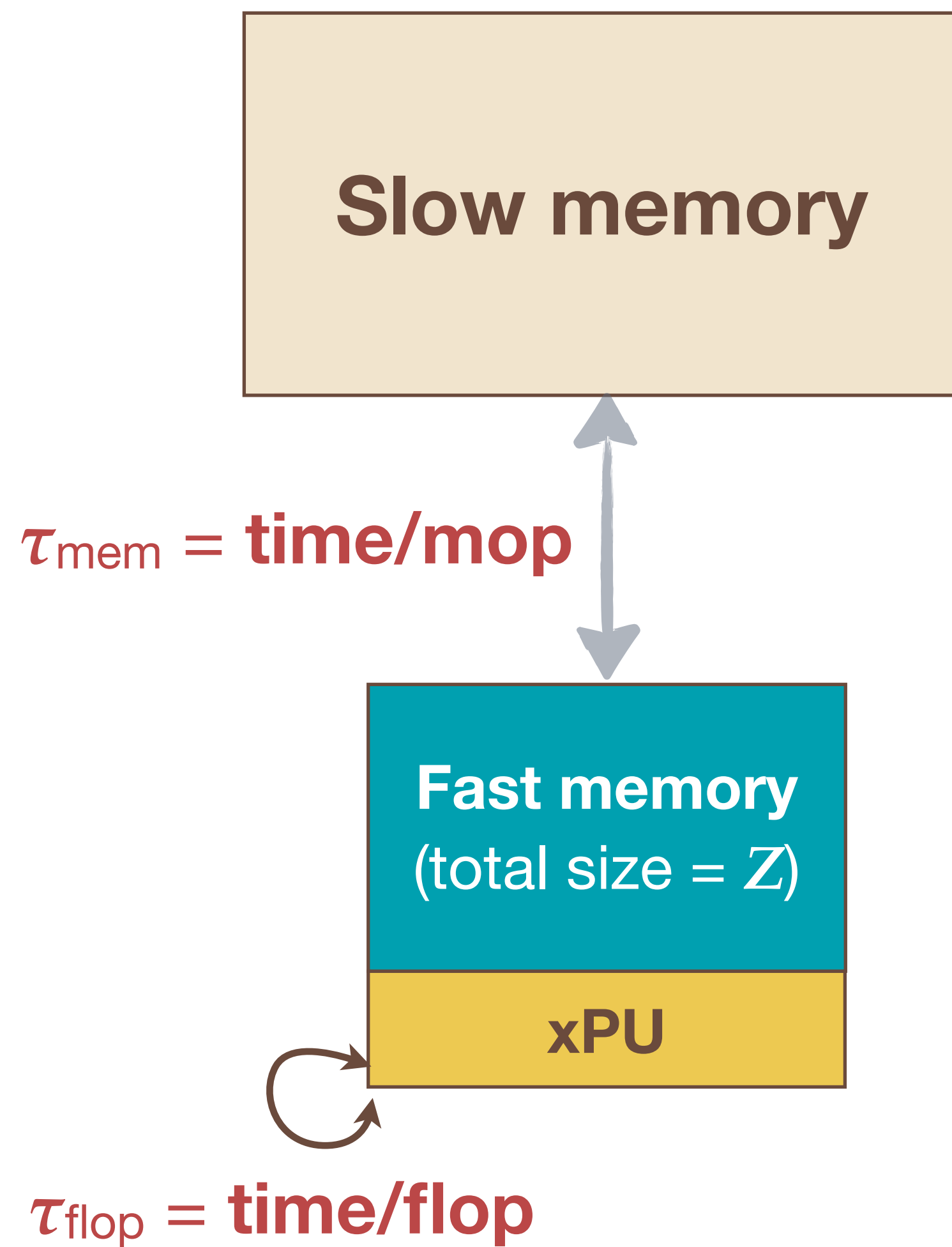
$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

Minimum time

Intensity
(flop : mop)

Balance
(flop : mop)

von Neumann bottleneck



von Neumann bottleneck

$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

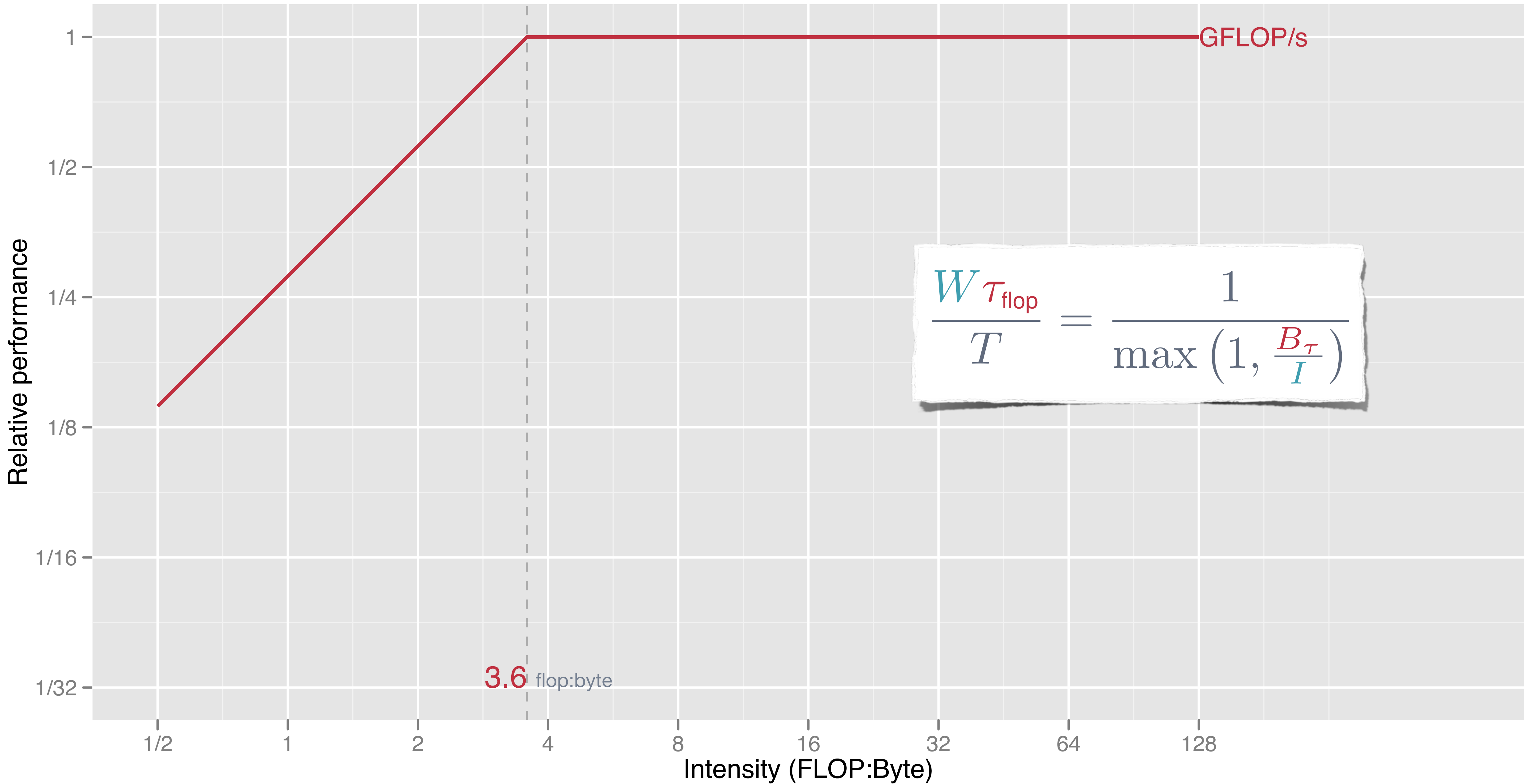
Intensity
(flop : mop)

Balance
(flop : mop)

Consider:

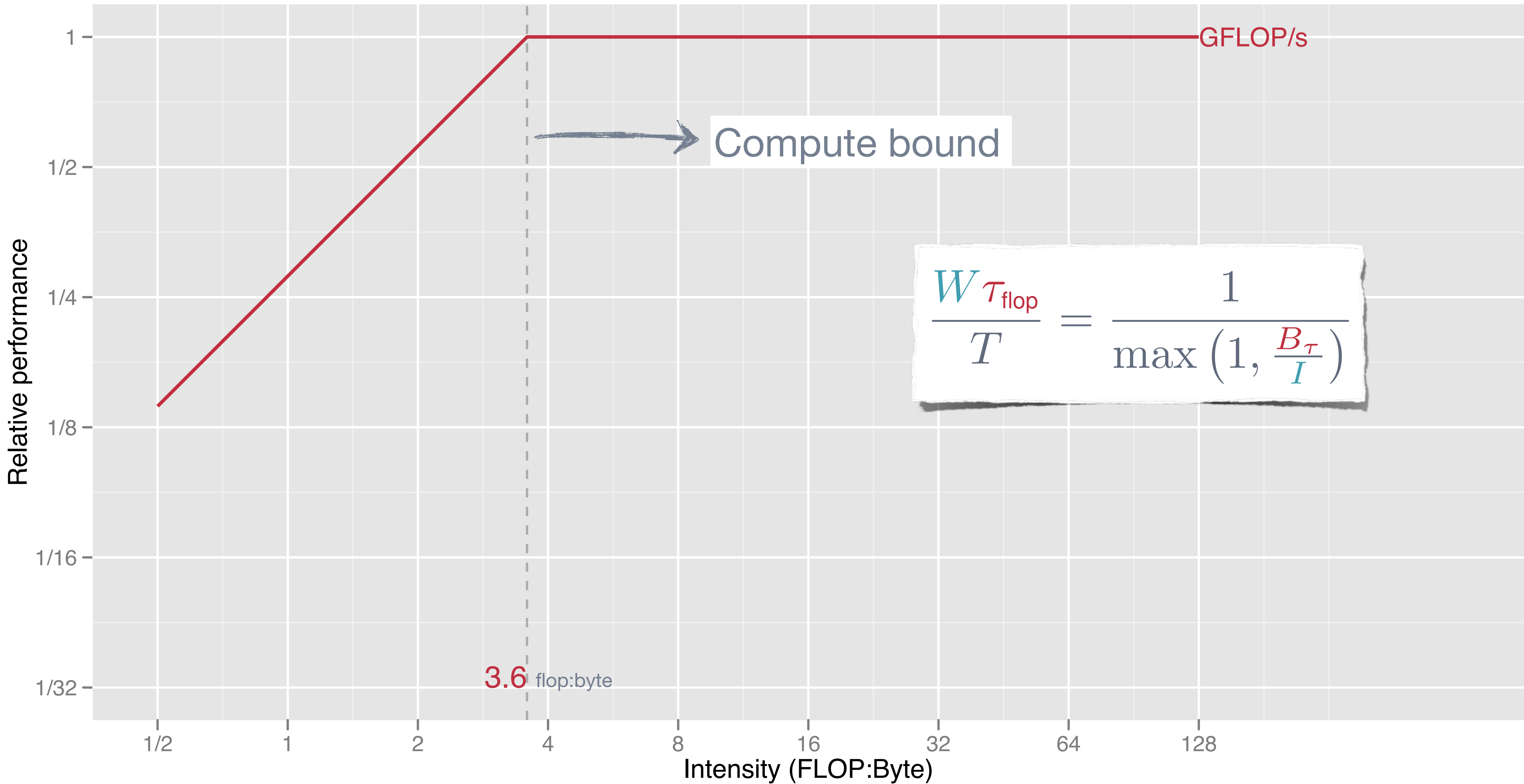
$$\frac{W \tau_{\text{flop}}}{T}$$

“Roofline” — Williams et al. (*Comm. ACM* '09)



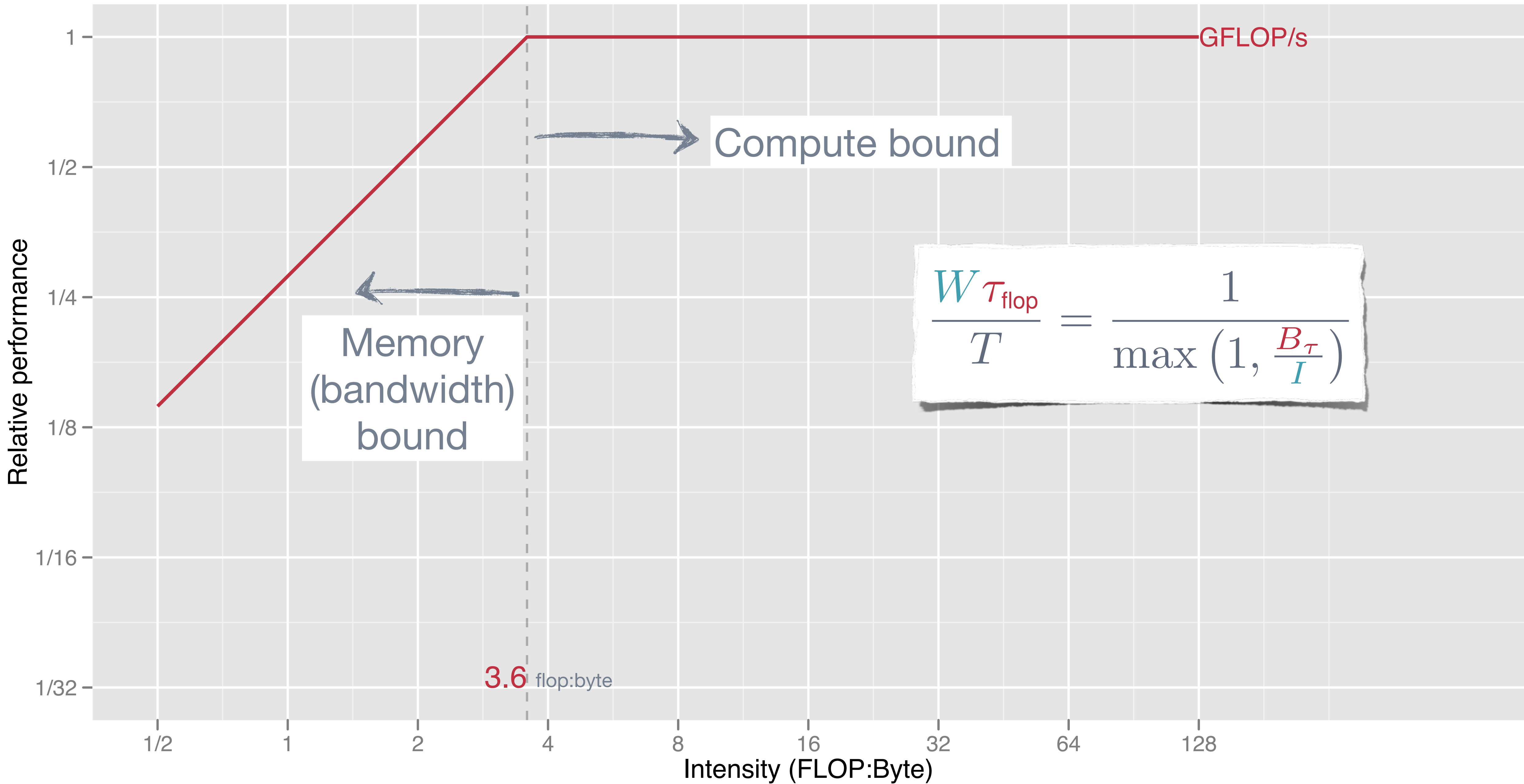
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

“Roofline” — Williams et al. (*Comm. ACM* '09)



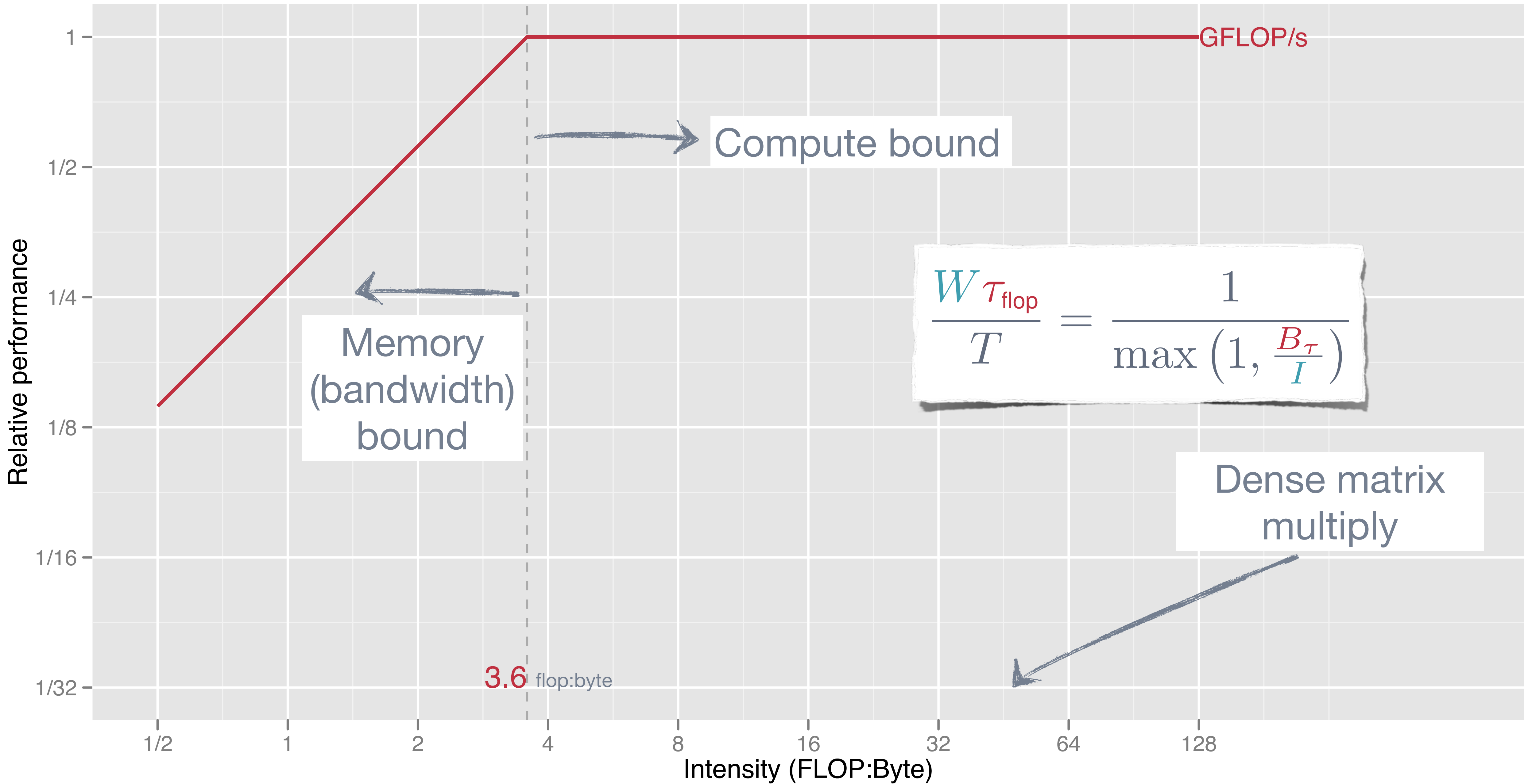
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

“Roofline” — Williams et al. (*Comm. ACM* '09)



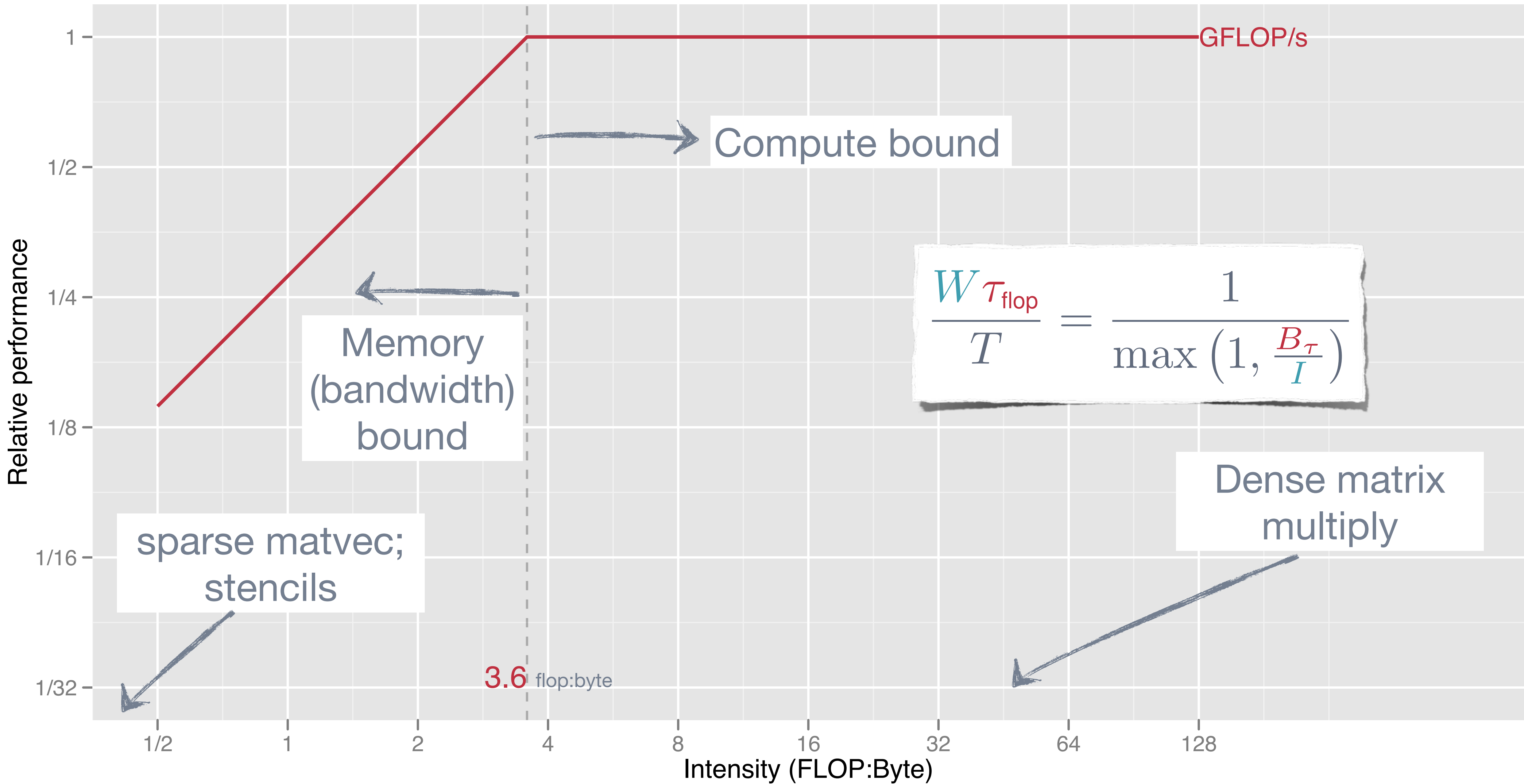
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

“Roofline” — Williams et al. (*Comm. ACM* '09)



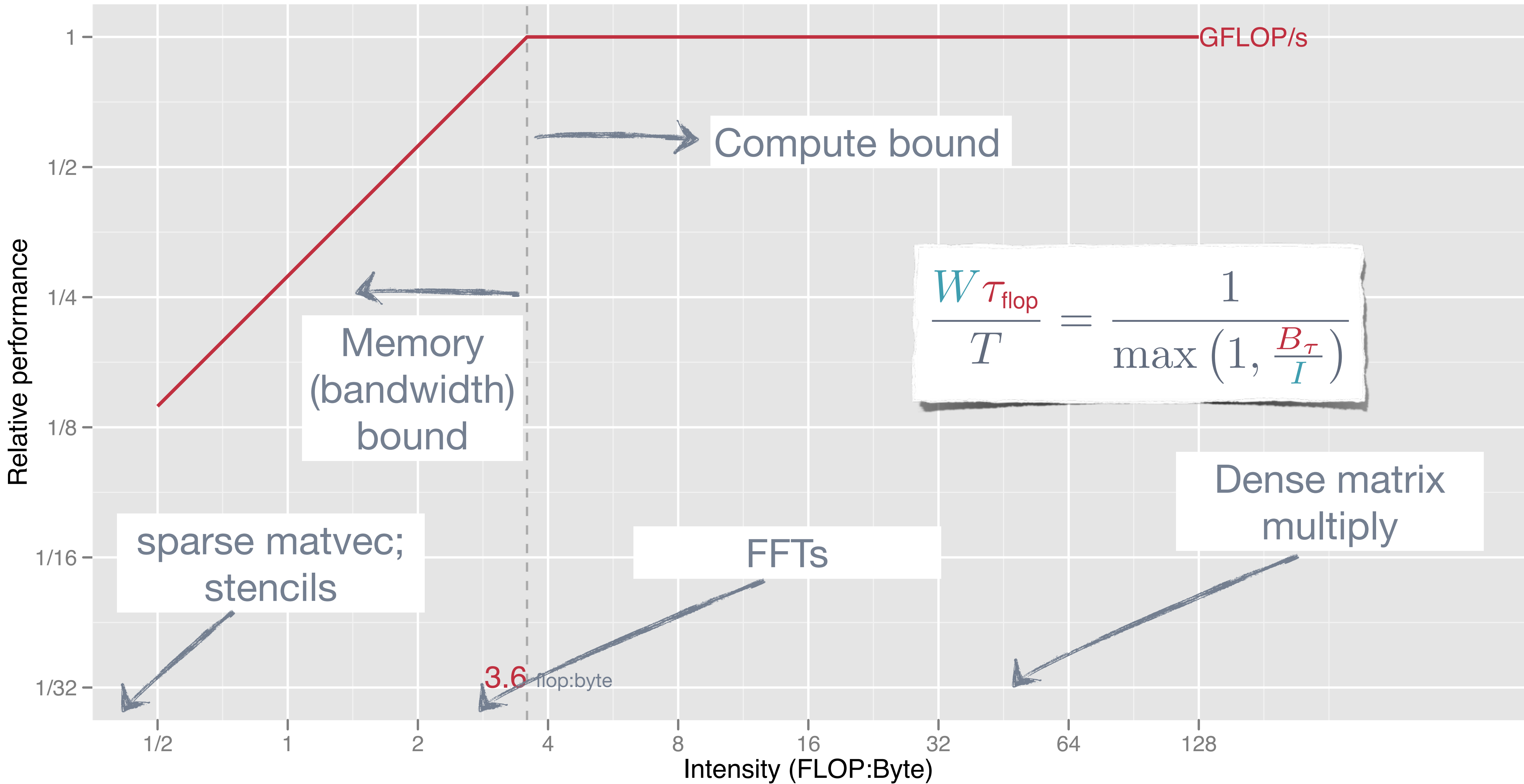
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

“Roofline” — Williams et al. (*Comm. ACM* '09)

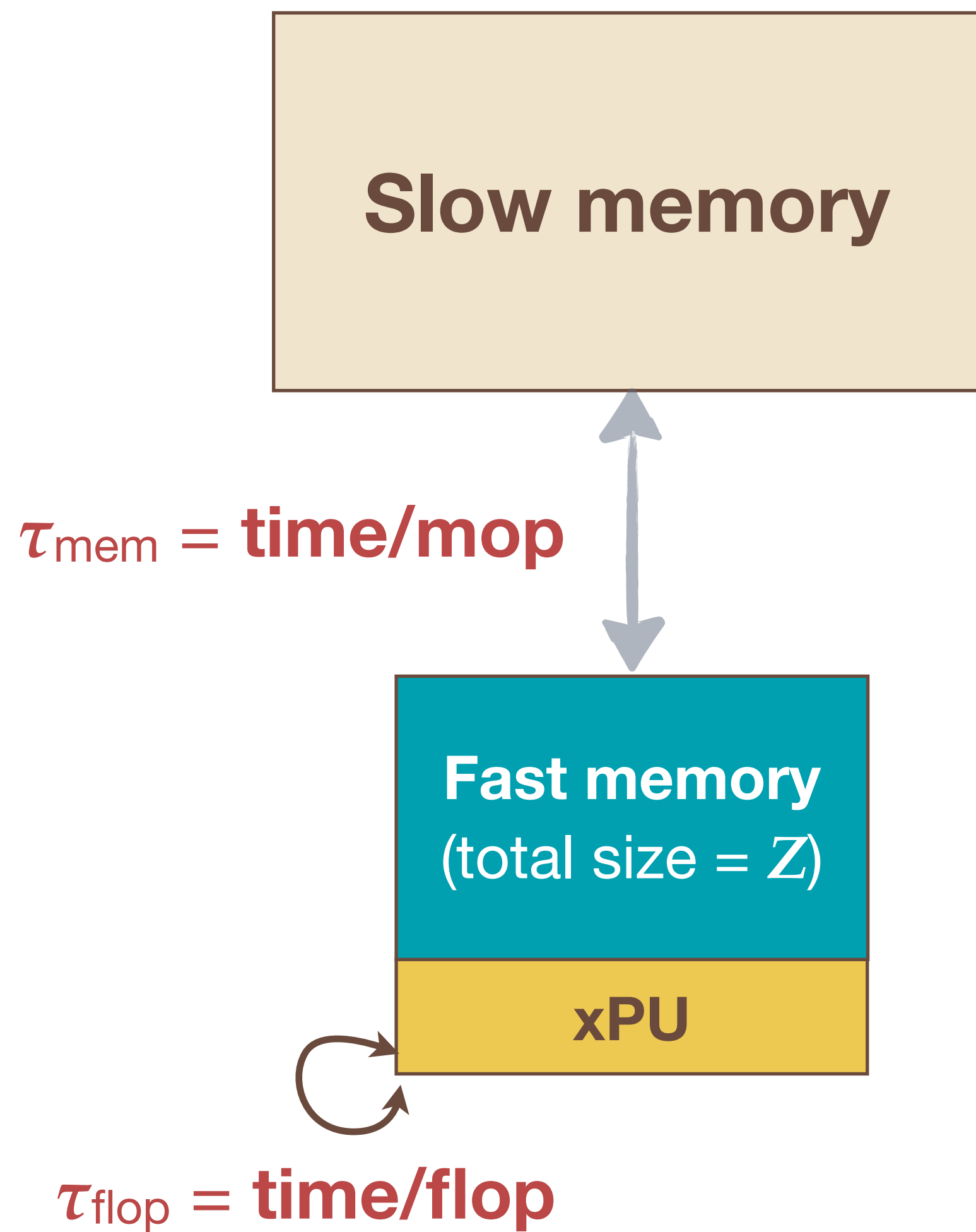


Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

“Roofline” — Williams et al. (*Comm. ACM* '09)

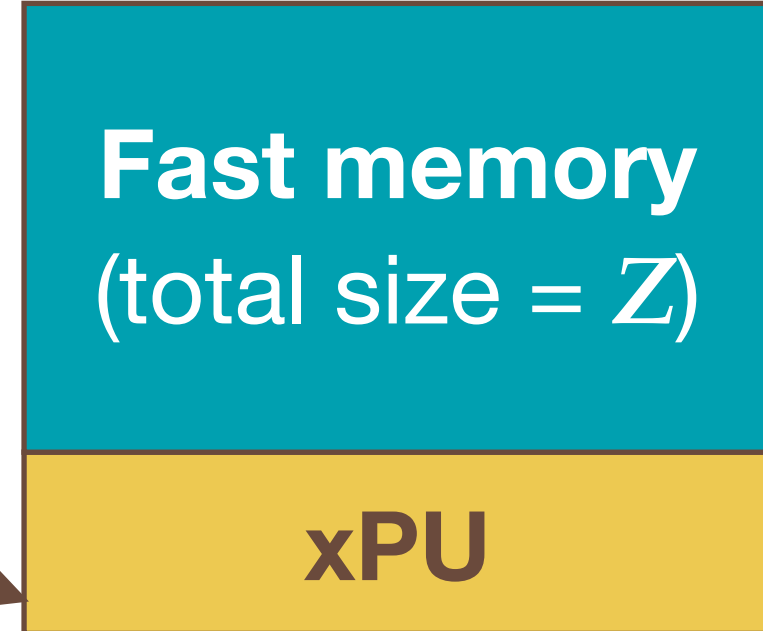
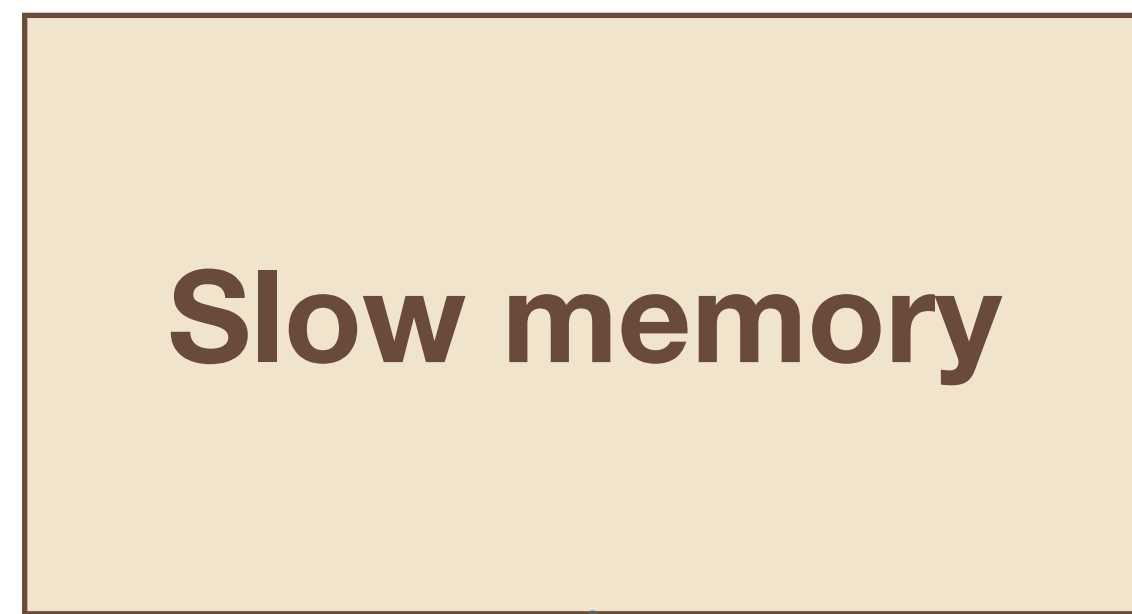


Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

An energy analogue

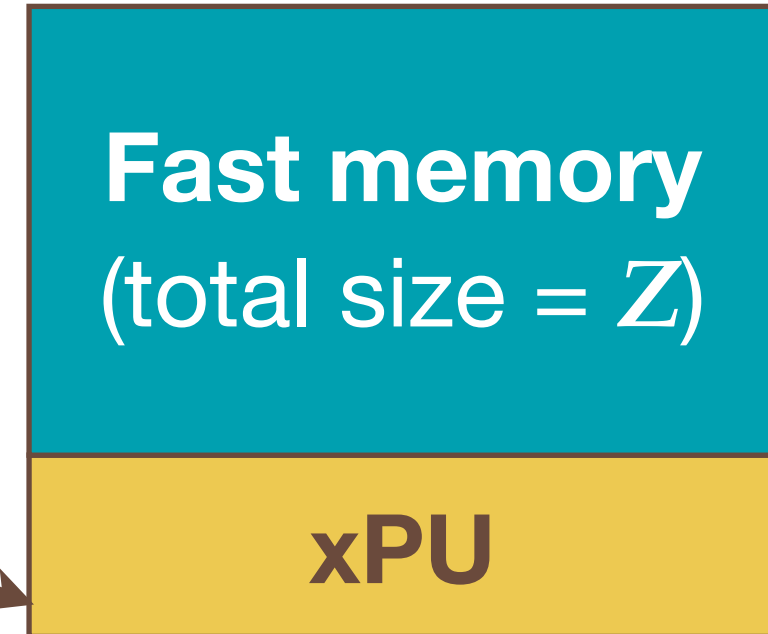
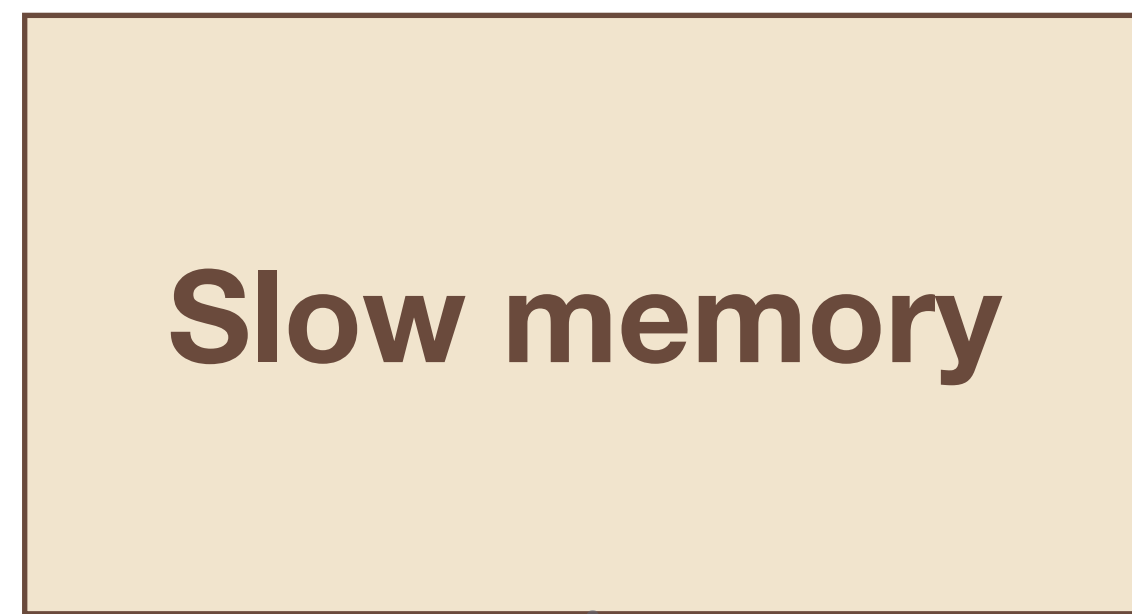


$\epsilon_{\text{mem}} = \text{energy/mop}$

$\epsilon_{\text{flop}} = \text{energy/flop}$

An energy analogue

$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right) \\
 E &= W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}} \\
 &= W \epsilon_{\text{flop}} \left(1 + \frac{B_{\epsilon}}{I}\right)
 \end{aligned}$$



$\epsilon_{\text{mem}} = \text{energy/mop}$

$\epsilon_{\text{flop}} = \text{energy/flop}$

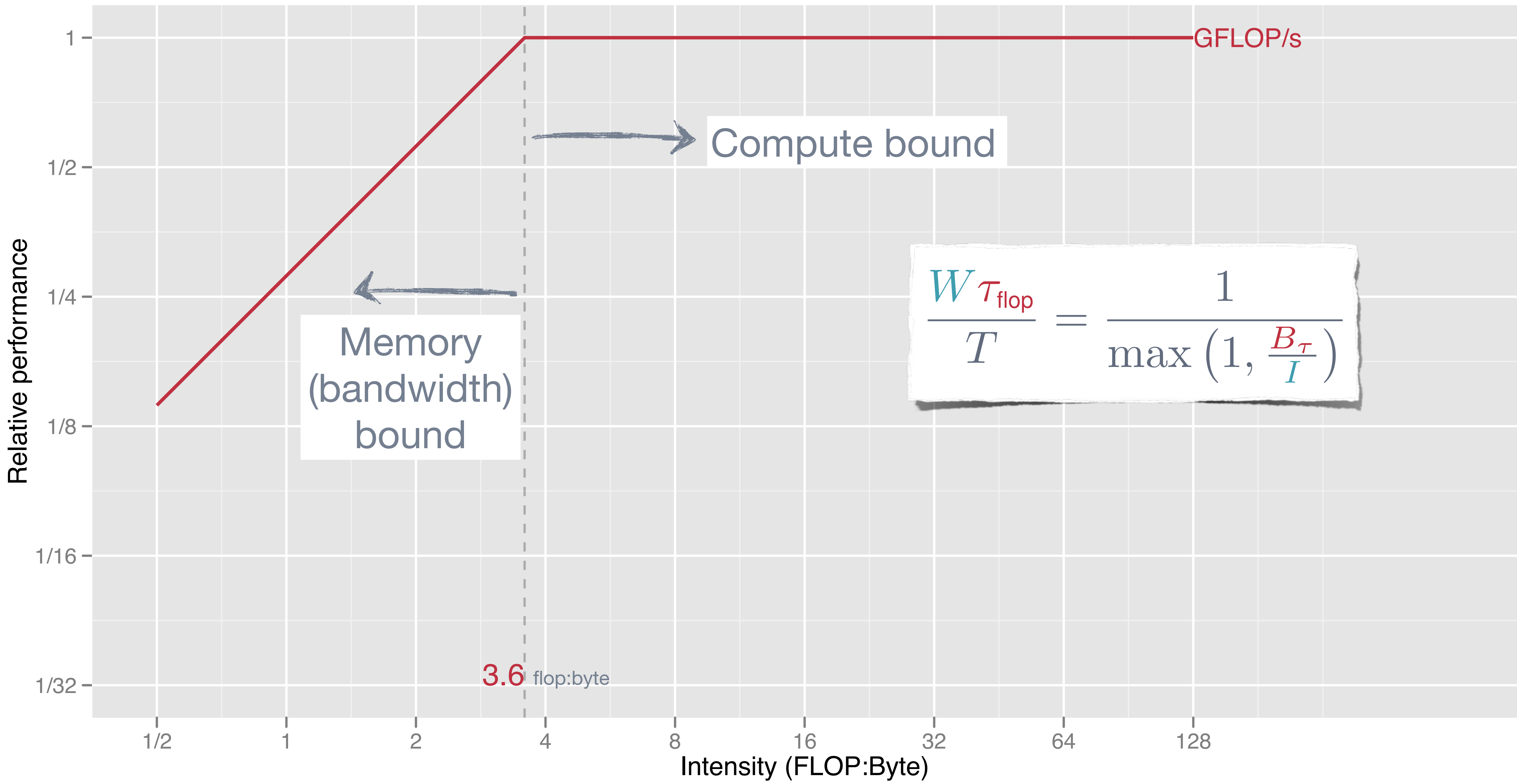
An energy analogue

$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

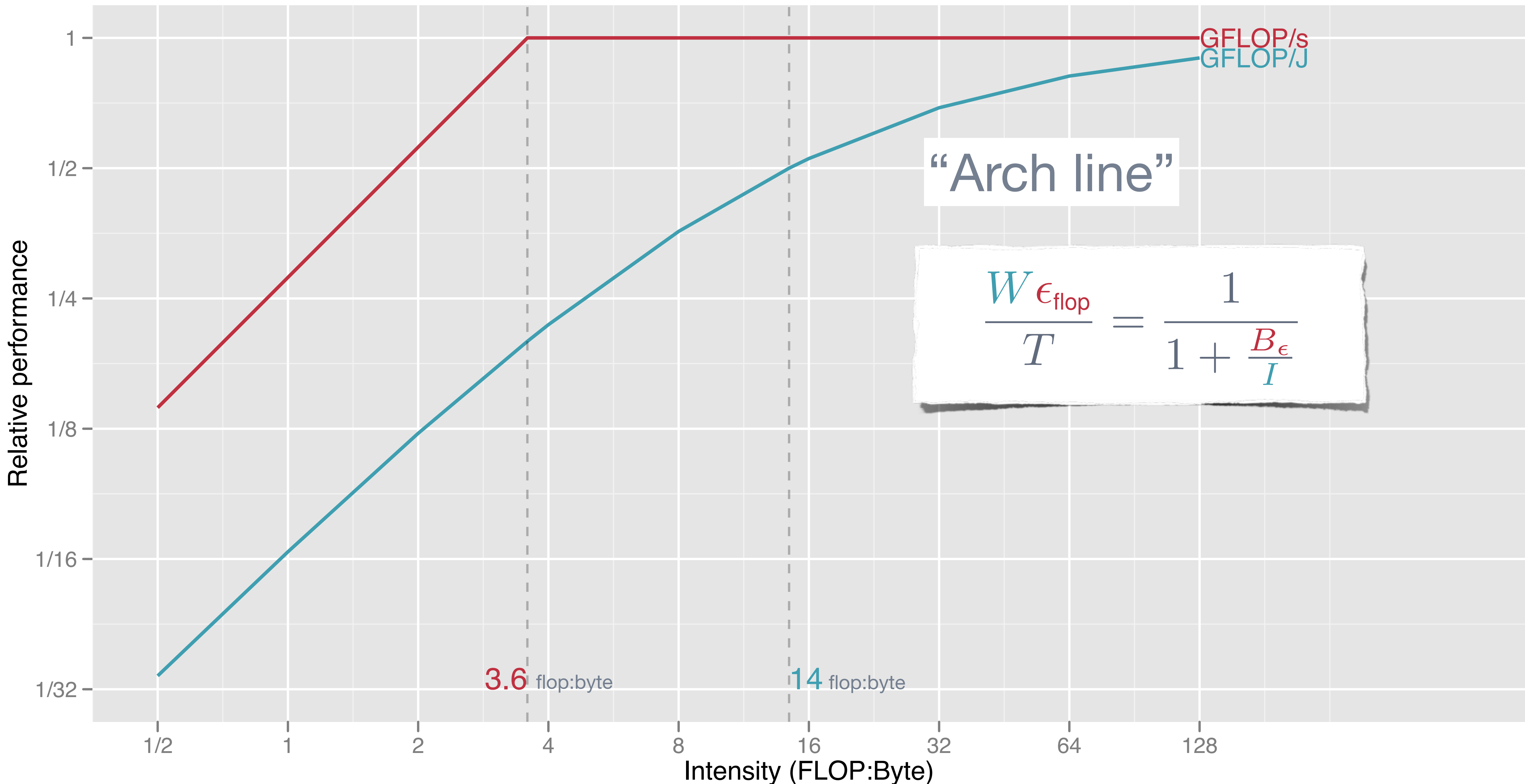
$$\begin{aligned}
 E &= W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}} \\
 &= W \epsilon_{\text{flop}} \left(1 + \frac{B_{\epsilon}}{I}\right)
 \end{aligned}$$

Consider:

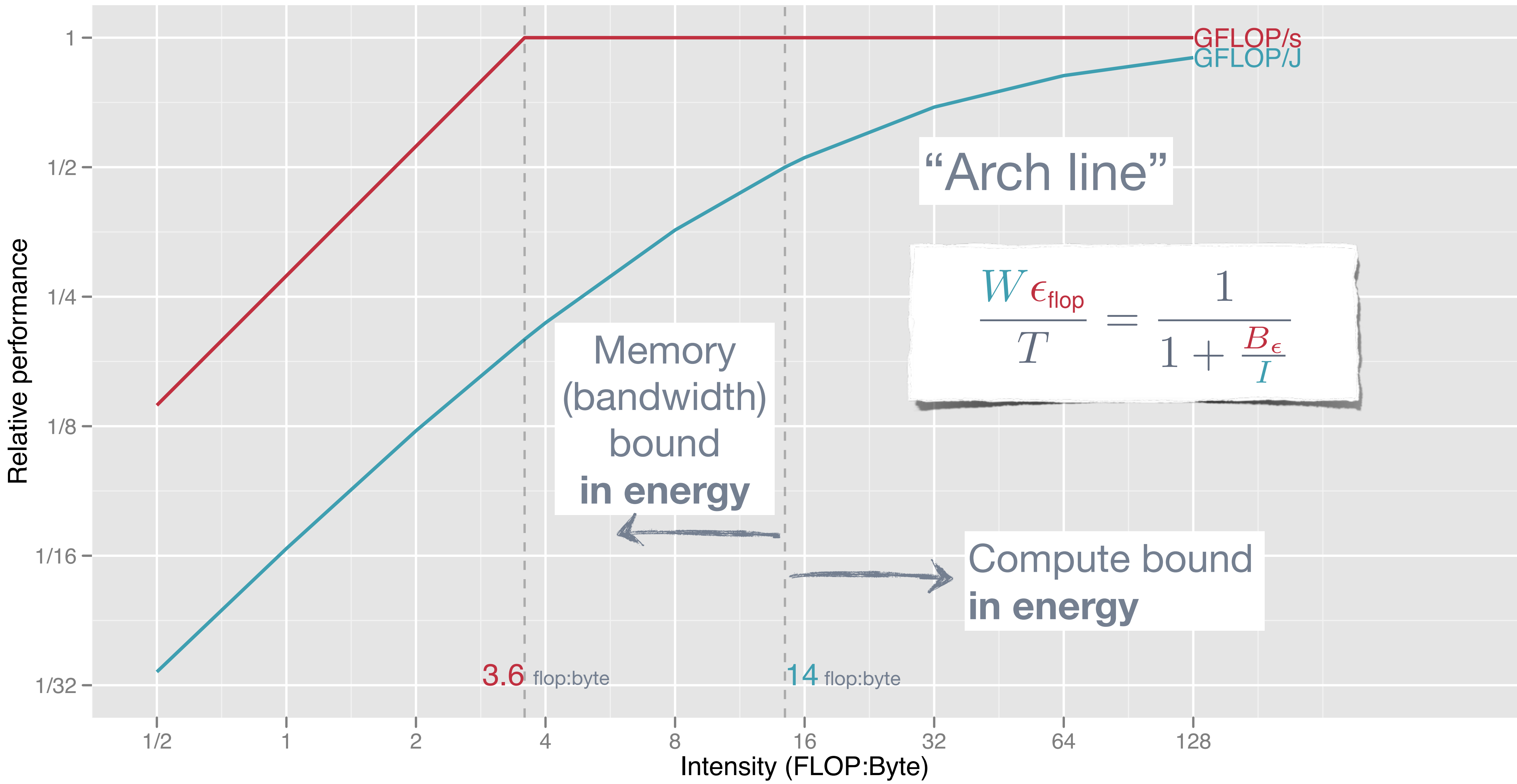
$$\frac{W \tau_{\text{flop}}}{T} \quad \text{and} \quad \frac{W \epsilon_{\text{flop}}}{E}$$



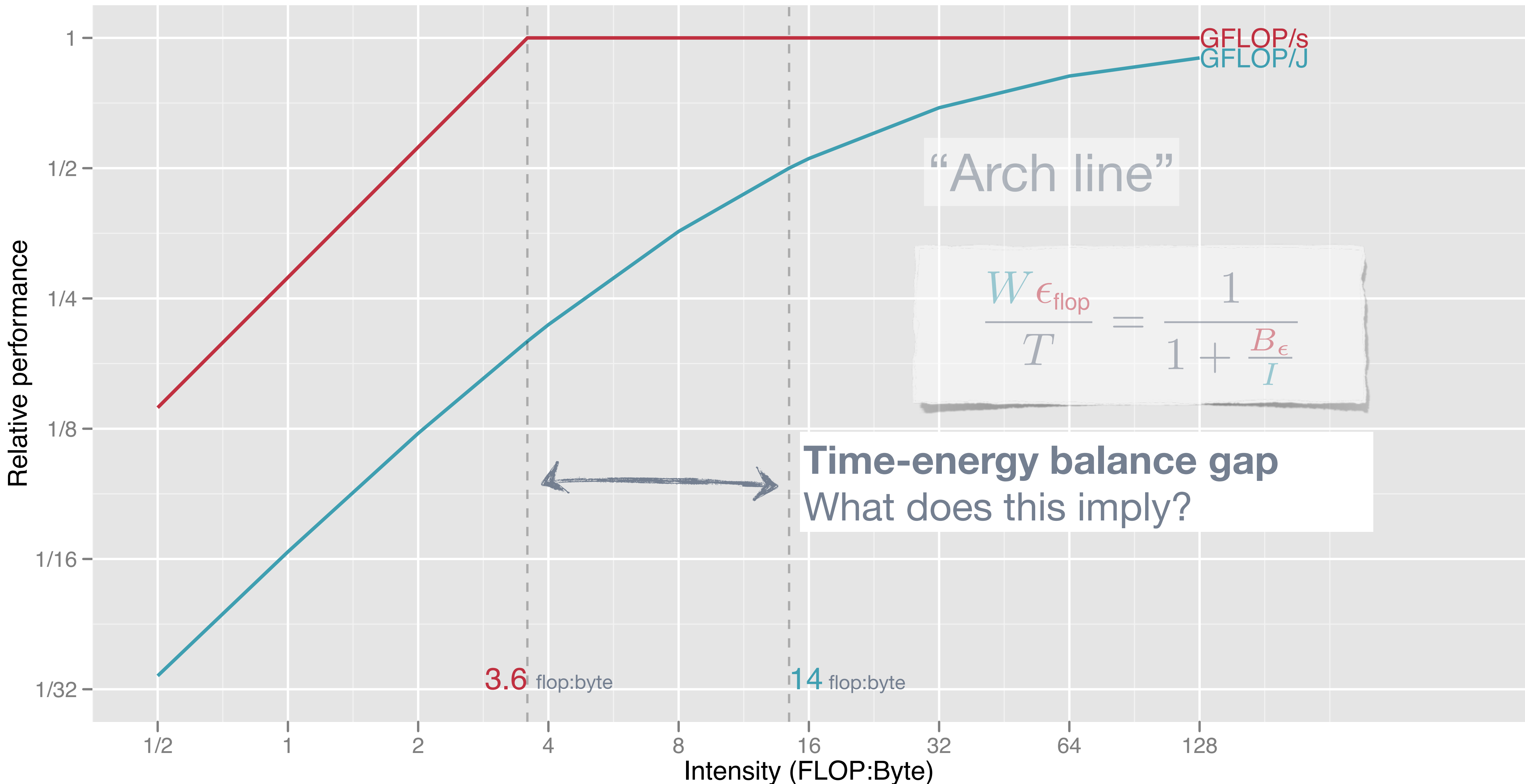
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



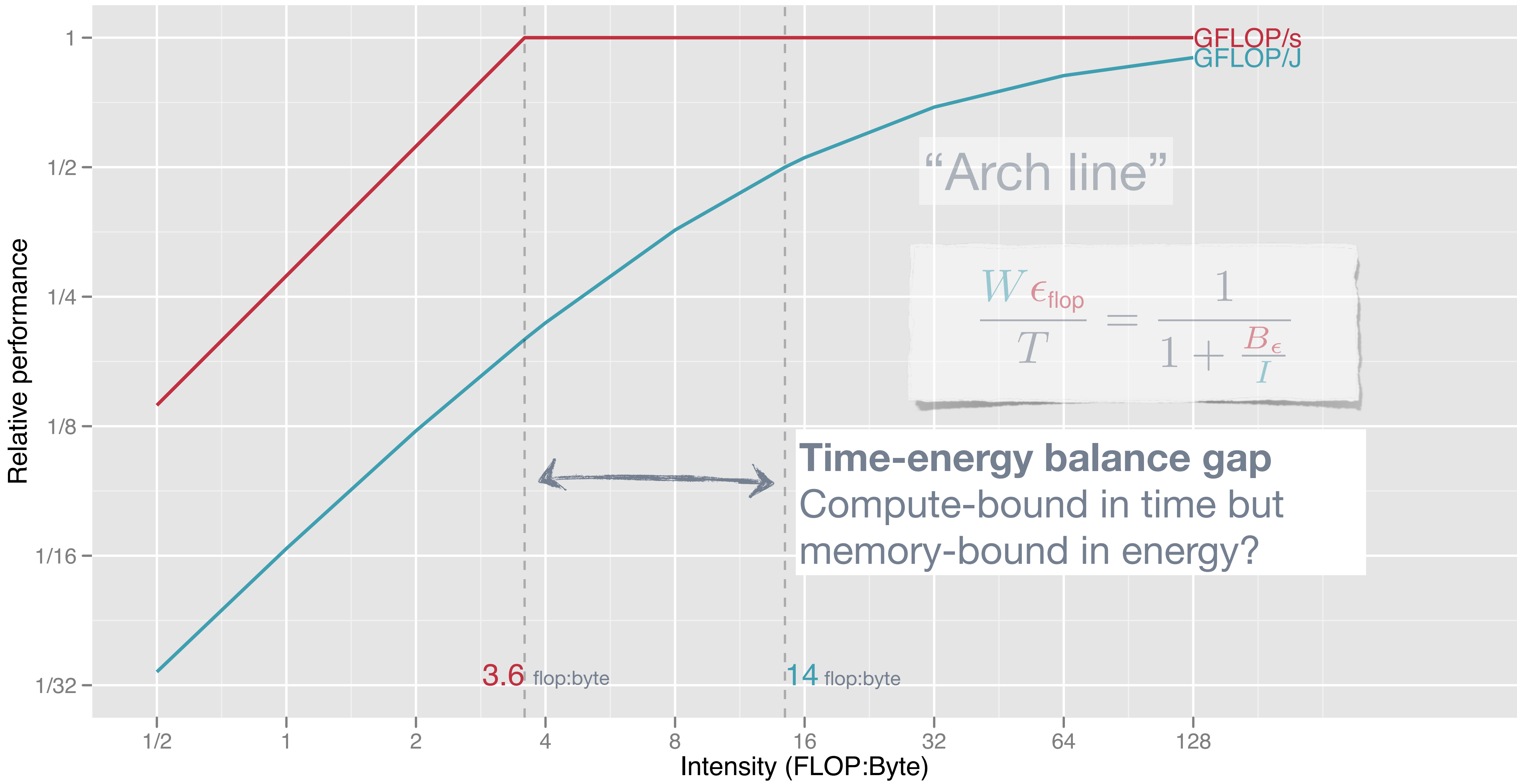
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



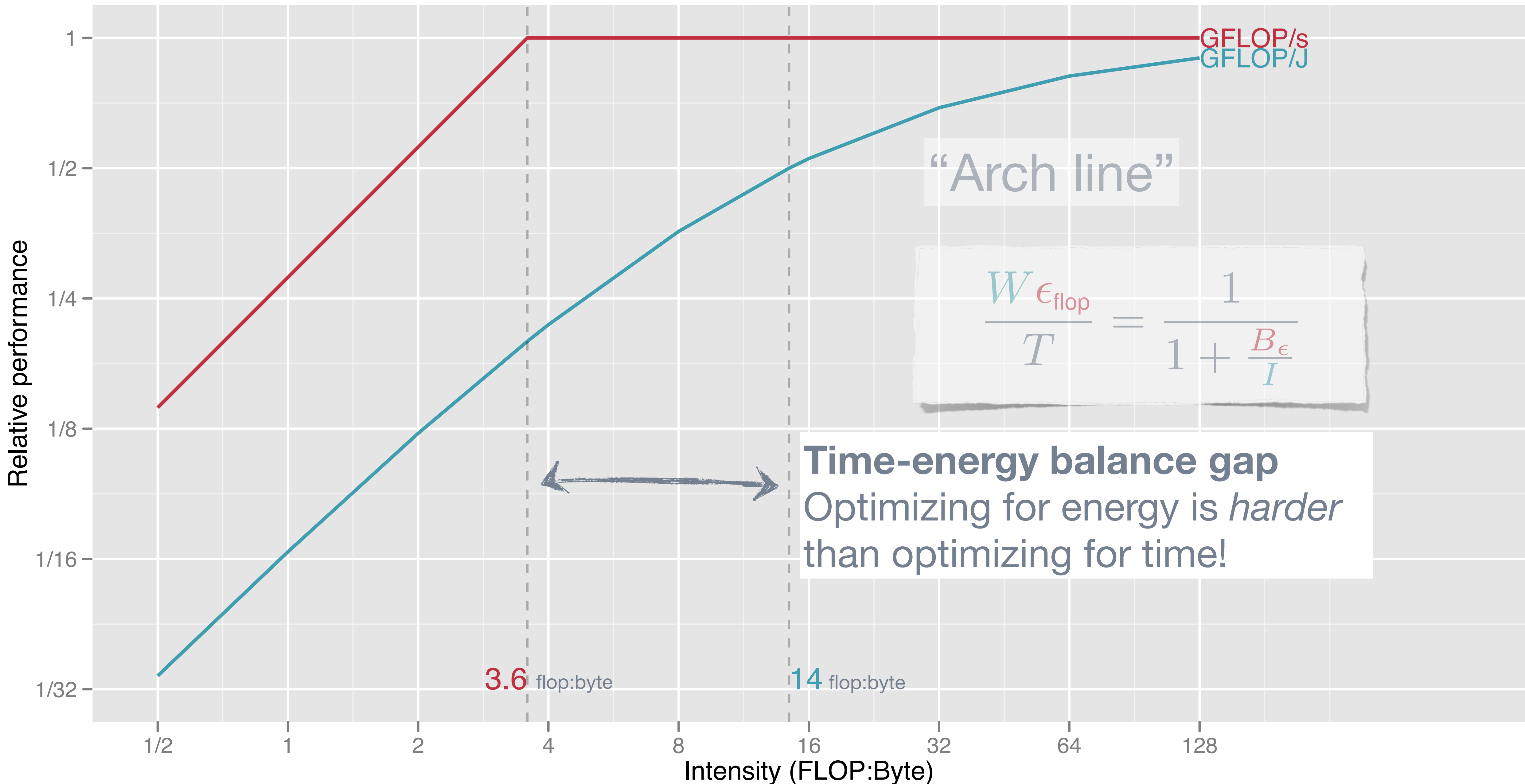
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



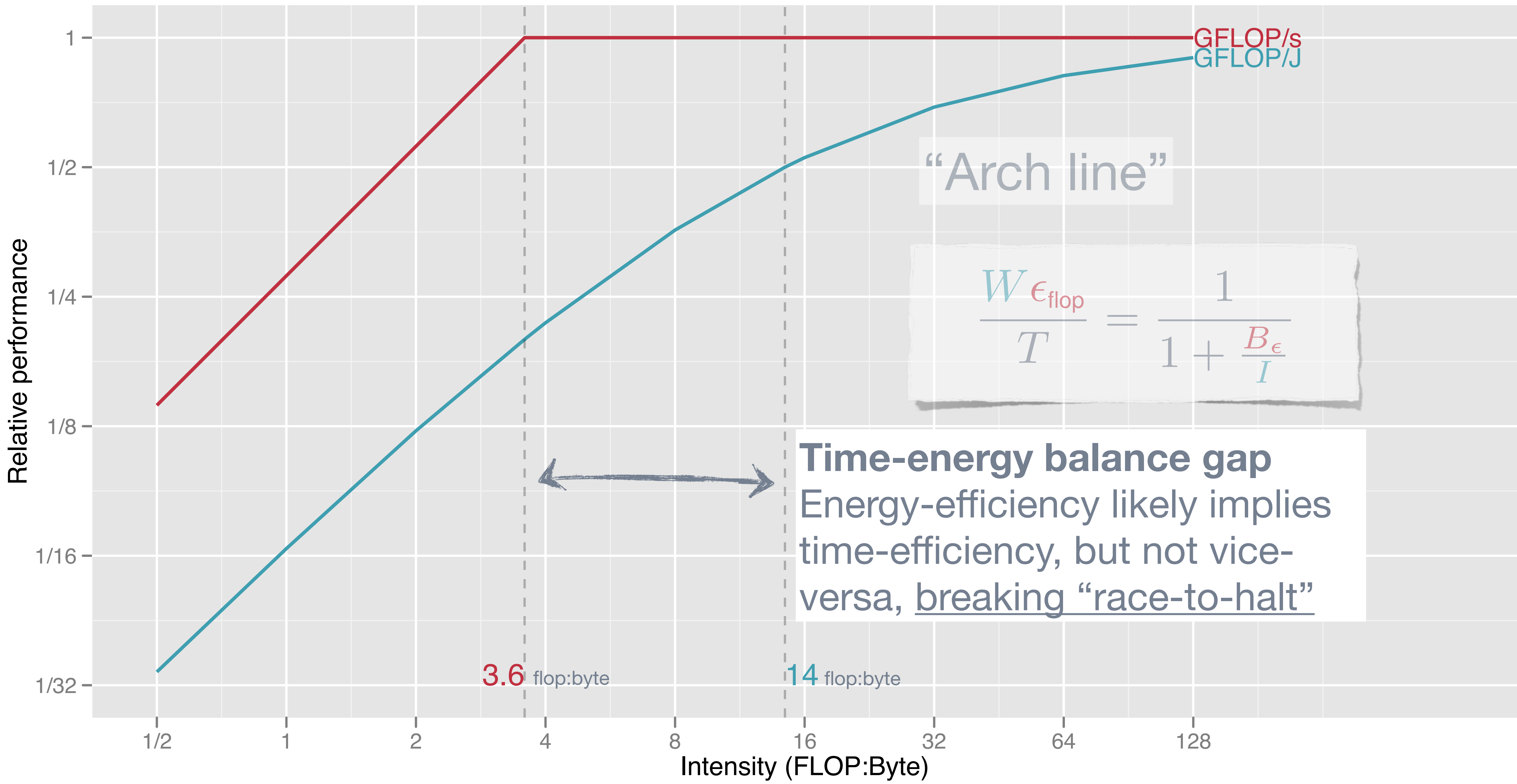
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



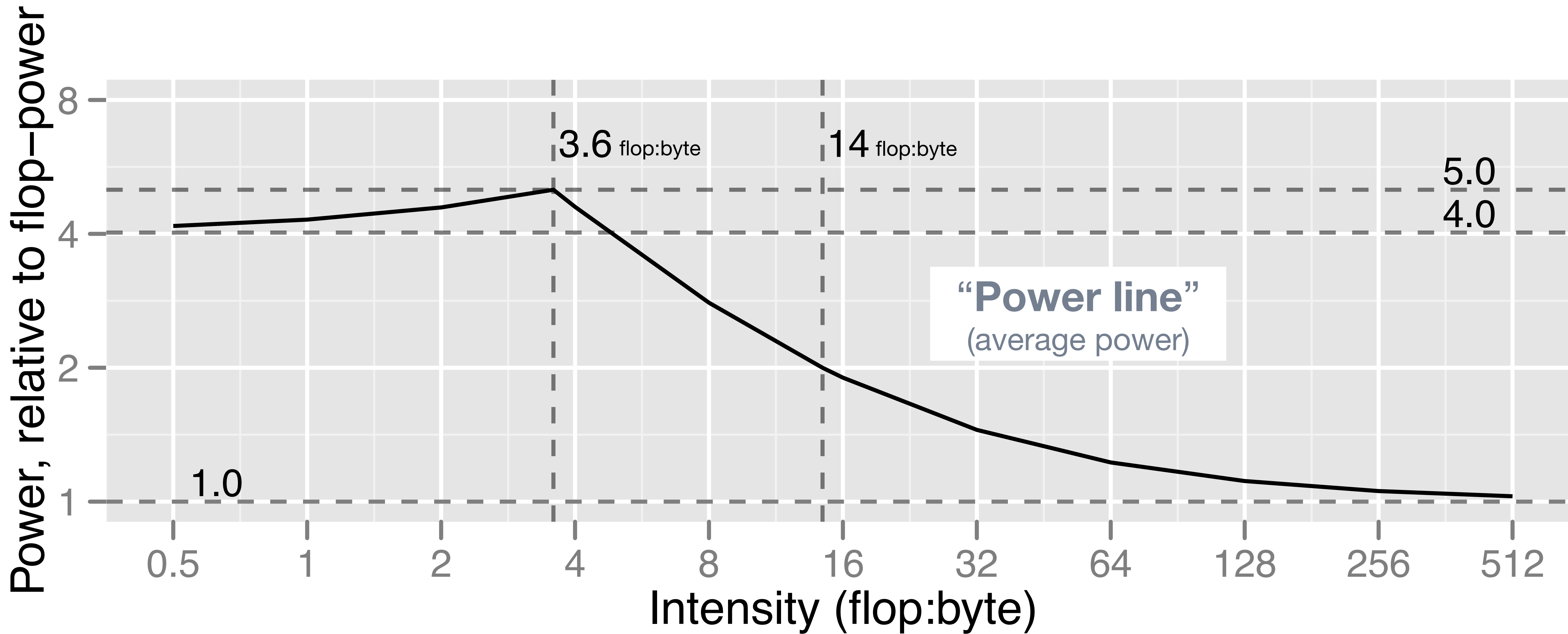
Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)



Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

Idea: Work-communication trade-offs

Algorithm 1 = (W, Q) versus Algorithm 2 = $(fW, \frac{Q}{m})$

$$I \equiv \frac{W}{Q}$$

Idea: Work-communication trade-offs

Algorithm 1 = (W, Q) versus Algorithm 2 = $(fW, \frac{Q}{m})$

$$I \equiv \frac{W}{Q}$$

$$\text{Speedup } \Delta T = \frac{T_{1,1}}{T_{f,m}}$$

$$\text{"Greenup" } \Delta E = \frac{E_{1,1}}{E_{f,m}}$$

Idea: Work-communication trade-offs

Algorithm 1 = (W, Q) versus Algorithm 2 = $(fW, \frac{Q}{m})$

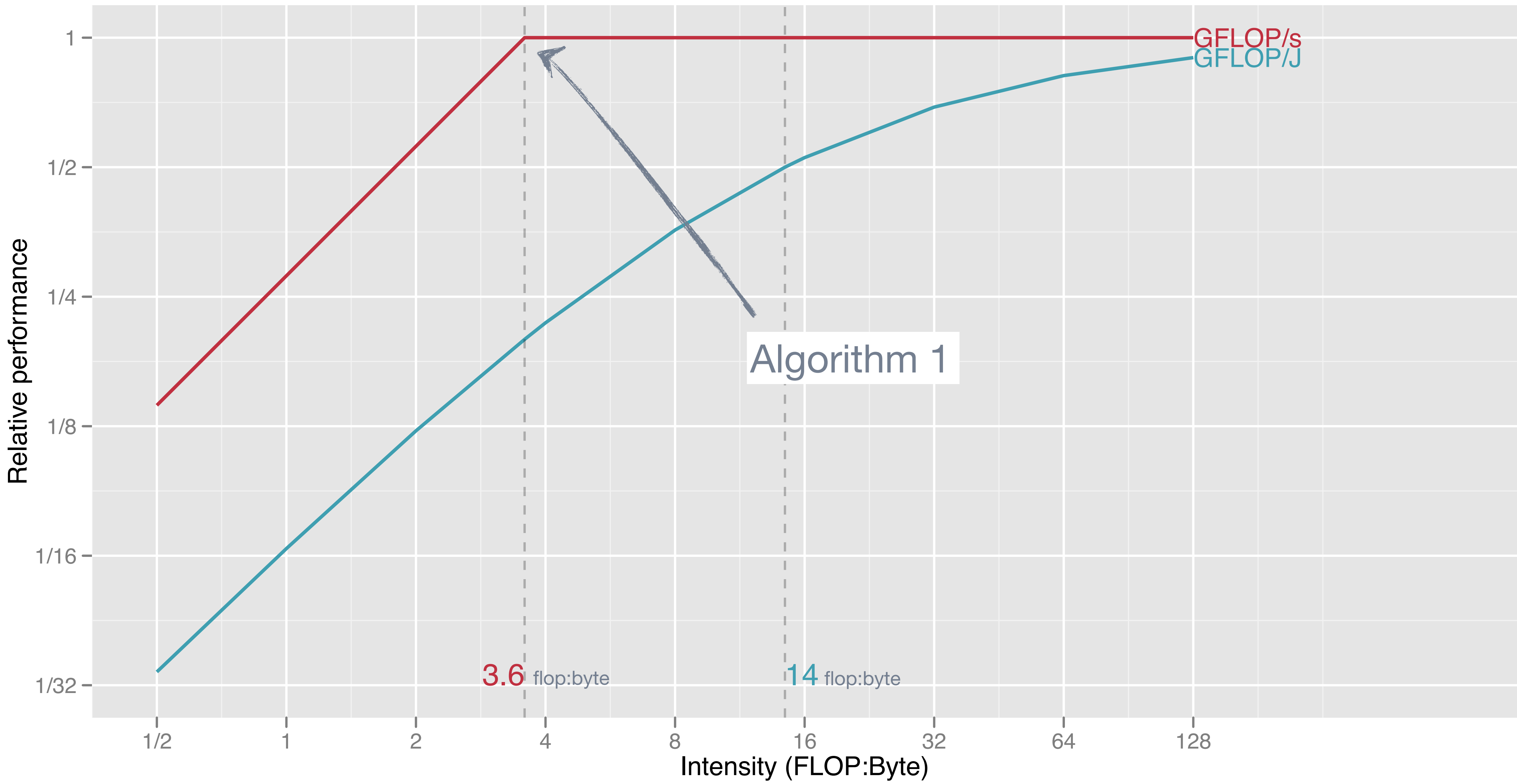
$$I \equiv \frac{W}{Q}$$

$$\text{Speedup } \Delta T = \frac{T_{1,1}}{T_{f,m}}$$

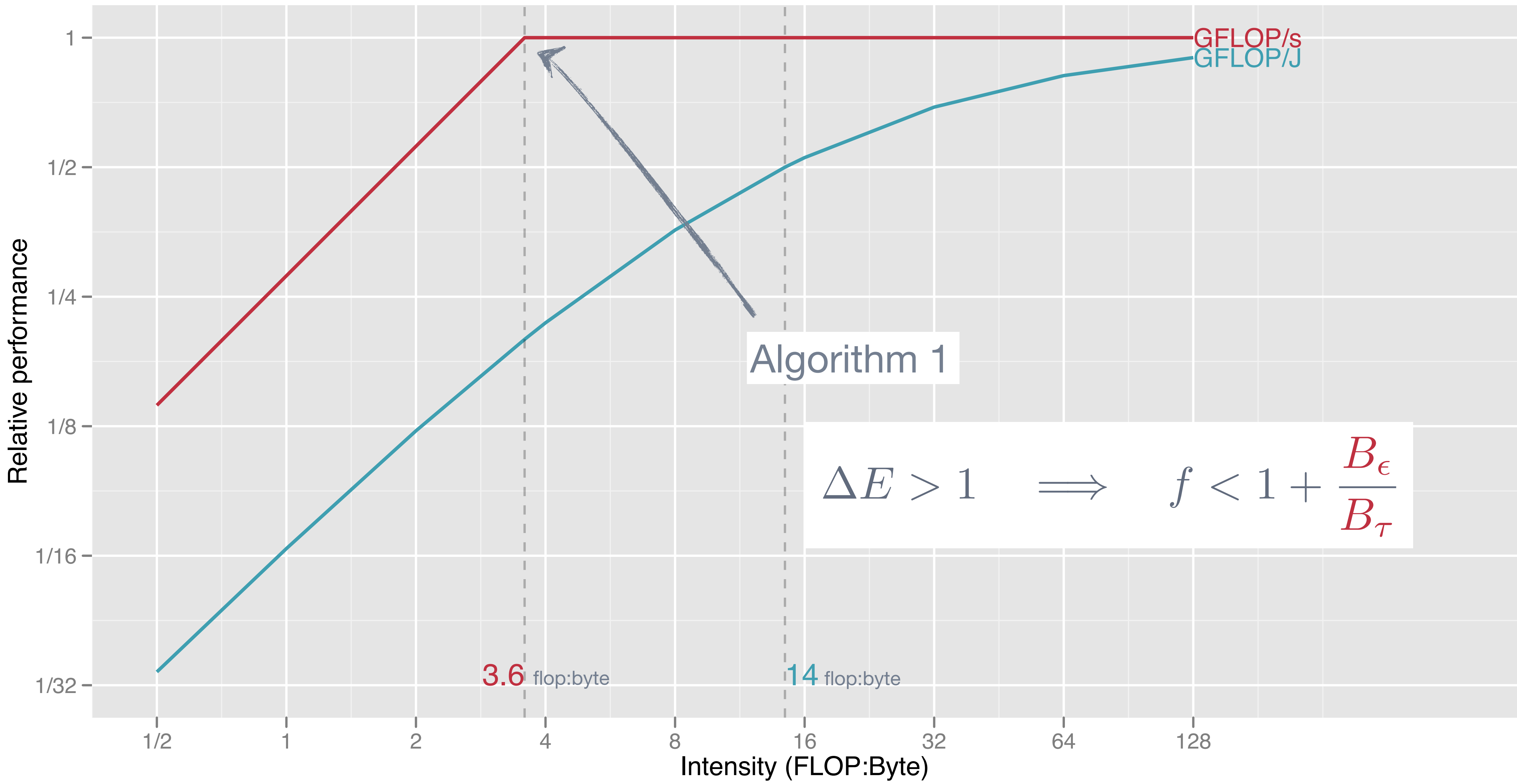
$$\text{"Greenup" } \Delta E = \frac{E_{1,1}}{E_{f,m}}$$

$$\Delta E > 1 \implies f < 1 + \frac{m-1}{m} \frac{B_\epsilon}{I}$$

A general "greenup" condition



Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

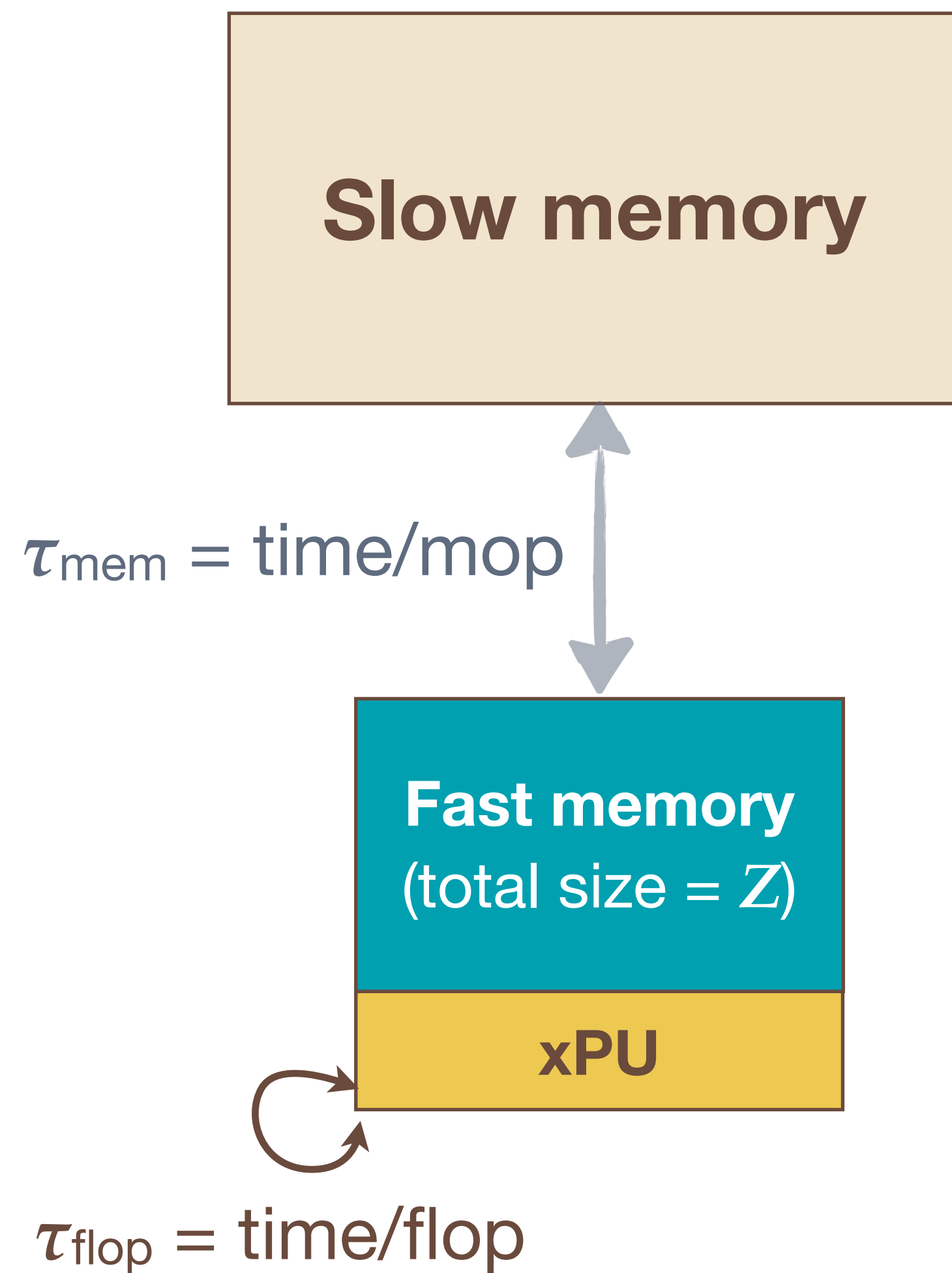


Balance estimates for a high-end NVIDIA Fermi in *double-precision*, according to Keckler et al. *IEEE Micro* (2011)

That was theory. What happens in practice?

⇒ Cannot ignore **constant power**.

Let's add it to our model and measure.

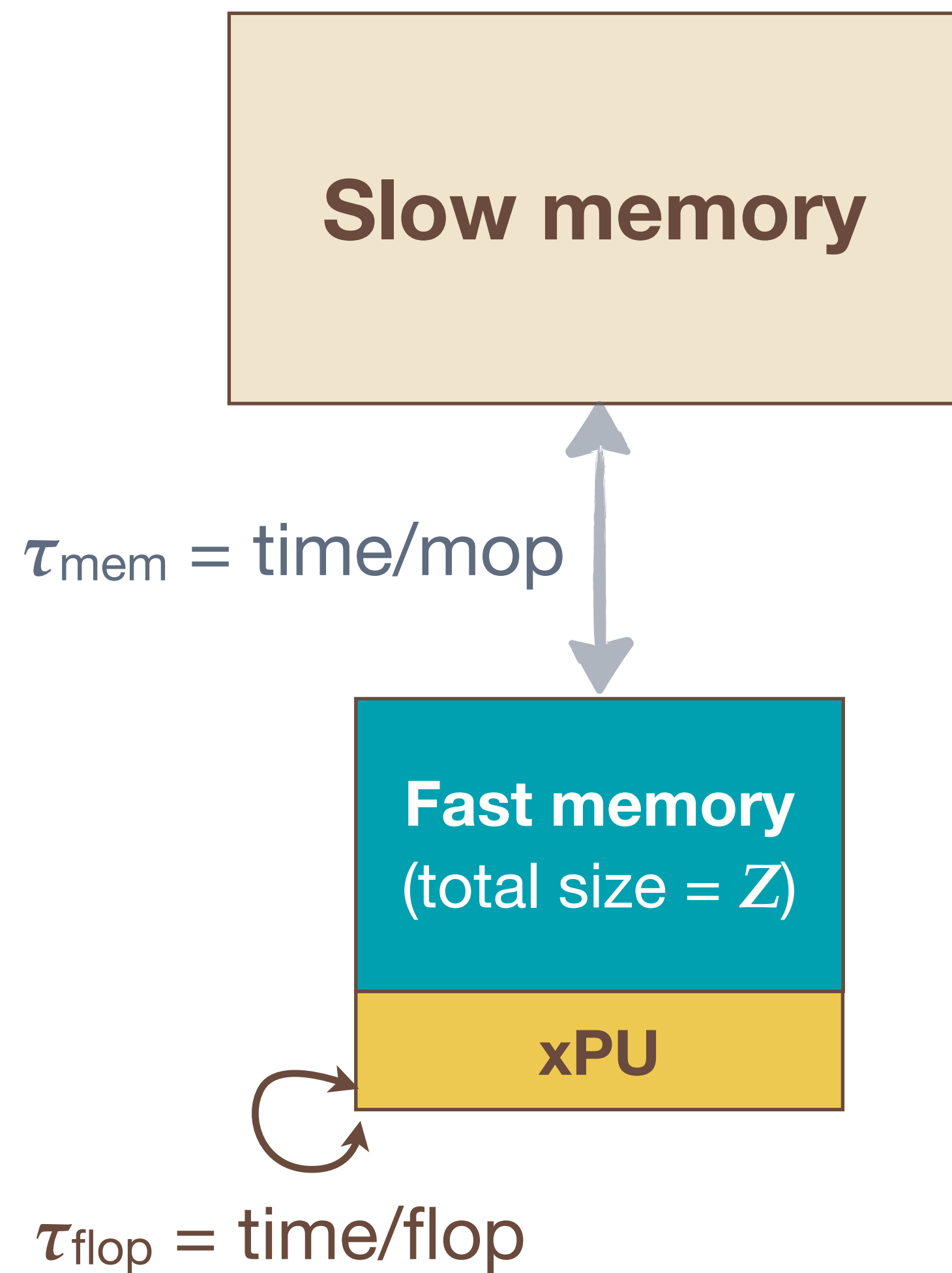


Constant power

$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right)
 \end{aligned}$$

$$\begin{aligned}
 E &= W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}} + T \pi_0 \\
 &= W \epsilon_{\text{flop}} \left(1 + \frac{B_{\epsilon}}{I} + \frac{\pi_0}{\epsilon_{\text{flop}}} \frac{T}{W}\right)
 \end{aligned}$$

Consider: $\frac{W \tau_{\text{flop}}}{T}$ and $\frac{W \epsilon_{\text{flop}}}{E}$



Constant power

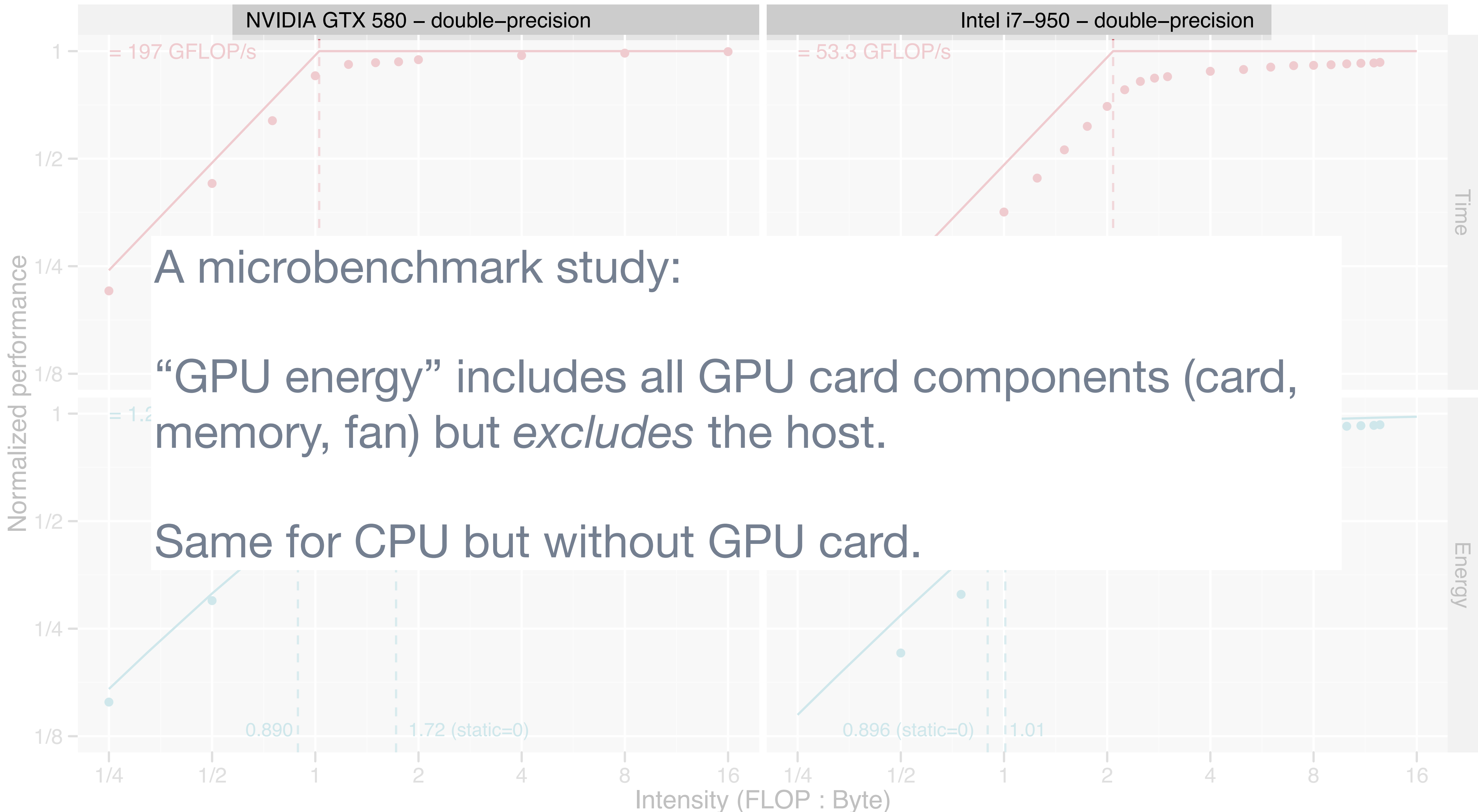
$$\begin{aligned}
 T &= \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}}) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{Q}{W} \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}}\right) \\
 &= W \tau_{\text{flop}} \max\left(1, \frac{B_{\tau}}{I}\right) \text{ Constant power}
 \end{aligned}$$

$$\begin{aligned}
 E &= W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}} + T \pi_0 \\
 &= W \epsilon_{\text{flop}} \left(1 + \frac{B_{\epsilon}}{I} + \frac{\pi_0}{\epsilon_{\text{flop}}} \frac{T}{W}\right)
 \end{aligned}$$

Consider: $\frac{W \tau_{\text{flop}}}{T}$ and $\frac{W \epsilon_{\text{flop}}}{E}$

NVIDIA GTX 580 – double-precision

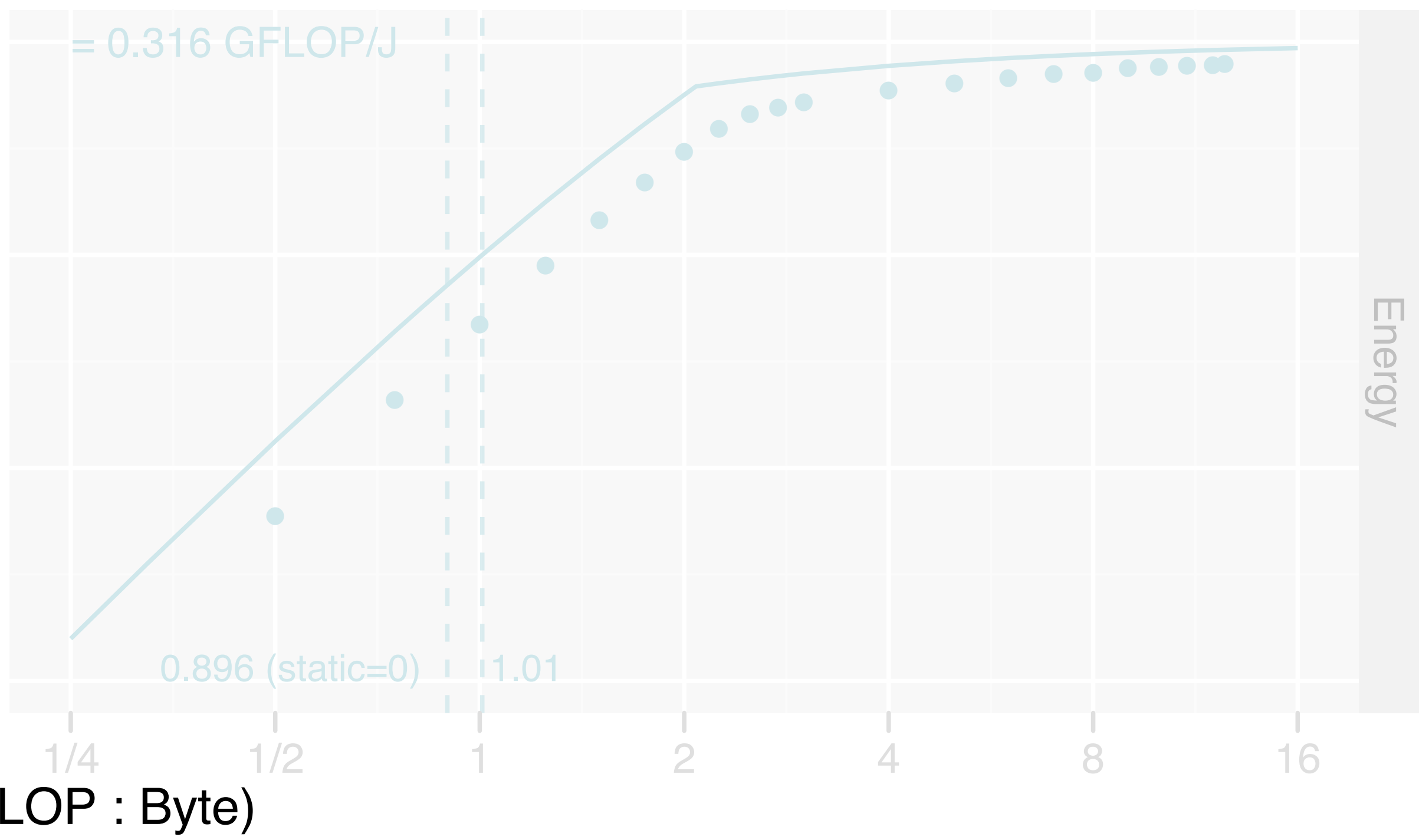
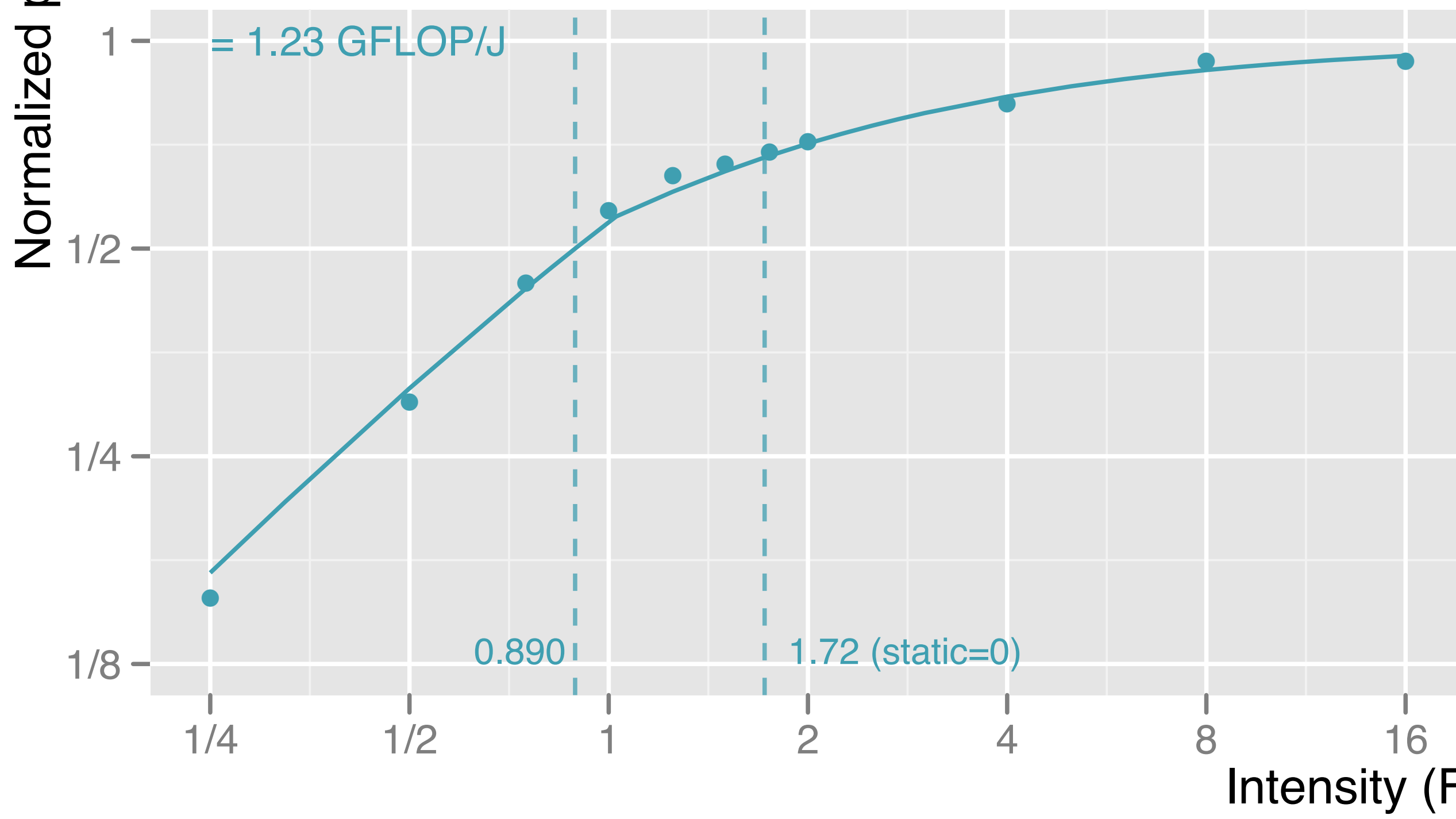
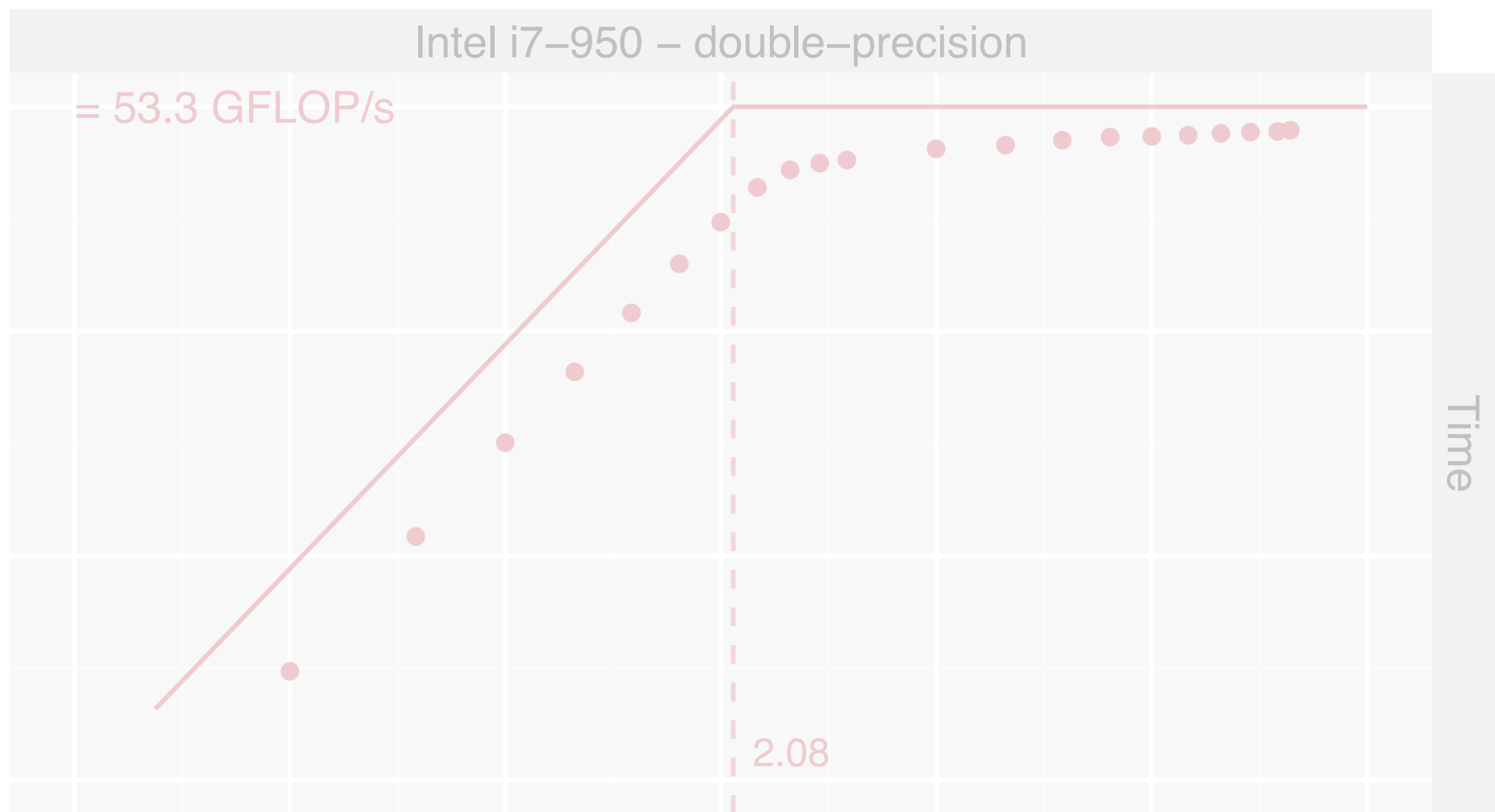
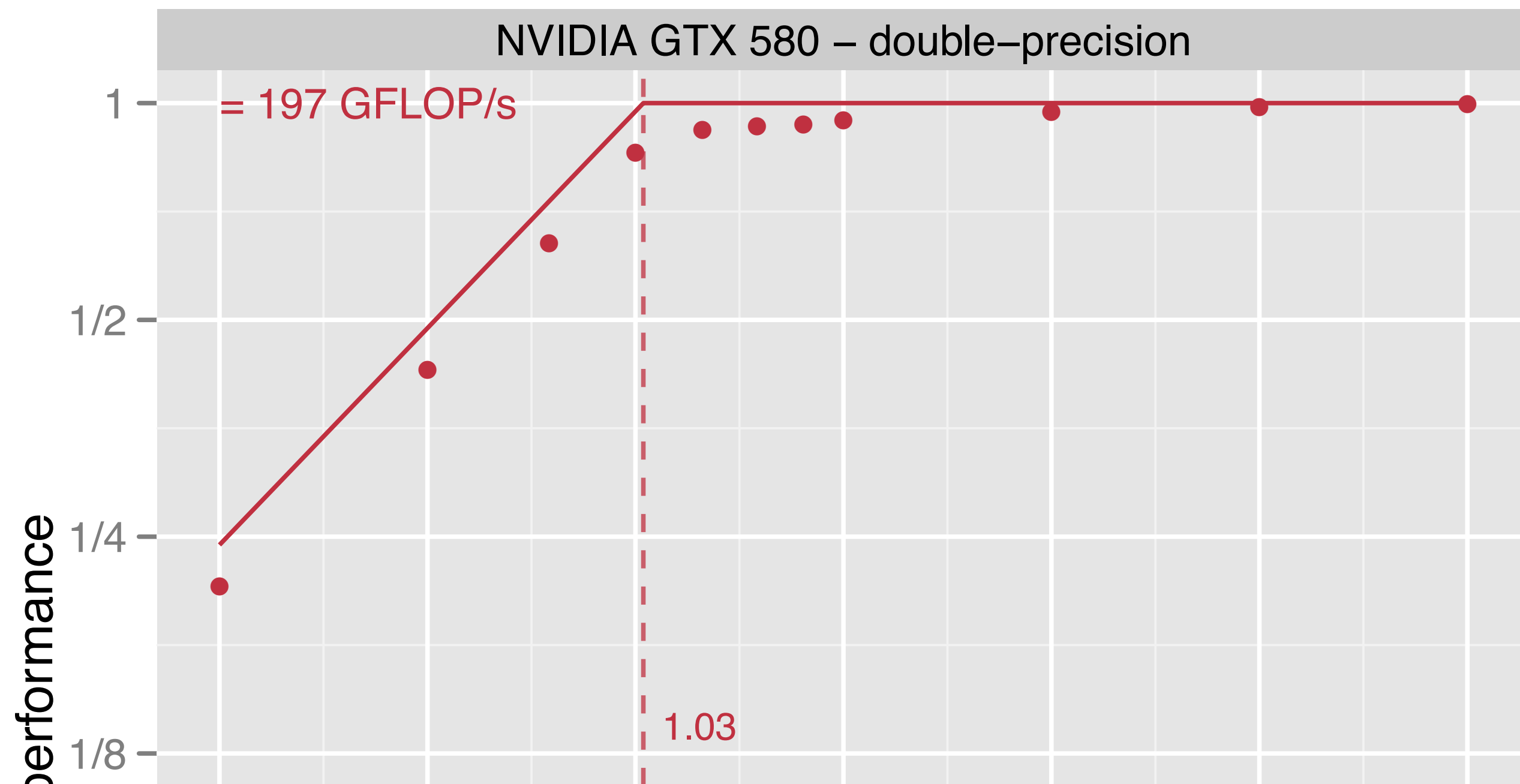
Intel i7-950 – double-precision



A microbenchmark study:

“GPU energy” includes all GPU card components (card, memory, fan) but *excludes* the host.

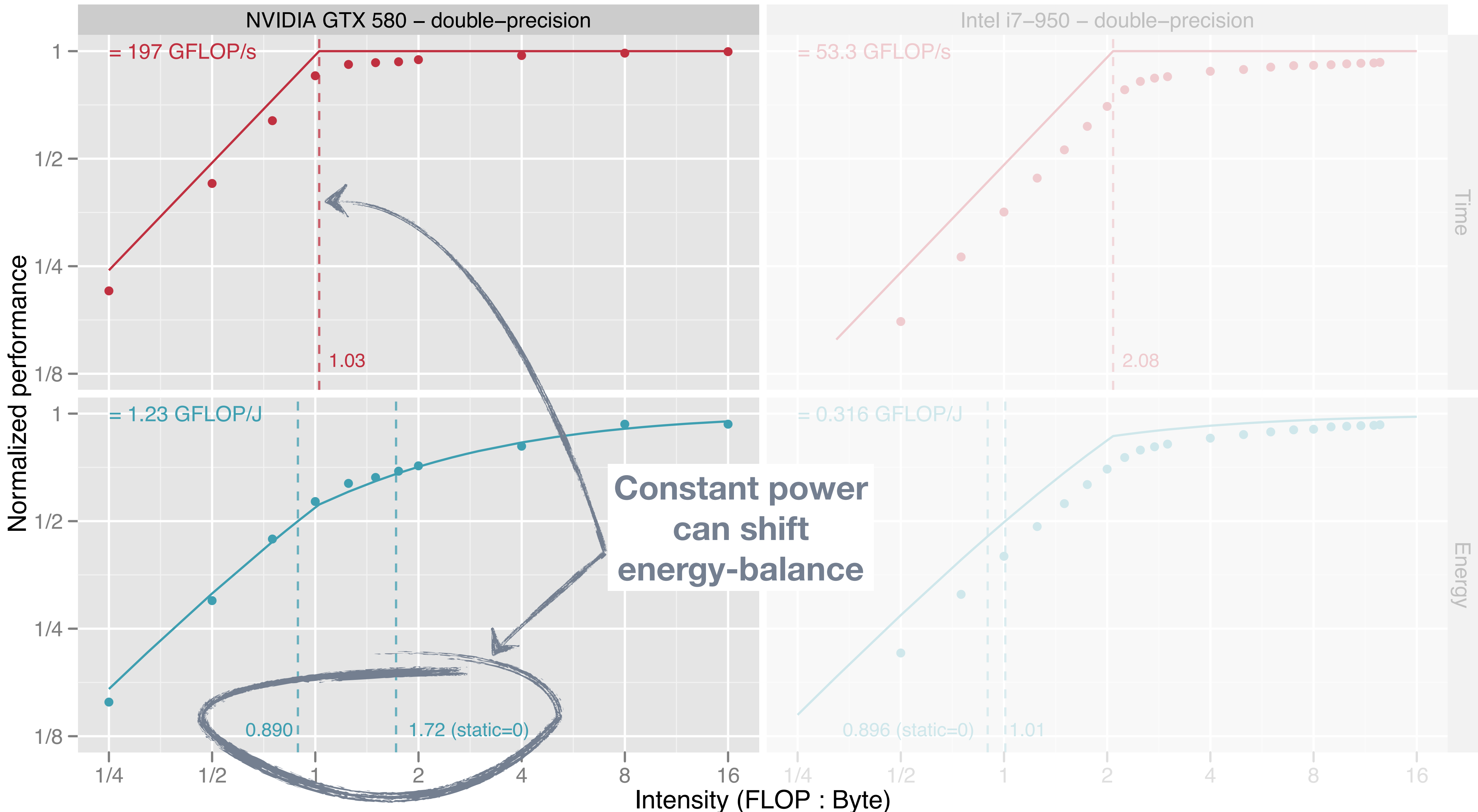
Same for CPU but without GPU card.

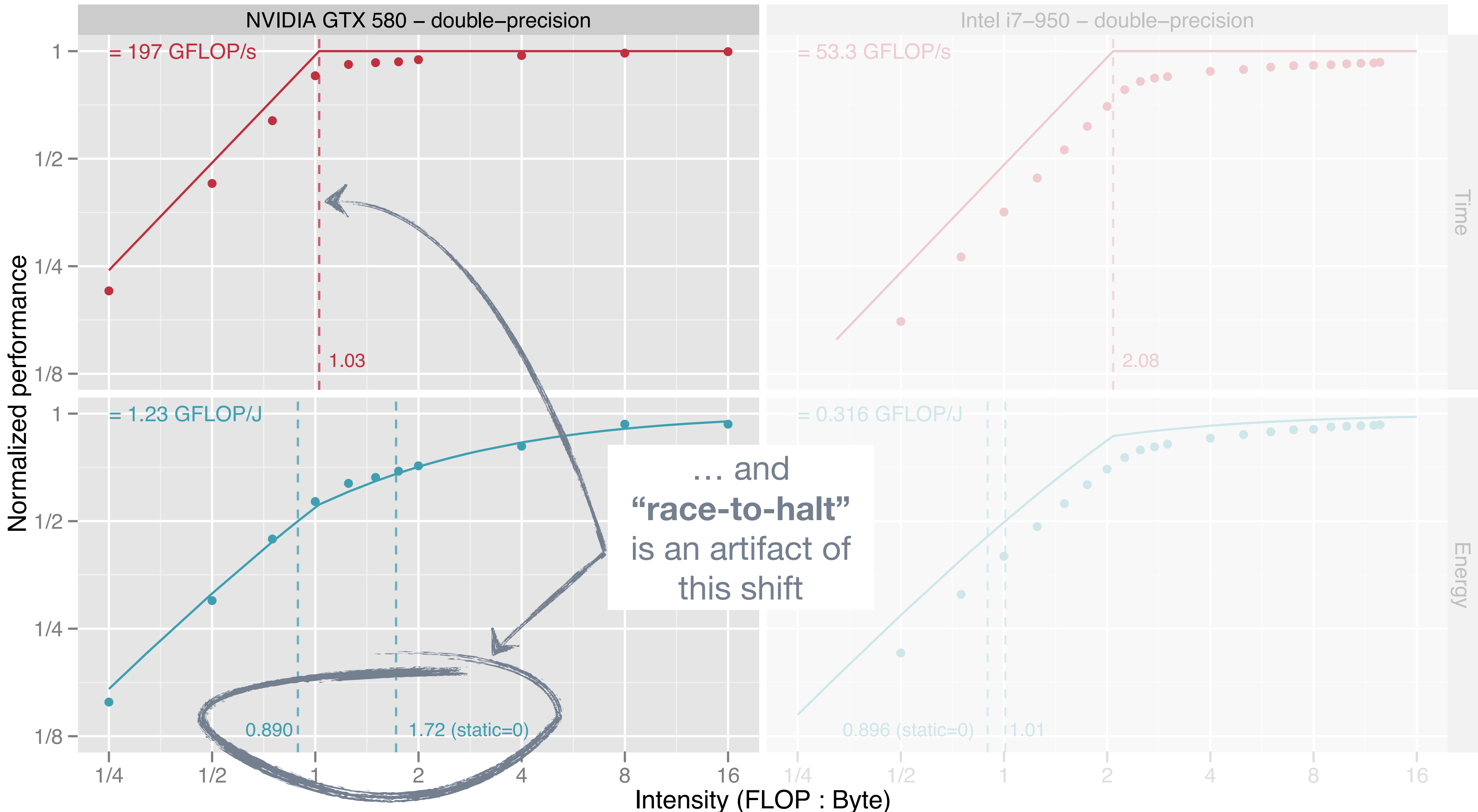


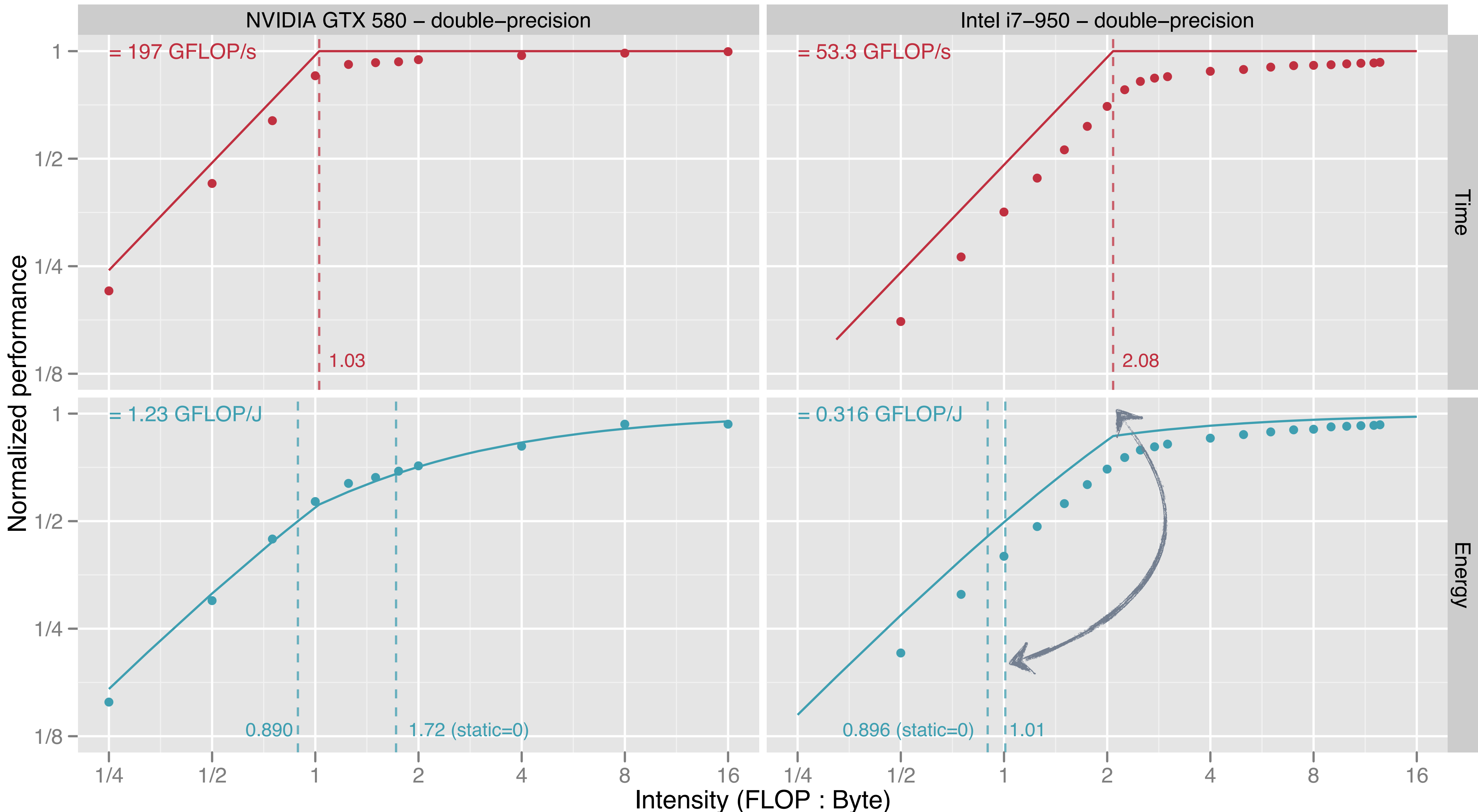
Time

Energy

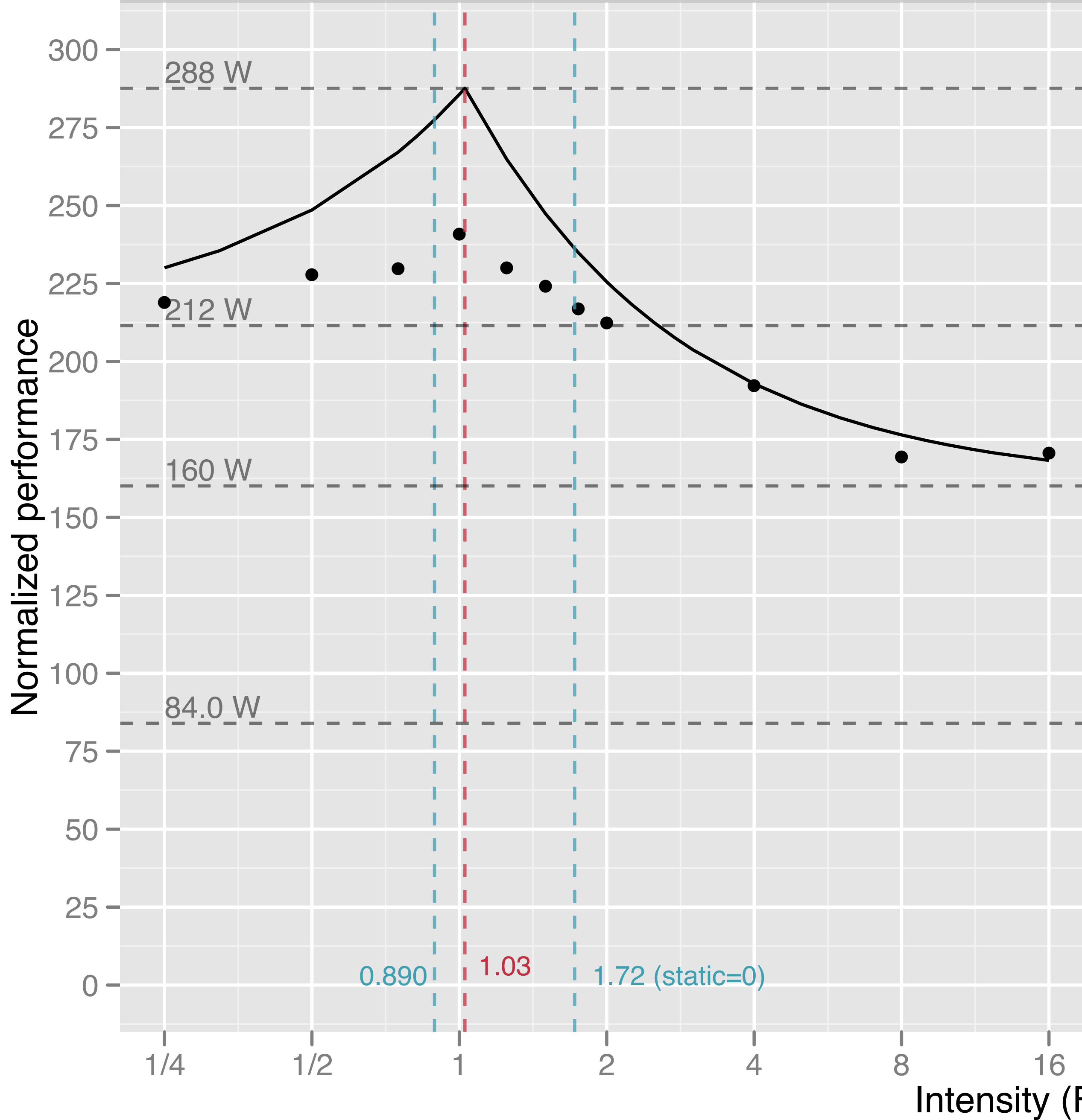
Intensity (FLOP : Byte)



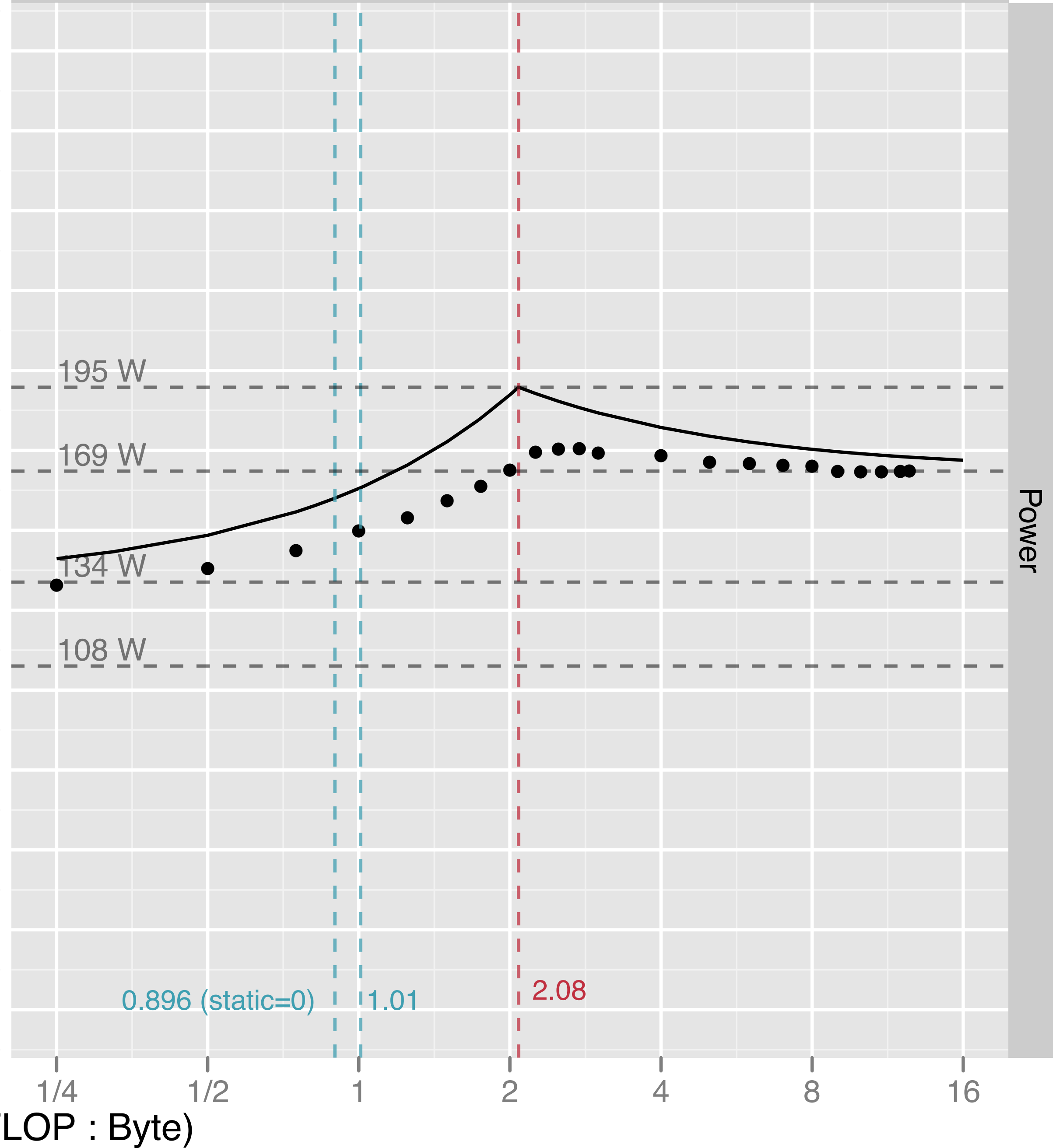




NVIDIA GTX 580 – double-precision



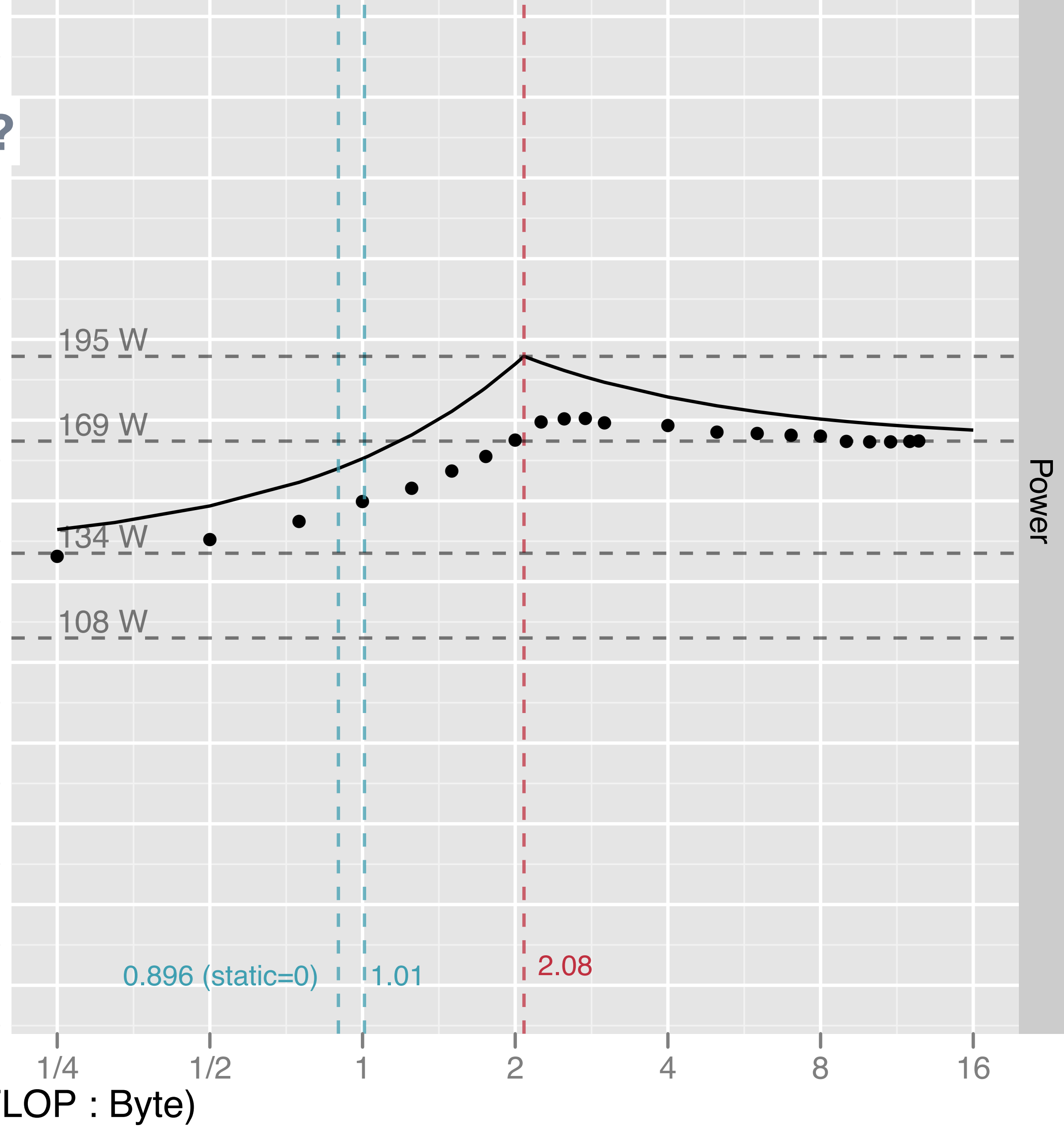
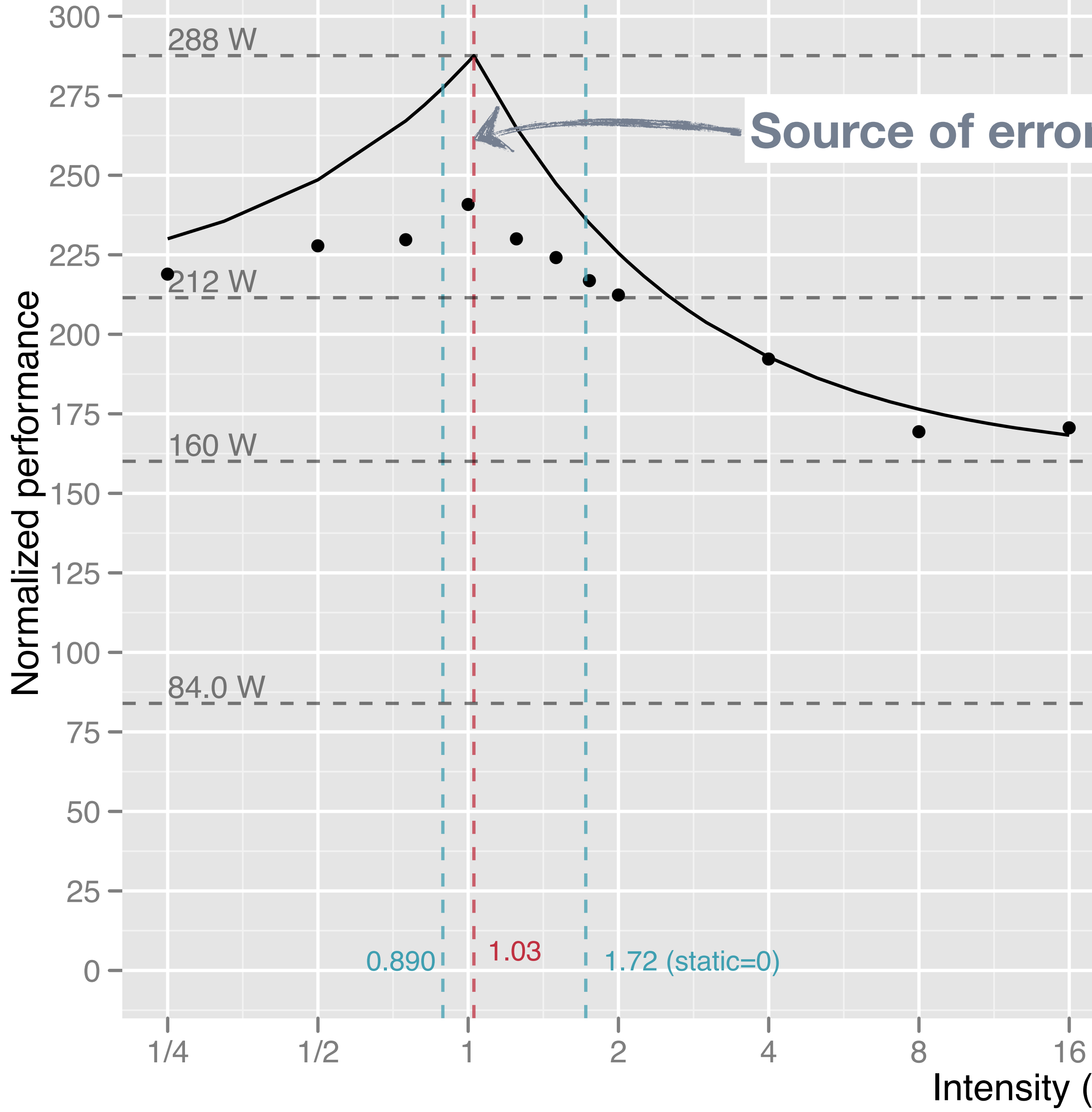
Intel i7-950 – double-precision



Power

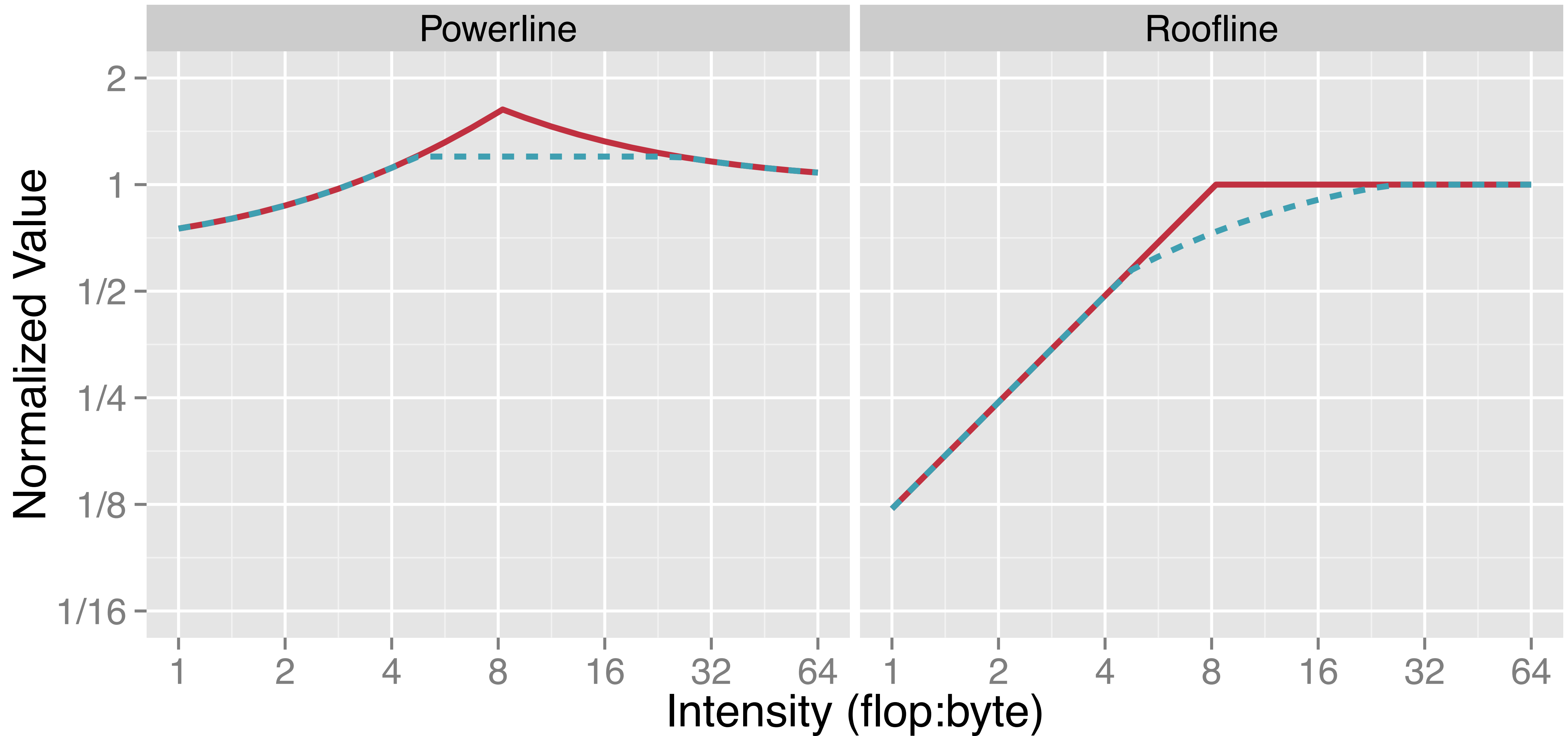
NVIDIA GTX 580 – double-precision

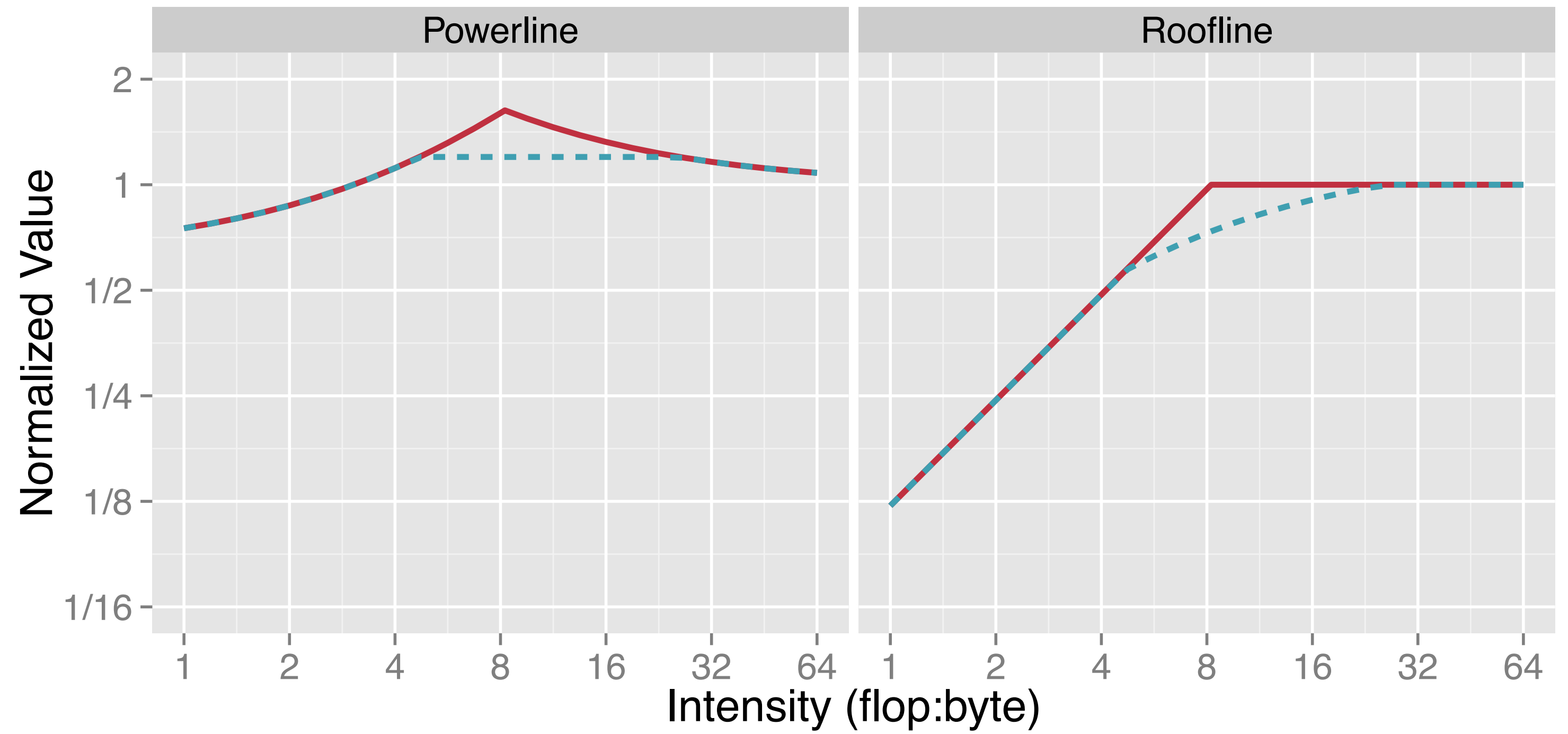
Intel i7-950 – double-precision



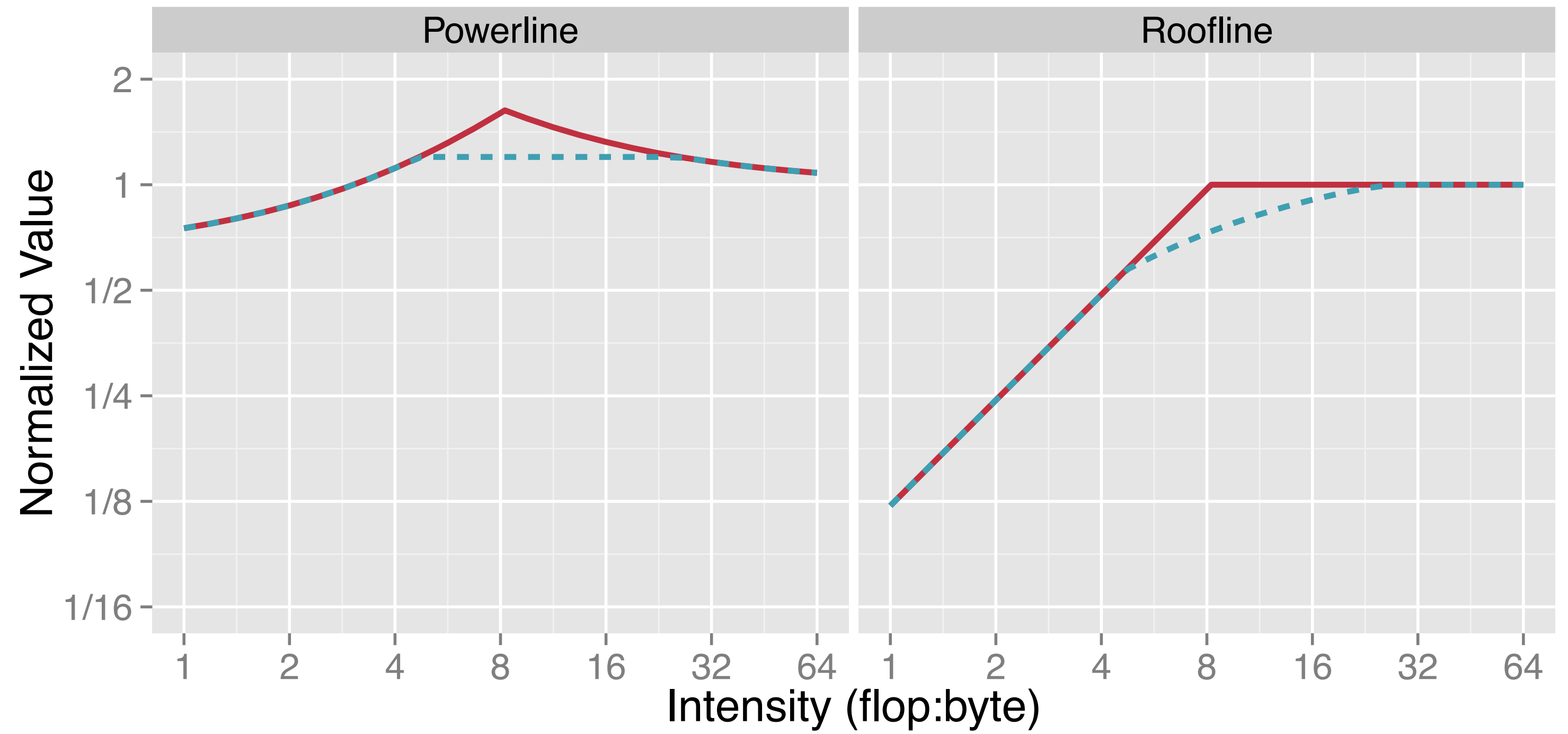
Power capping is critical.
It's also easy to add.

What might a power cap look like?





Adding a power cap

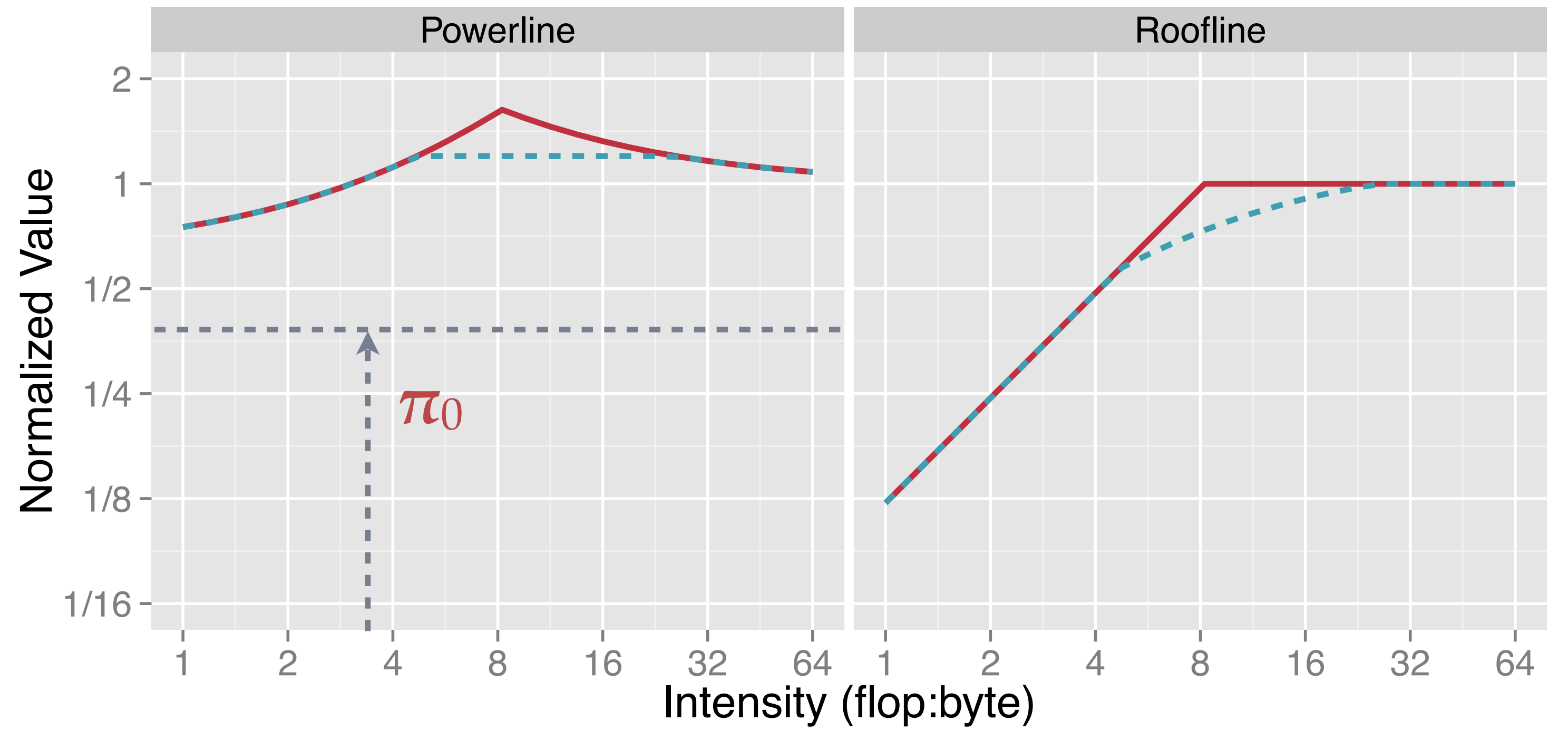


π_0 : constant

Adding a power cap



Adding a power cap

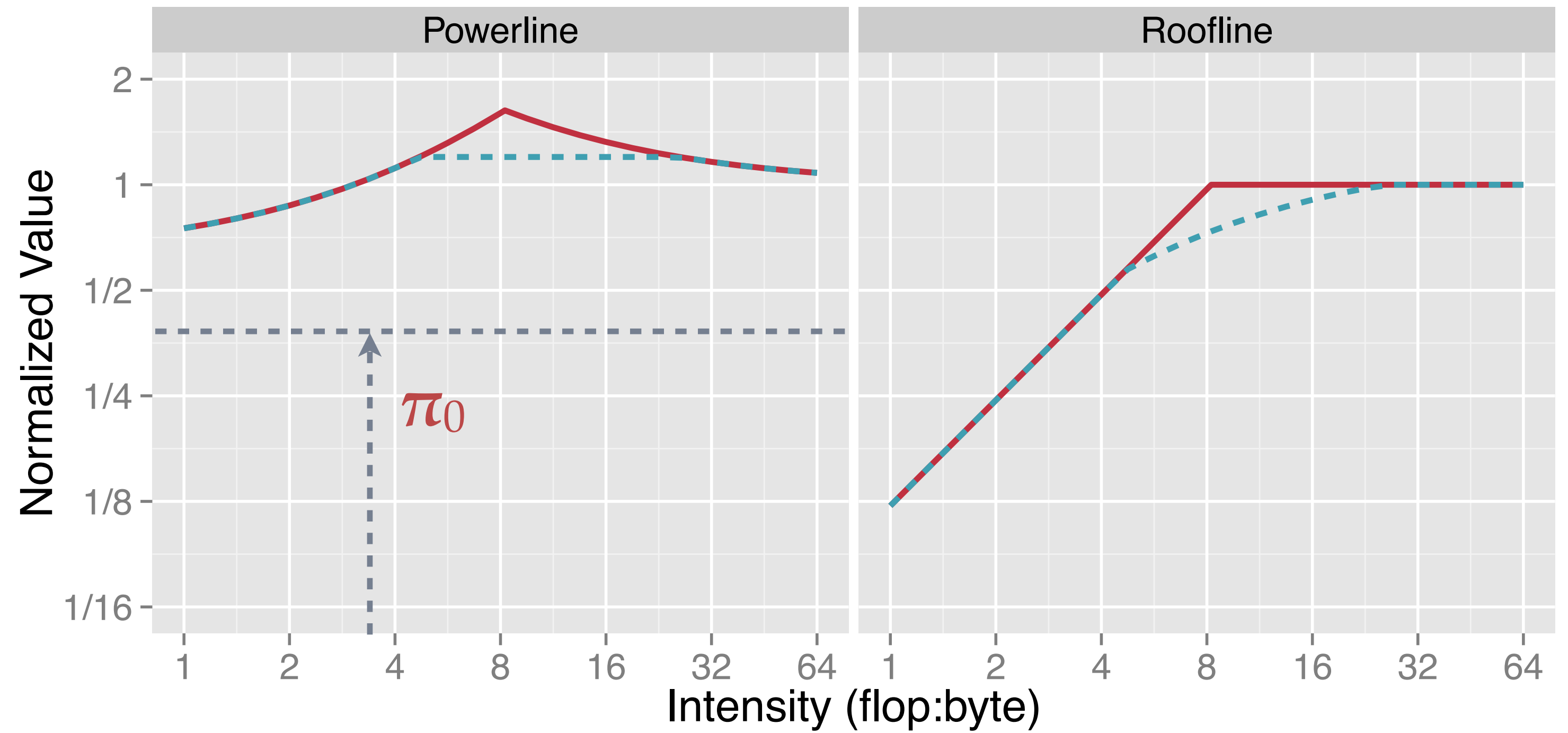


$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

Adding a power cap

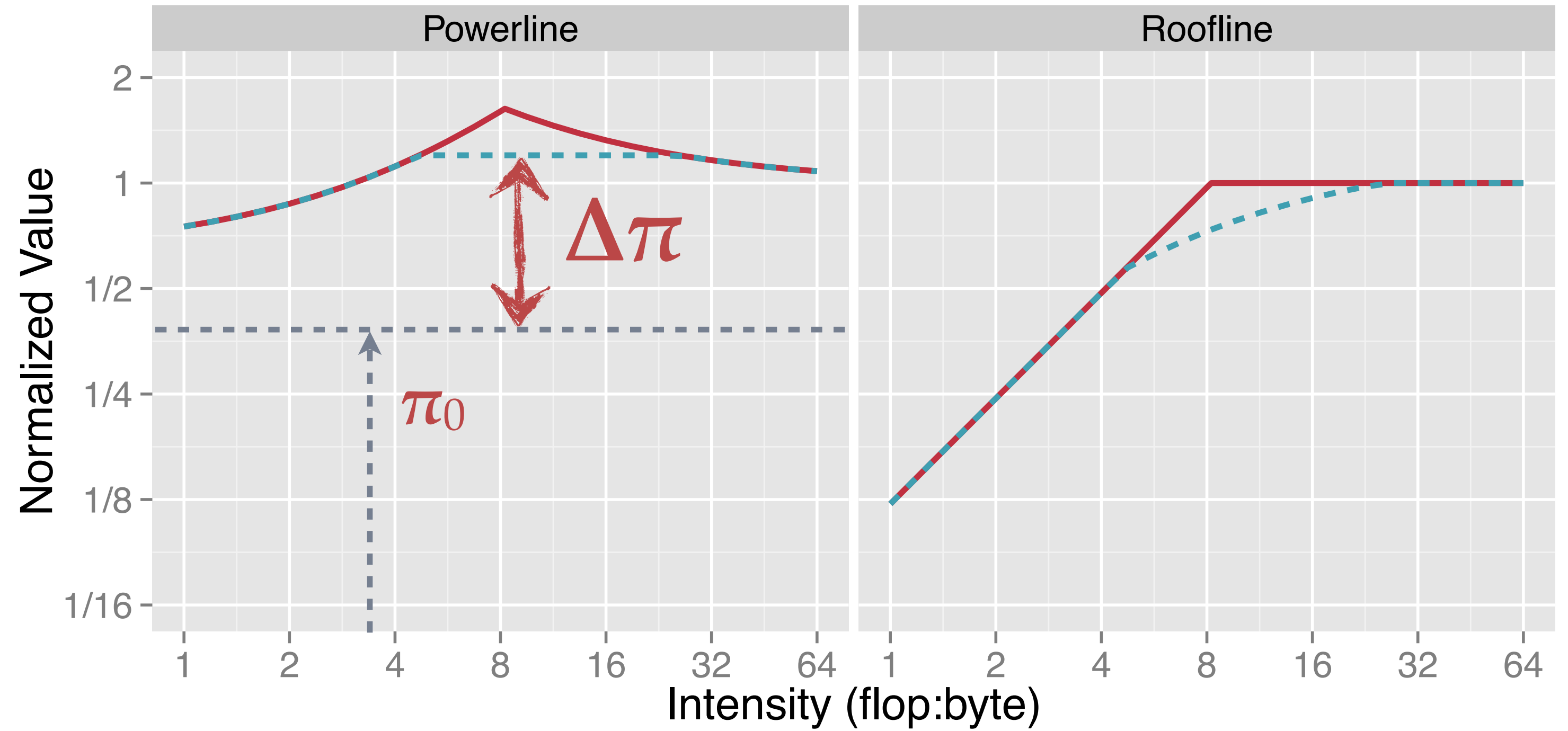


$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

Adding a power cap

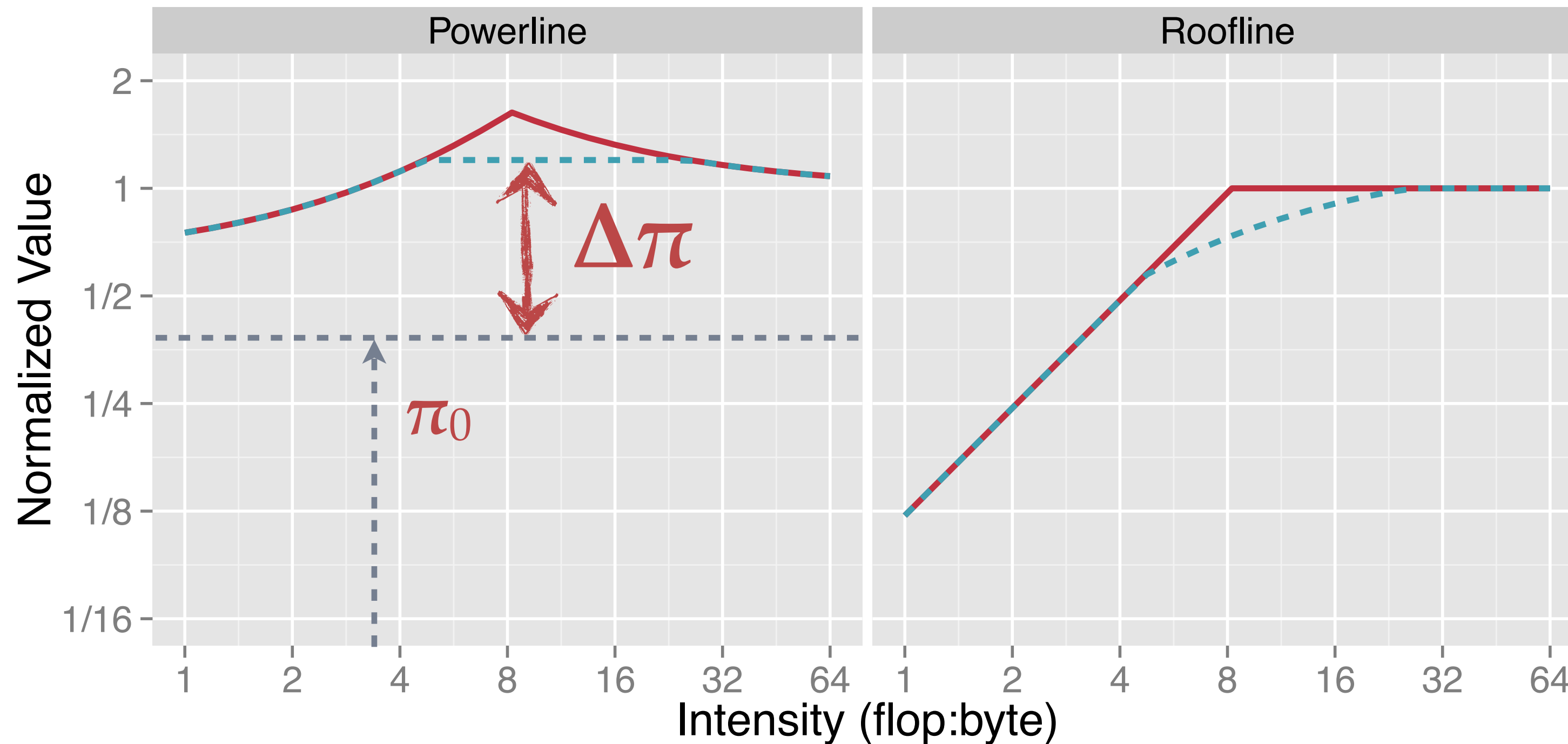


$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

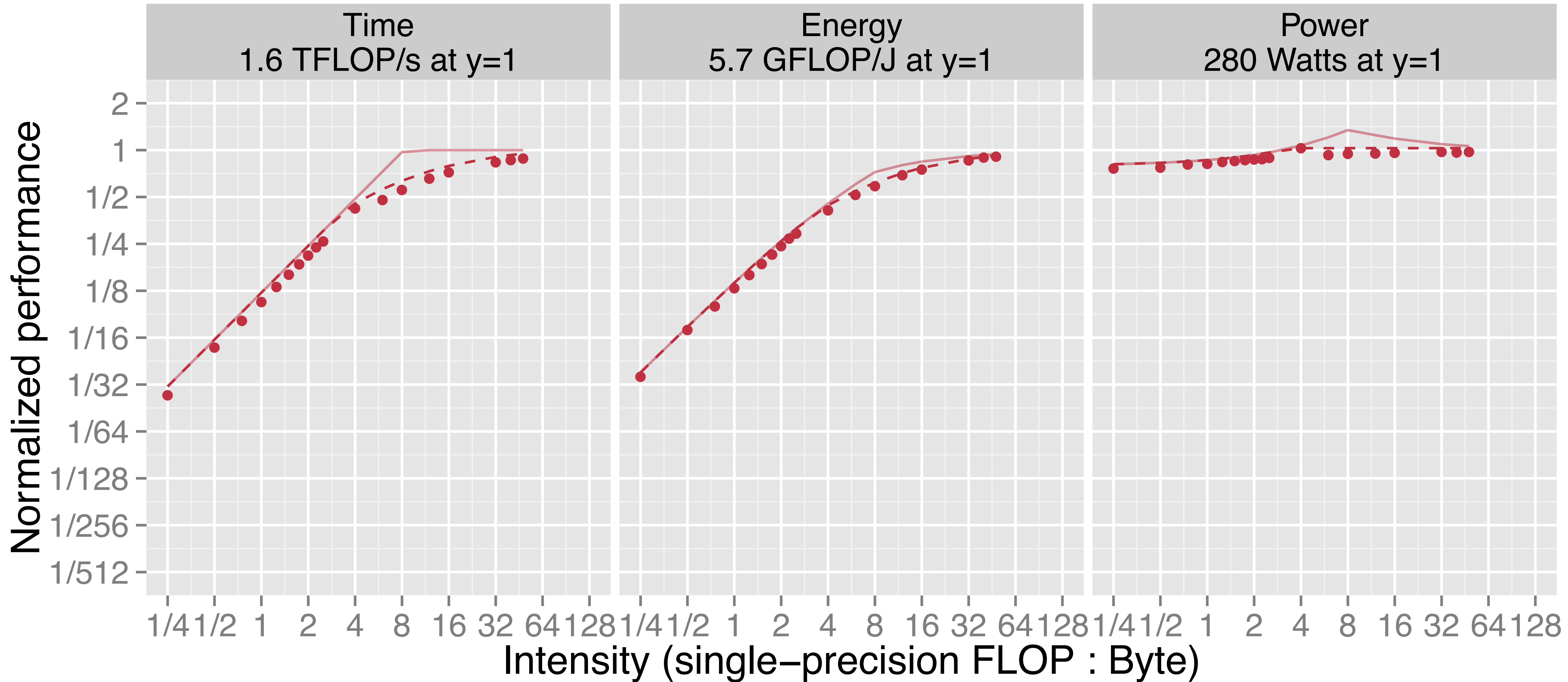
Adding a power cap



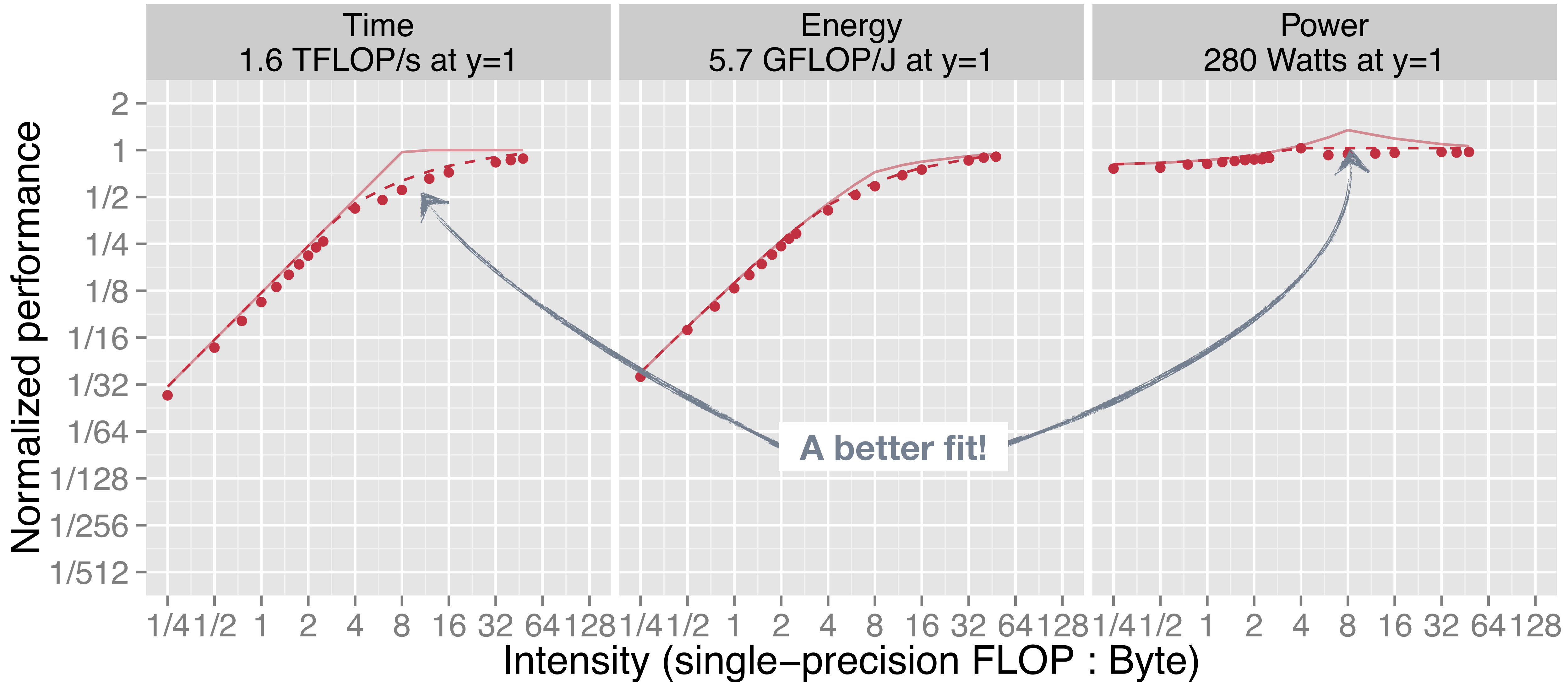
$$T_{\text{free}} = \max(W \tau_{\text{flop}}, Q \tau_{\text{mem}})$$

\Downarrow

$$T = \max\left(W \tau_{\text{flop}}, Q \tau_{\text{mem}}, \frac{W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}}}{\Delta\pi}\right)$$

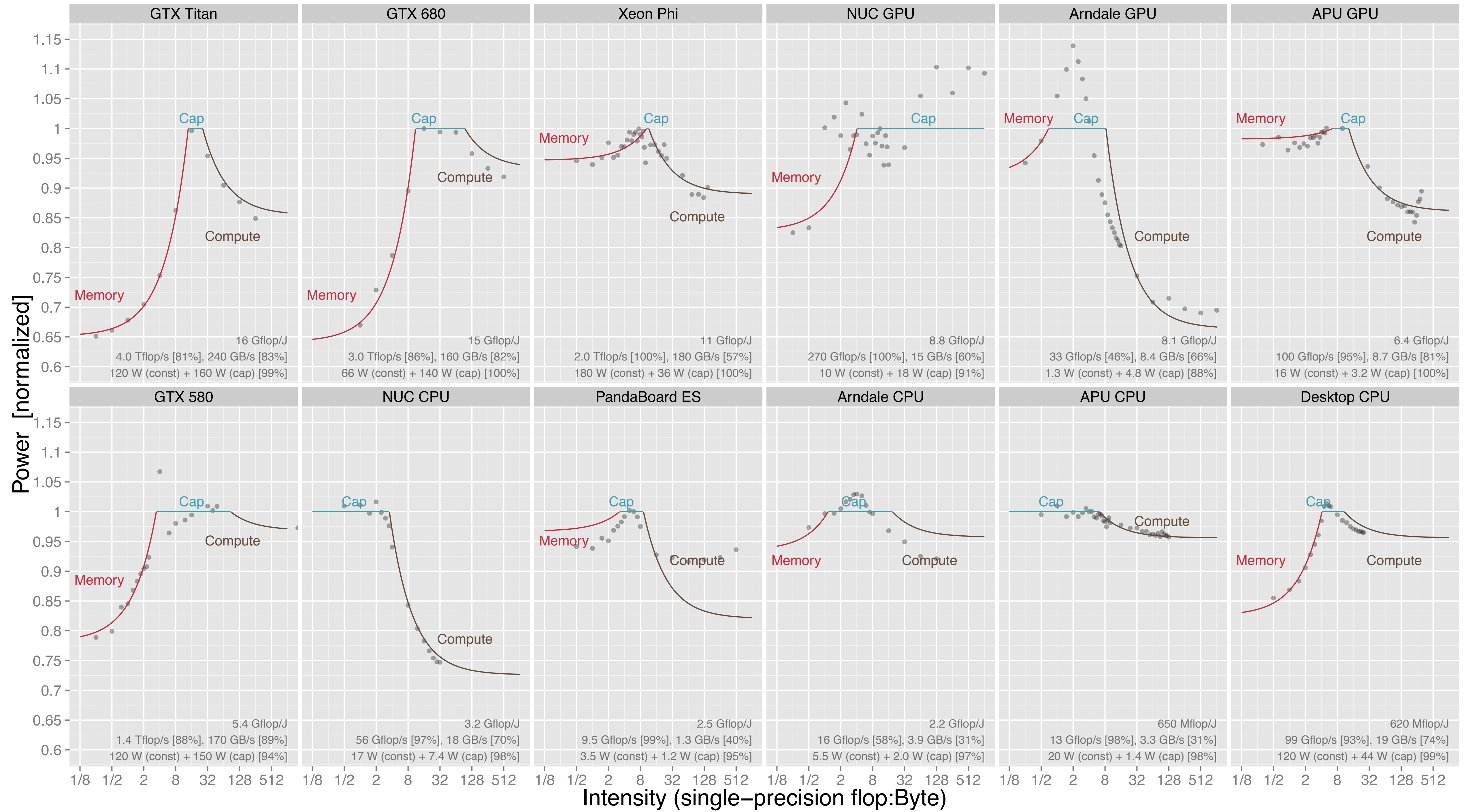


NVIDIA
 GTX 580
 GF100
 Fermi



■ NVIDIA
 ● GTX 580
 ● GF100
 ■ Fermi

The model suggests a **structure** in the time, energy, and power relationships. It also facilitates **analysis**.



$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

Example:
Caps imply throttling!

$$T = \max \left(W \tau_{\text{flop}}, Q \tau_{\text{mem}}, \frac{W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}}}{\Delta\pi} \right)$$

↓

$\tilde{\tau}_{\text{flop}}$

≡

$$\frac{T}{W} \equiv \tau_{\text{flop}} s_{\text{flop}}$$

$\tilde{\tau}_{\text{mem}}$

≡

$$\frac{T}{Q} \equiv \tau_{\text{mem}} s_{\text{mem}}$$

↓

s_{flop}

≡

$$\max \left\{ 1, \frac{B_\tau}{I}, \frac{\epsilon_{\text{flop}} / \tau_{\text{flop}}}{\Delta\pi} \left(1 + \frac{B_\epsilon}{I} \right) \right\}$$

s_{mem}

≡

$$s_{\text{flop}} \frac{I}{B_\tau}$$

Throttling factors
i.e., allowable slowdown

$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

Example:
Caps imply throttling!

$$T = \max \left(W \tau_{\text{flop}}, Q \tau_{\text{mem}}, \frac{W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}}}{\Delta\pi} \right)$$

↓

$$\tilde{\tau}_{\text{flop}} \equiv \frac{T}{W} \equiv \tau_{\text{flop}} s_{\text{flop}}$$

$$\tilde{\tau}_{\text{mem}} \equiv \frac{T}{Q} \equiv \tau_{\text{mem}} s_{\text{mem}}$$

Throttling factors
i.e., allowable slowdown

↓

$$s_{\text{flop}} \equiv \max \left\{ 1, \frac{B_\tau}{I}, \frac{\epsilon_{\text{flop}} / \tau_{\text{flop}}}{\Delta\pi} \left(1 + \frac{B_\epsilon}{I} \right) \right\}$$

$$s_{\text{mem}} \equiv s_{\text{flop}} \frac{I}{B_\tau}$$

$\pi_0 + \Delta\pi$: max power

$\Delta\pi$: usable

π_0 : constant

Example:
Caps imply throttling!

$$T = \max \left(W \tau_{\text{flop}}, Q \tau_{\text{mem}}, \frac{W \epsilon_{\text{flop}} + Q \epsilon_{\text{mem}}}{\Delta\pi} \right)$$

\Downarrow

$$\tilde{\tau}_{\text{flop}} \equiv \frac{T}{W} \equiv \tau_{\text{flop}} s_{\text{flop}}$$

$$\tilde{\tau}_{\text{mem}} \equiv \frac{T}{Q} \equiv \tau_{\text{mem}} s_{\text{mem}}$$

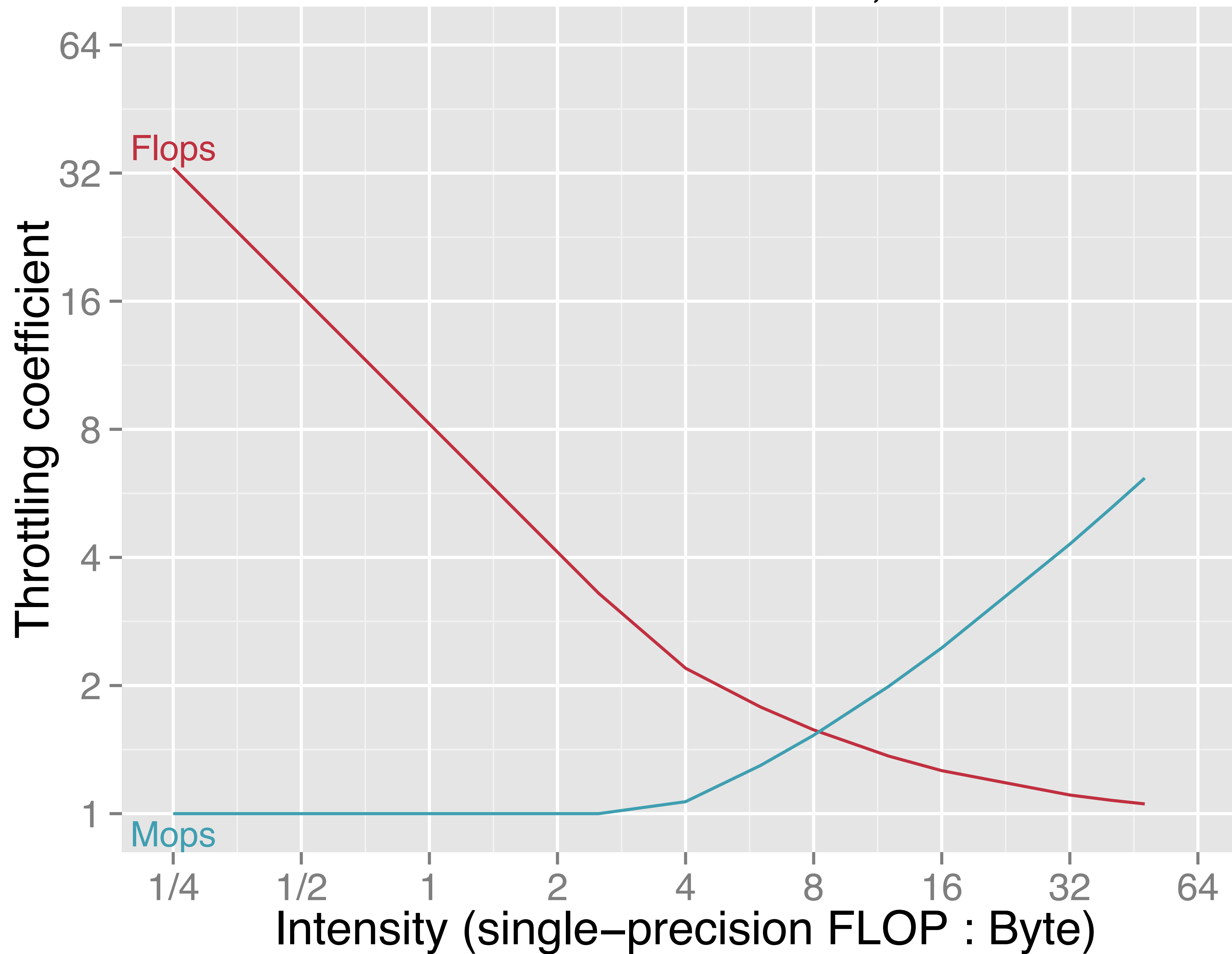
\Downarrow

$$s_{\text{flop}} \equiv \max \left\{ 1, \frac{B_\tau}{I}, \frac{\epsilon_{\text{flop}} / \tau_{\text{flop}}}{\Delta\pi} \left(1 + \frac{B_\epsilon}{I} \right) \right\}$$

$$s_{\text{mem}} \equiv s_{\text{flop}} \frac{I}{B_\tau}$$

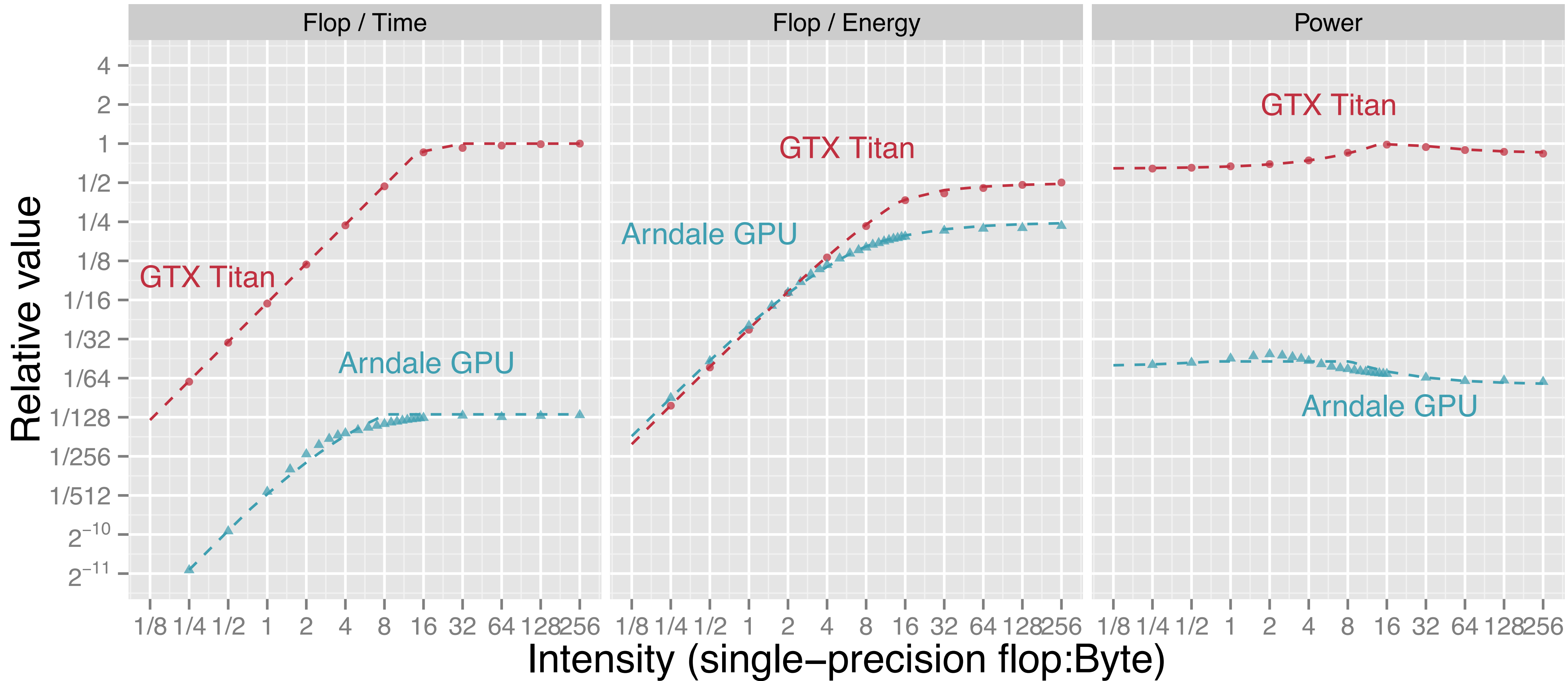
Throttling factors
i.e., allowable slowdown

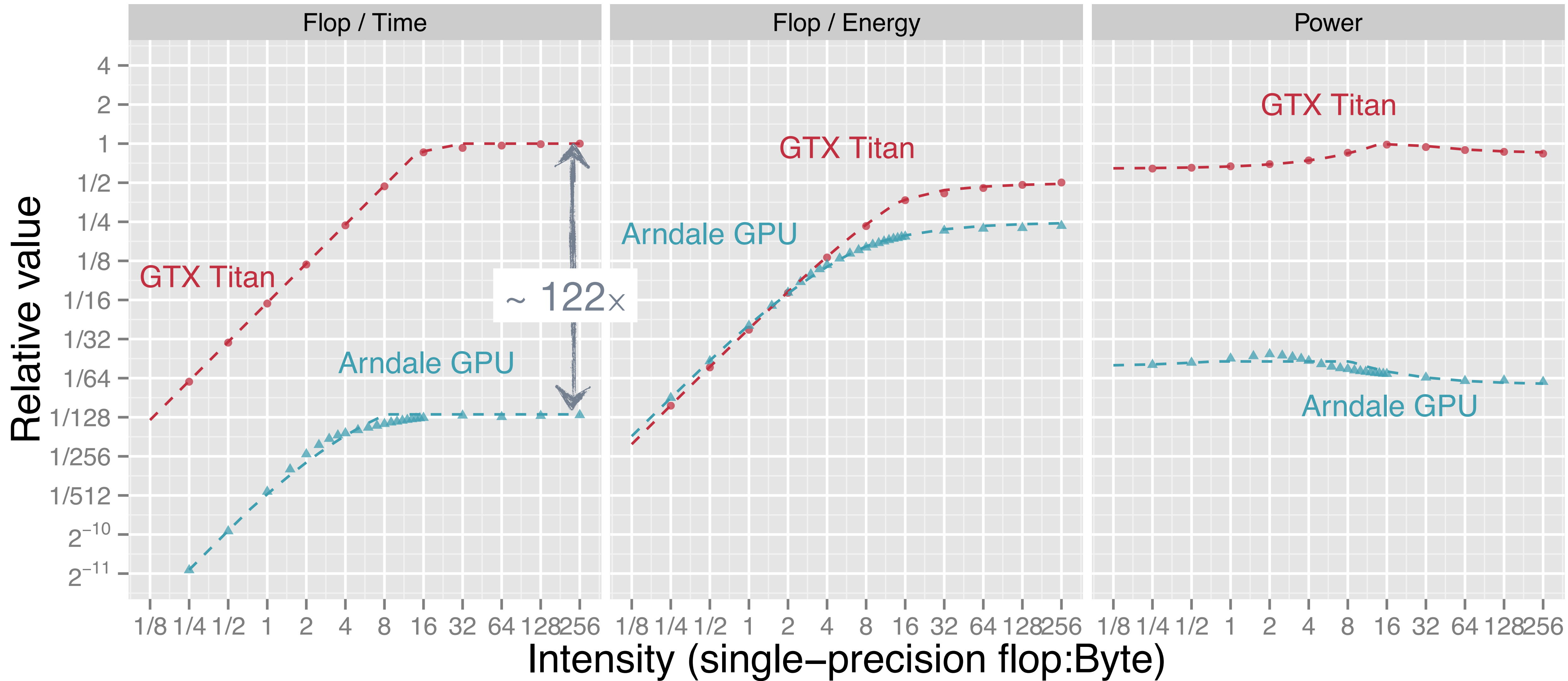
NVIDIA GTX 580 GF100; Fermi

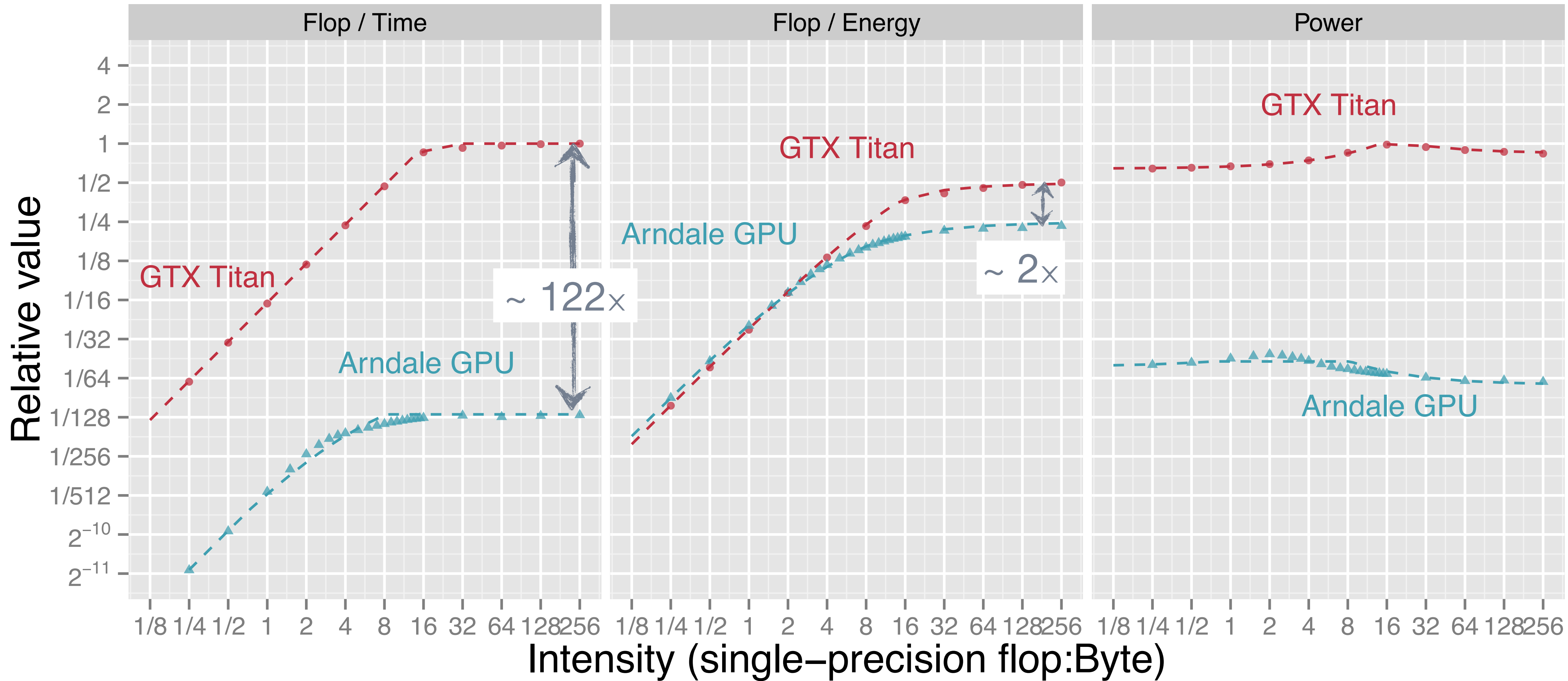


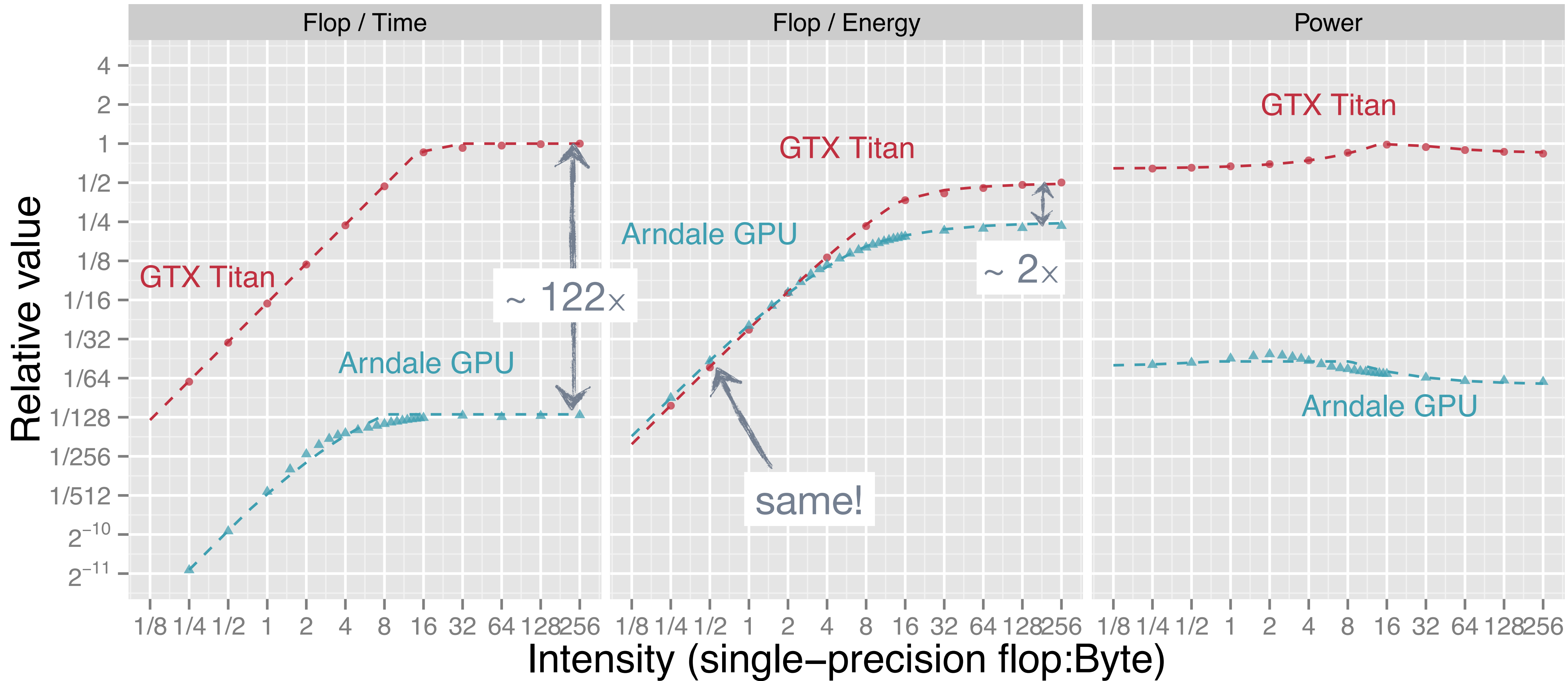
A “what if ...” comparison of system building blocks:

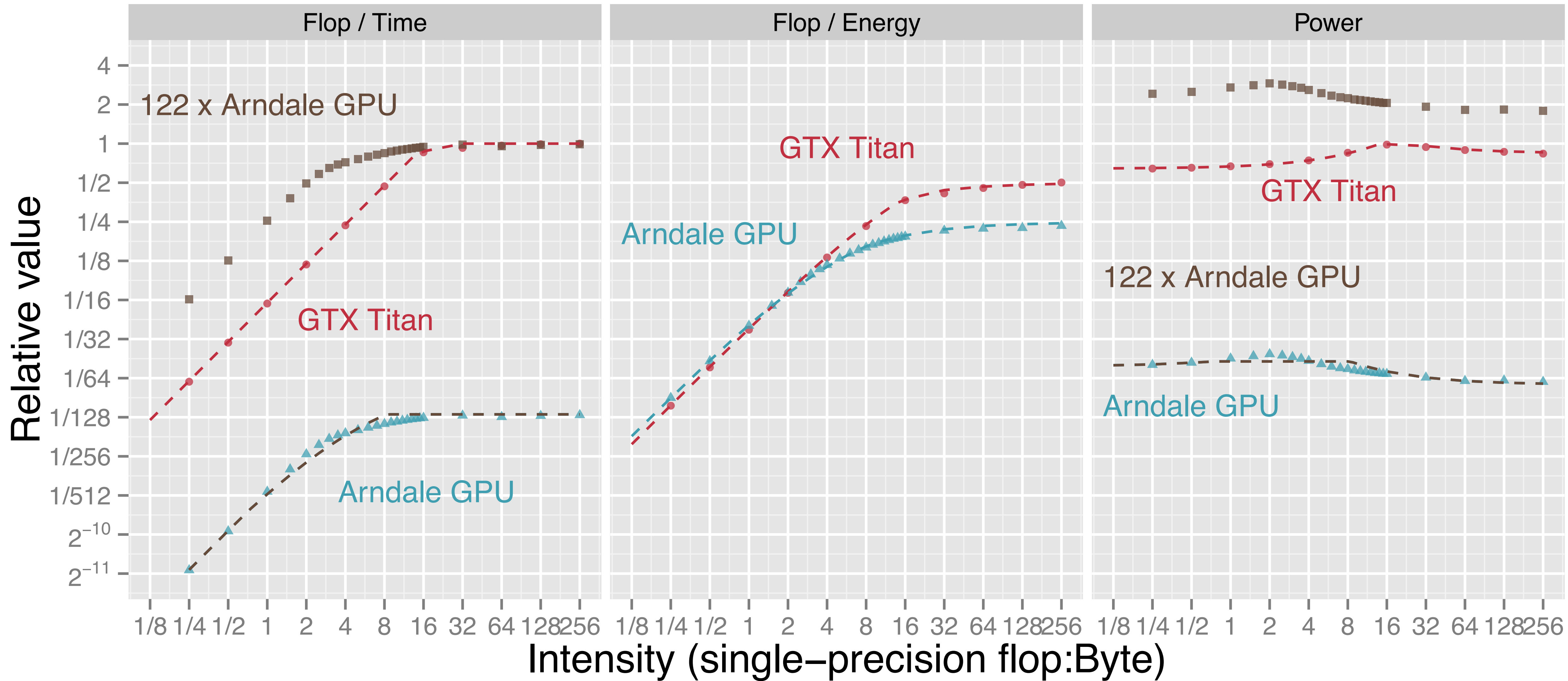
→ **GTX Titan** vs. **Samsung Exynos 5** (GPU only; Arndale)

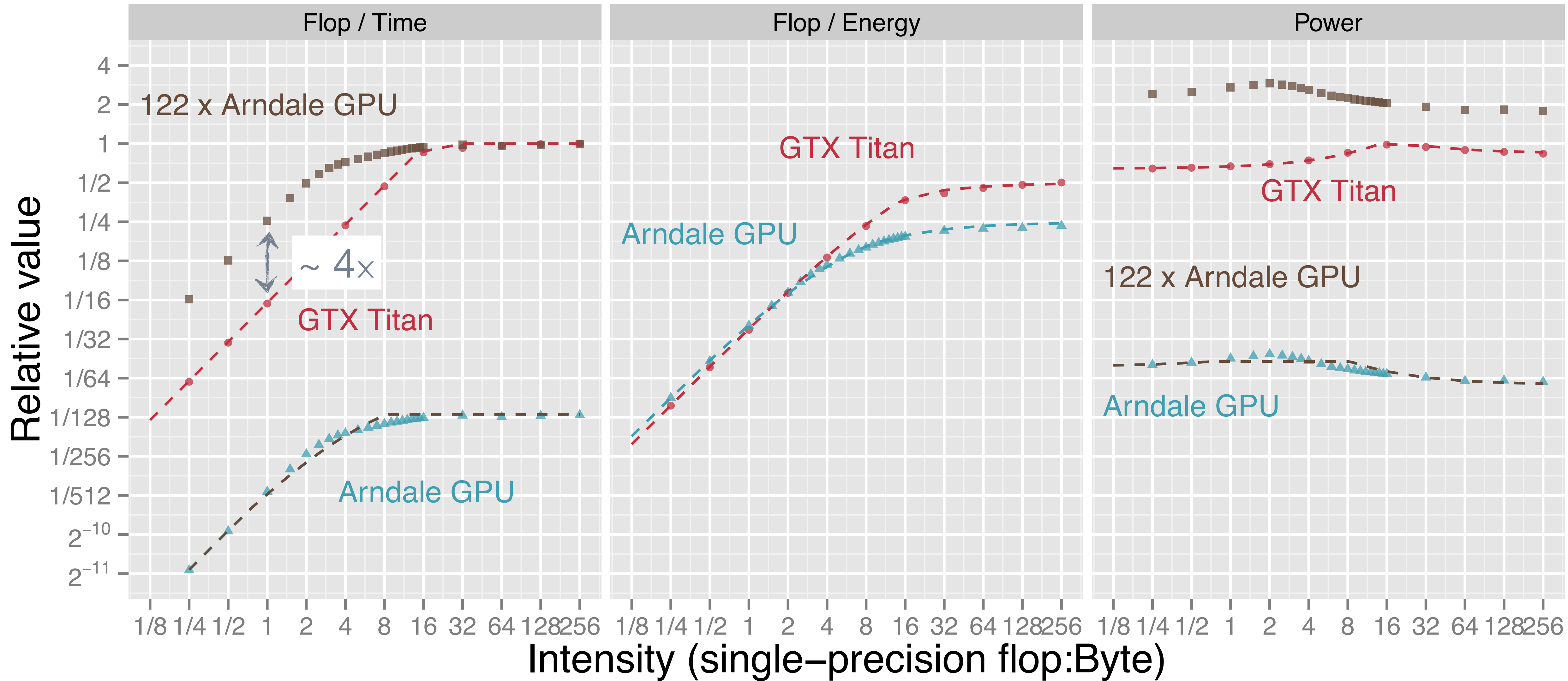


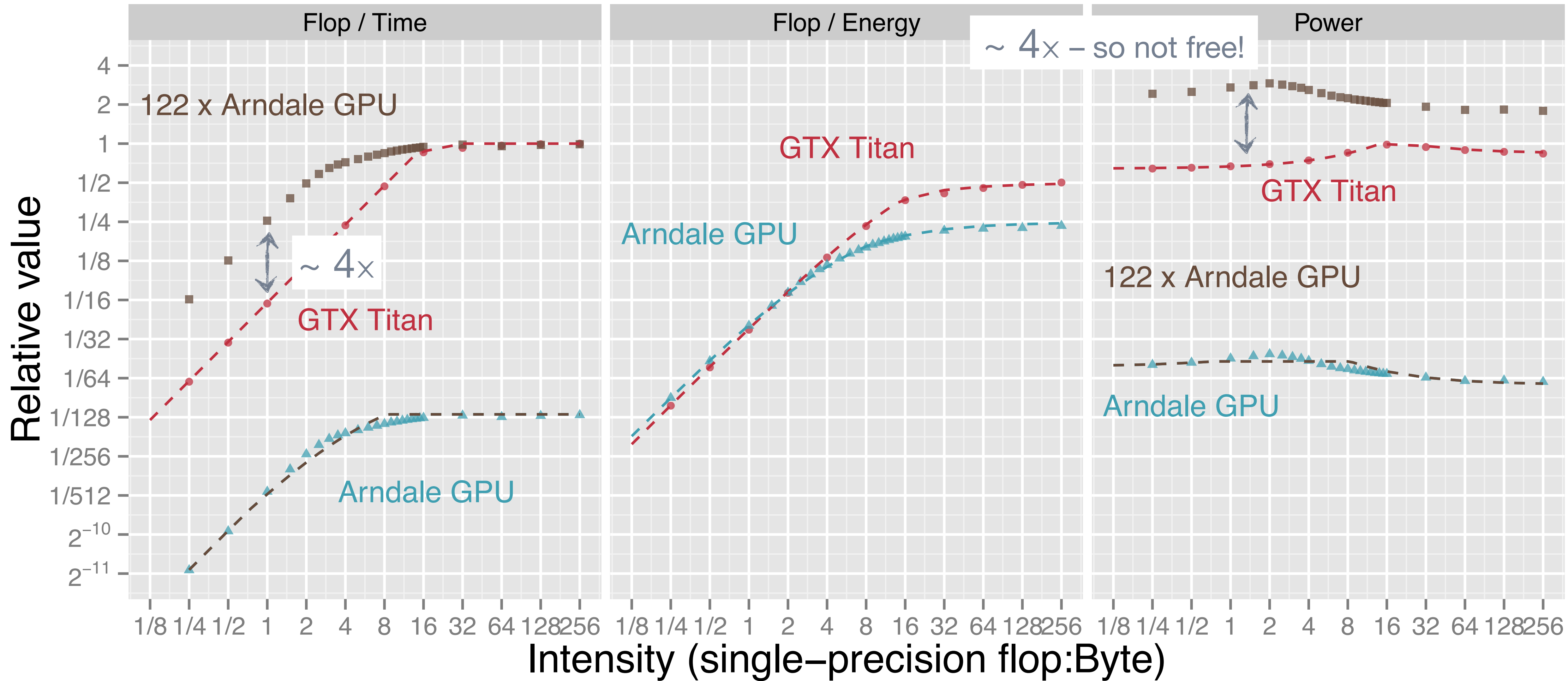












Rebutting Hillis, Part 2:

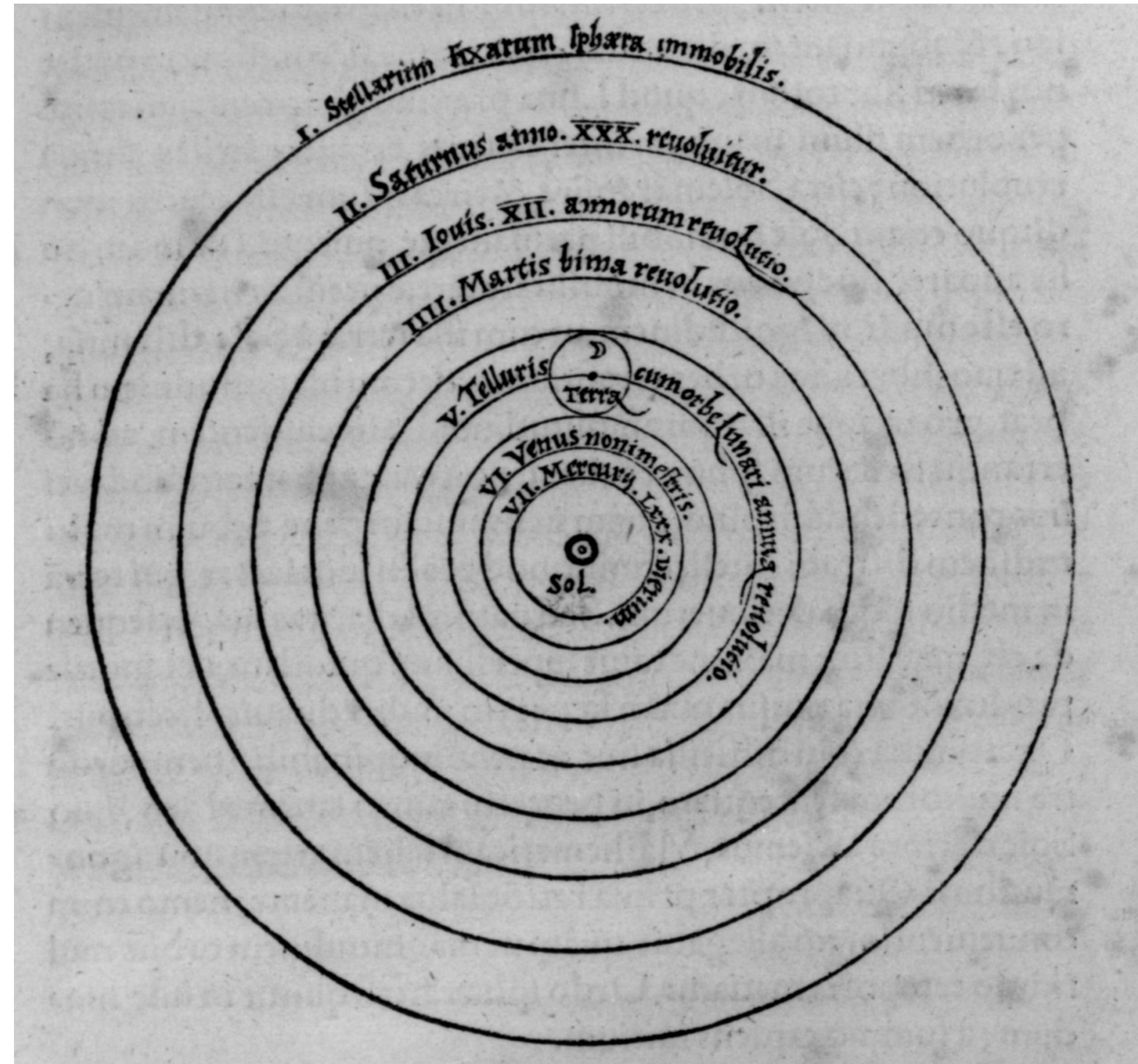
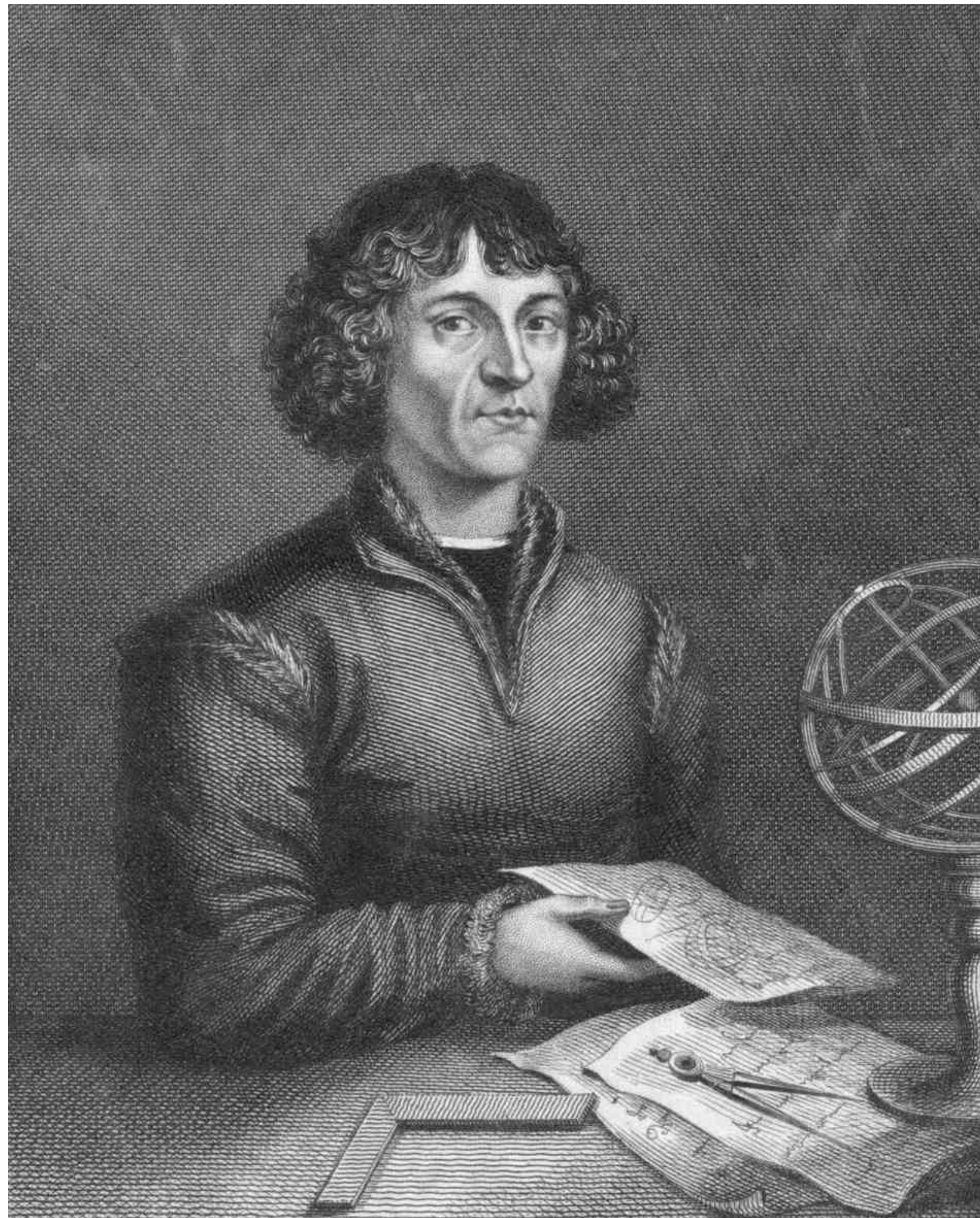


Kent Czechowski

“Algorithmic” power and die-area constraints

Can first principles of algorithms direct architectures toward bigger wins?

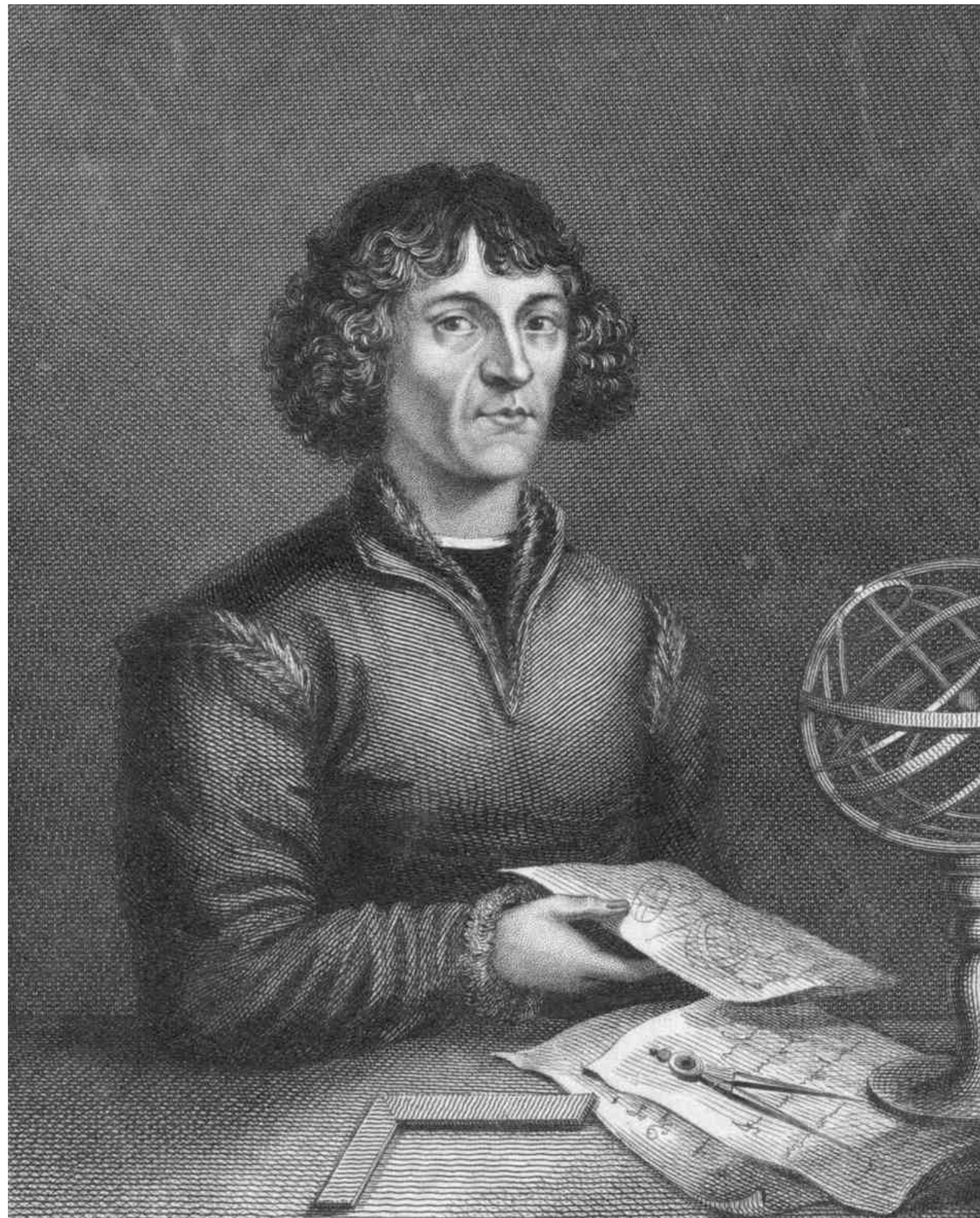
Kent: “What problem are we really trying to solve?”



http://media.npr.org/assets/img/2011/11/02/copernicus_custom-54ef1ed5bcadf2be8a9c44d36dc7b273696946b2-s6-c30.jpg

<http://www.hps.cam.ac.uk/starry/copernicuslrg.jpg>

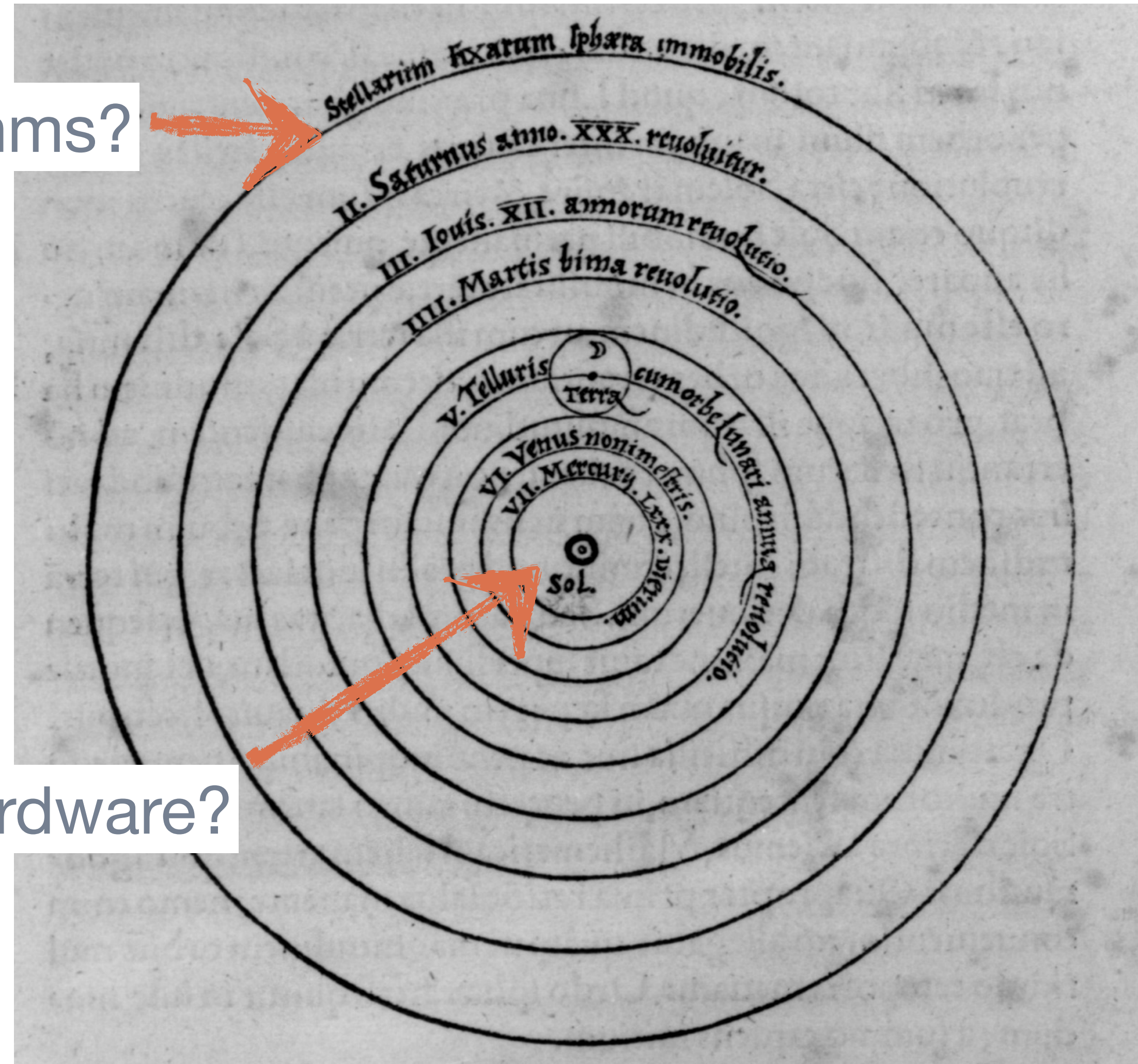
Kent: “What problem are we really trying to solve?”



Algorithms?

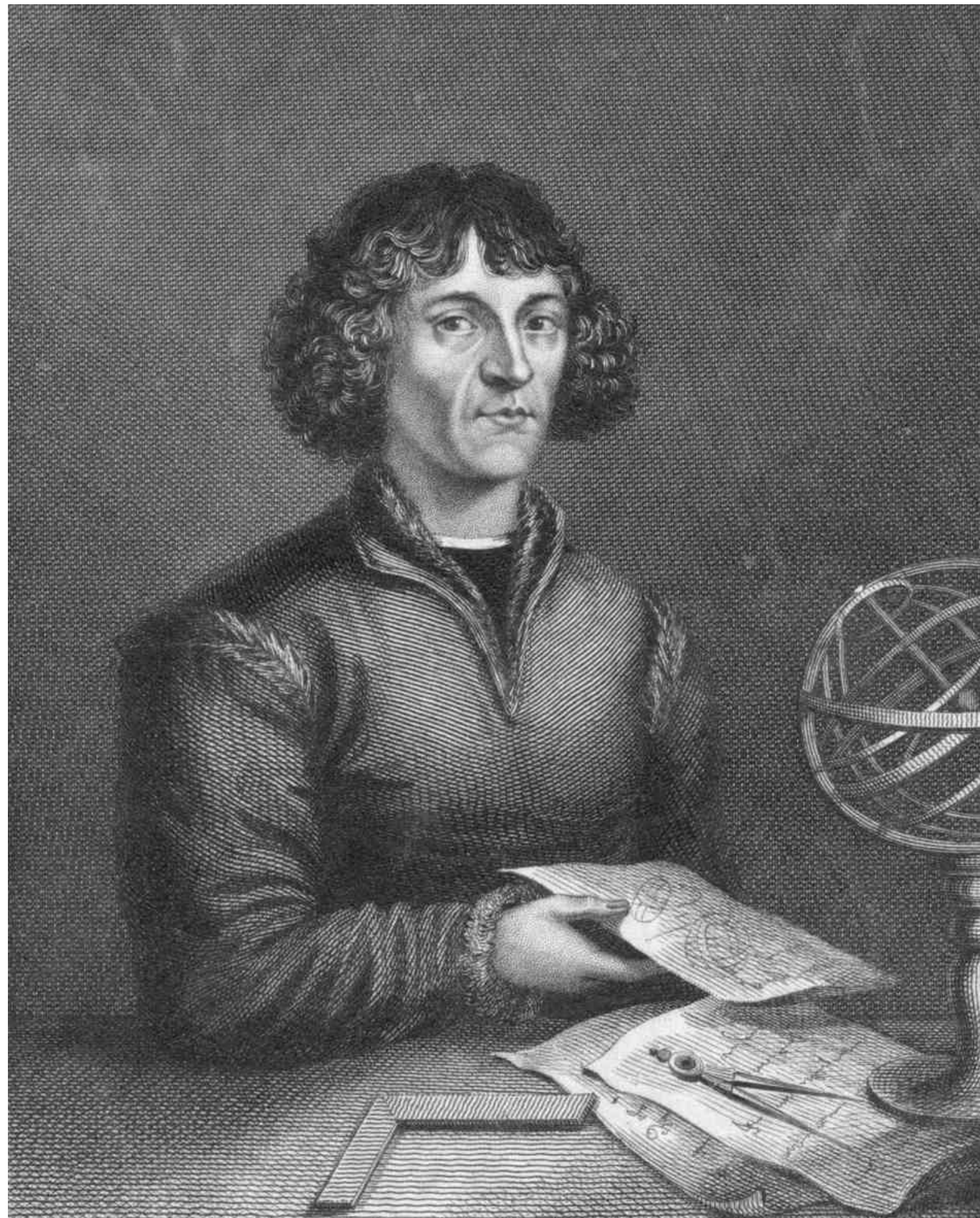


Hardware?



http://media.npr.org/assets/img/2011/11/02/copernicus_custom-54ef1ed5bcadf2be8a9c44d36dc7b273696946b2-s6-c30.jpg
<http://www.hps.cam.ac.uk/starry/copernicus1rg.jpg>

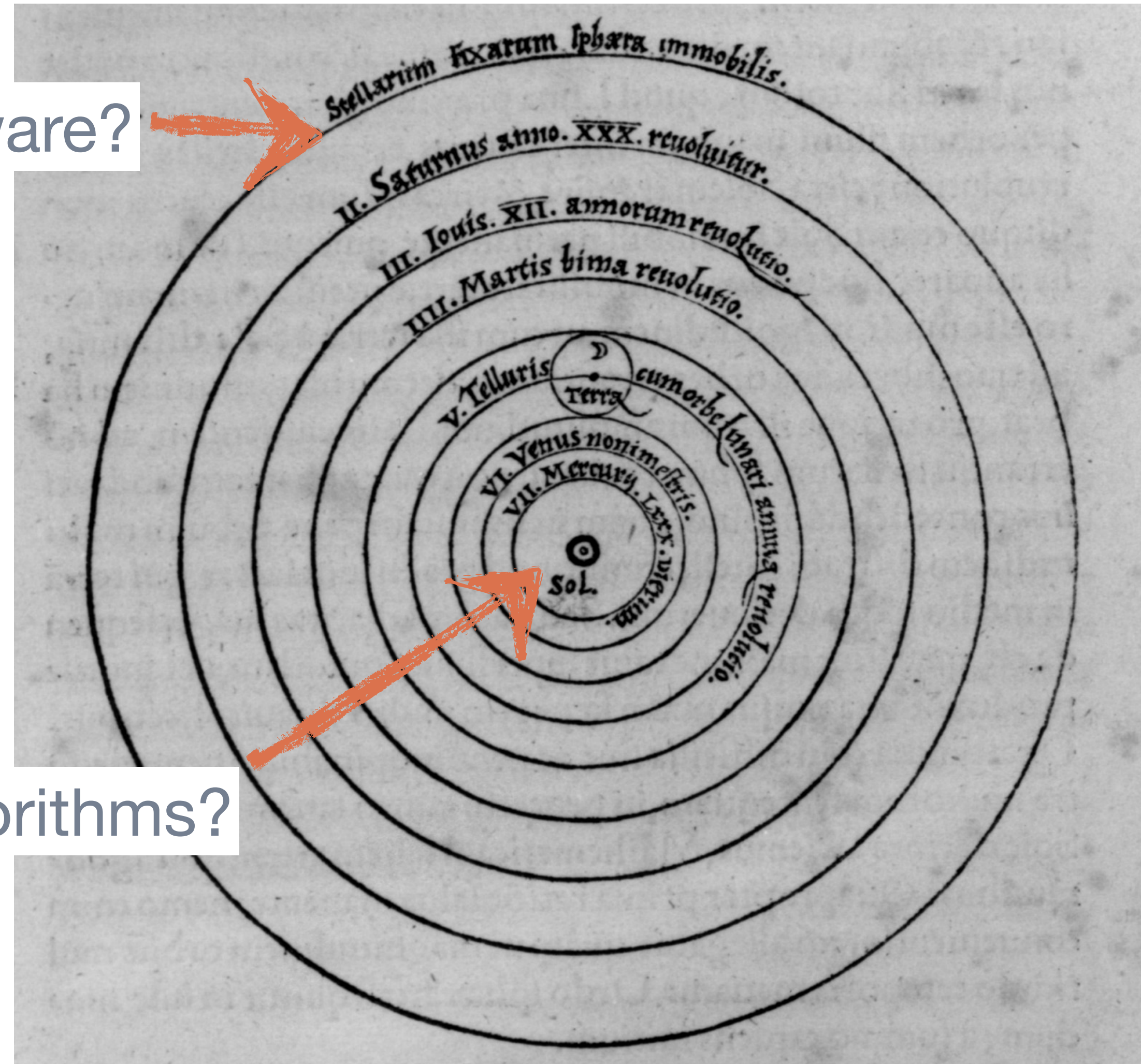
Kent: “What problem are we really trying to solve?”



Hardware?



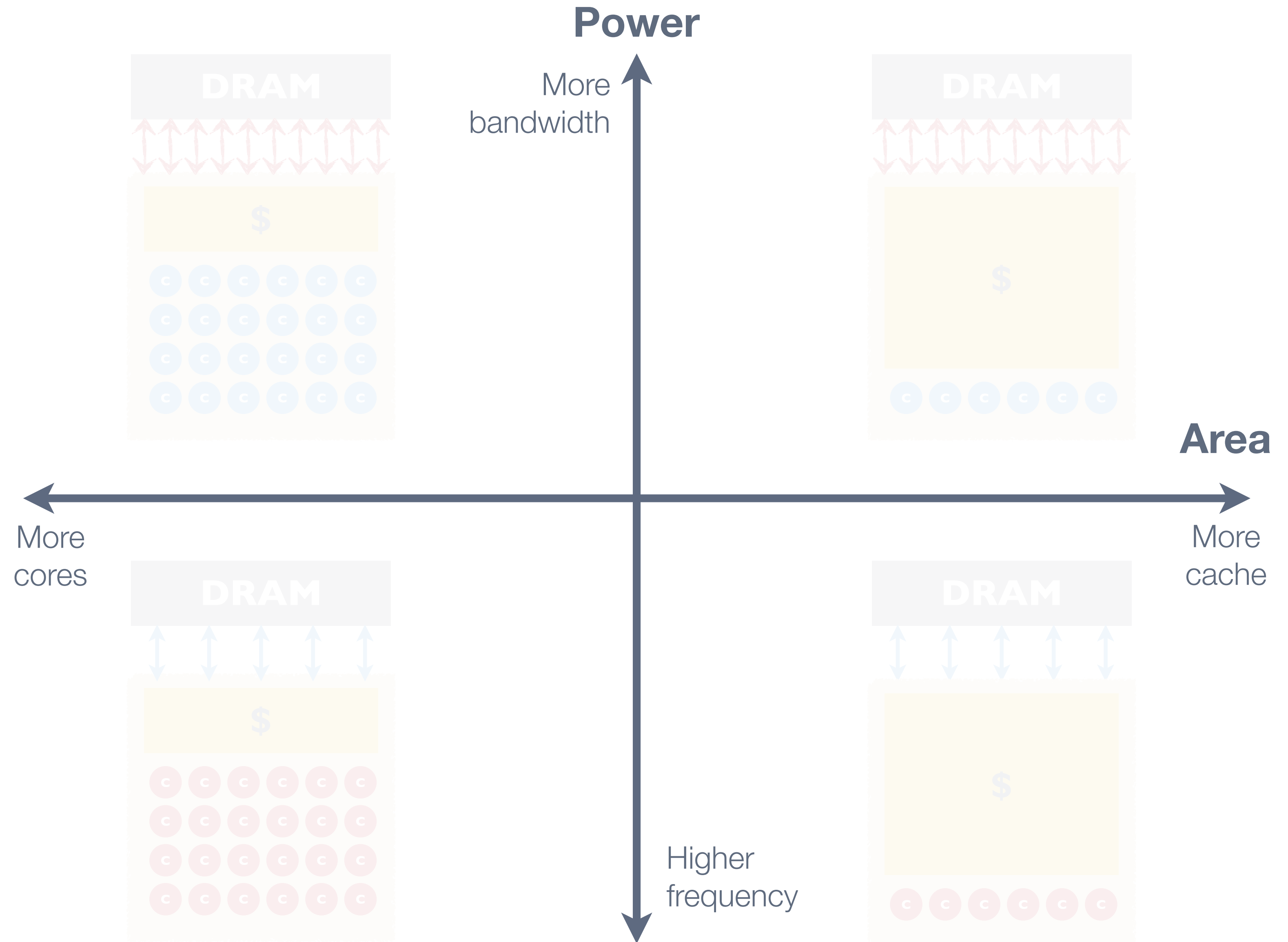
Algorithms?



http://media.npr.org/assets/img/2011/11/02/copernicus_custom-54ef1ed5bcadf2be8a9c44d36dc7b273696946b2-s6-c30.jpg
<http://www.hps.cam.ac.uk/starry/copernicus1rg.jpg>

A Notional Design Problem (1 processor-node example)

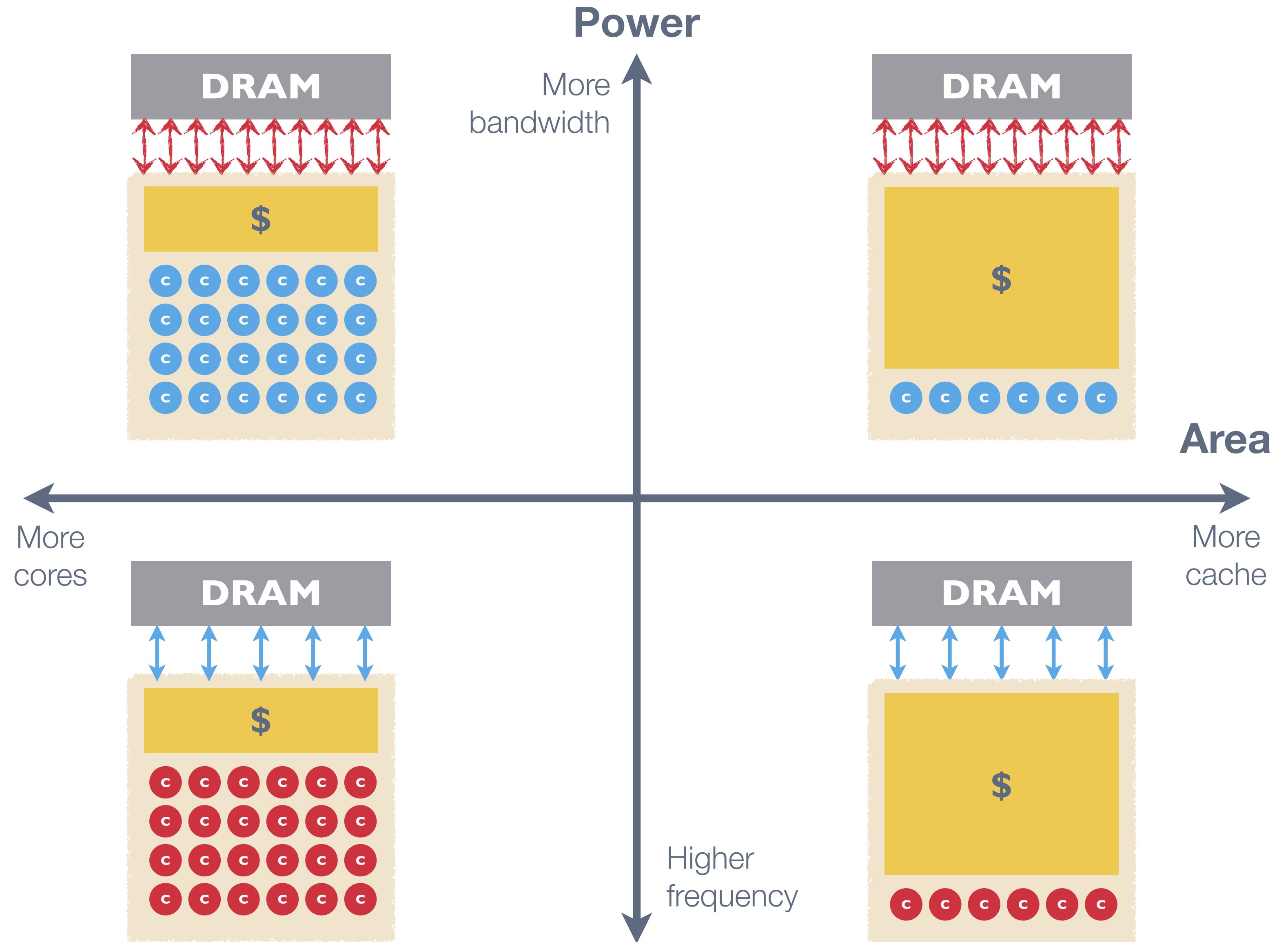
For fixed:
* **die area** (transistors)
* **power** budget
* computation



A Notional Design Problem (1 processor-node example)

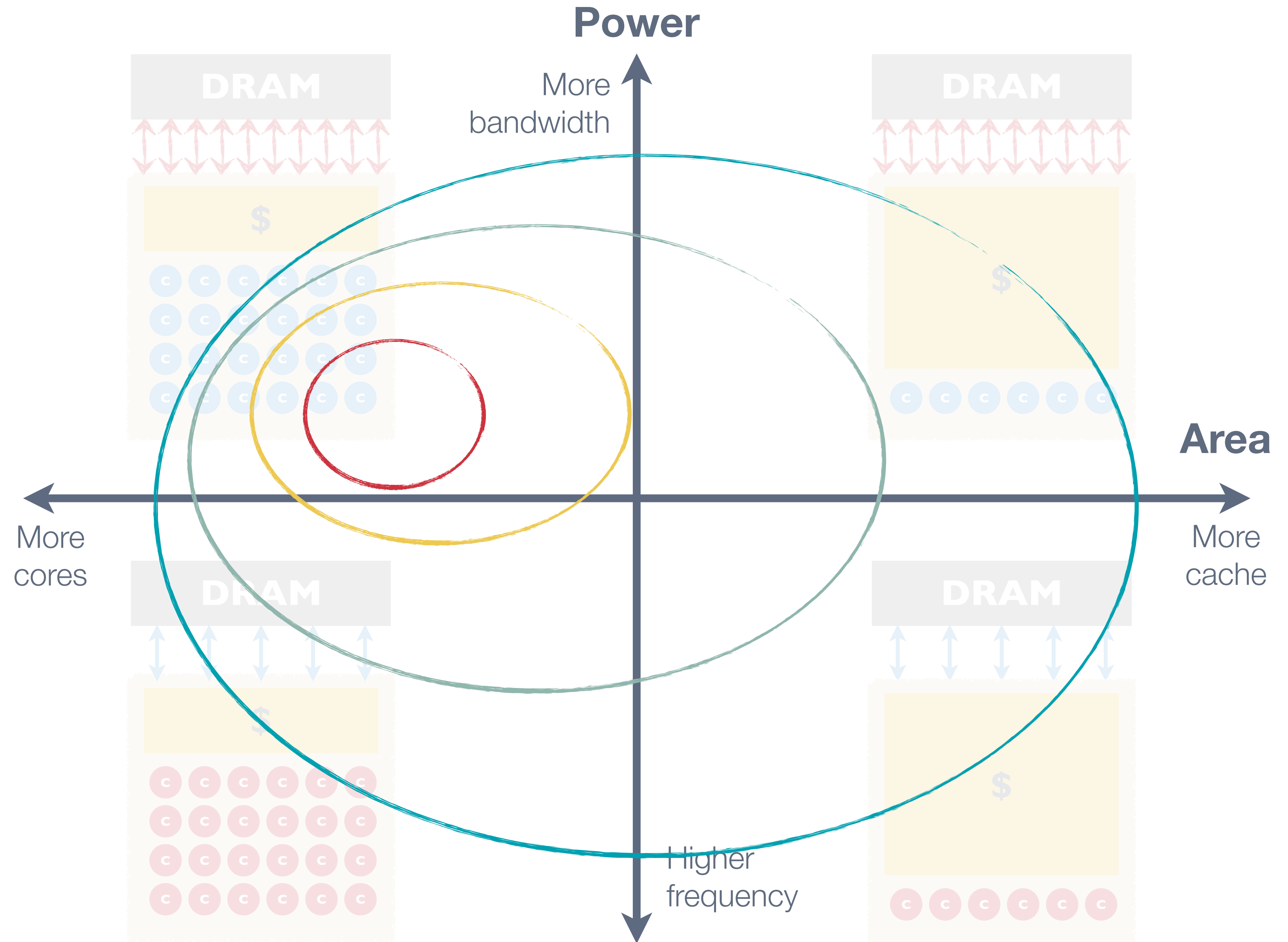
For fixed:

- * **die area** (transistors)
- * **power** budget
- * **computation**



A Notional Design Problem (1 processor-node example)

For fixed:
* die area (transistors)
* power budget
* **computation**

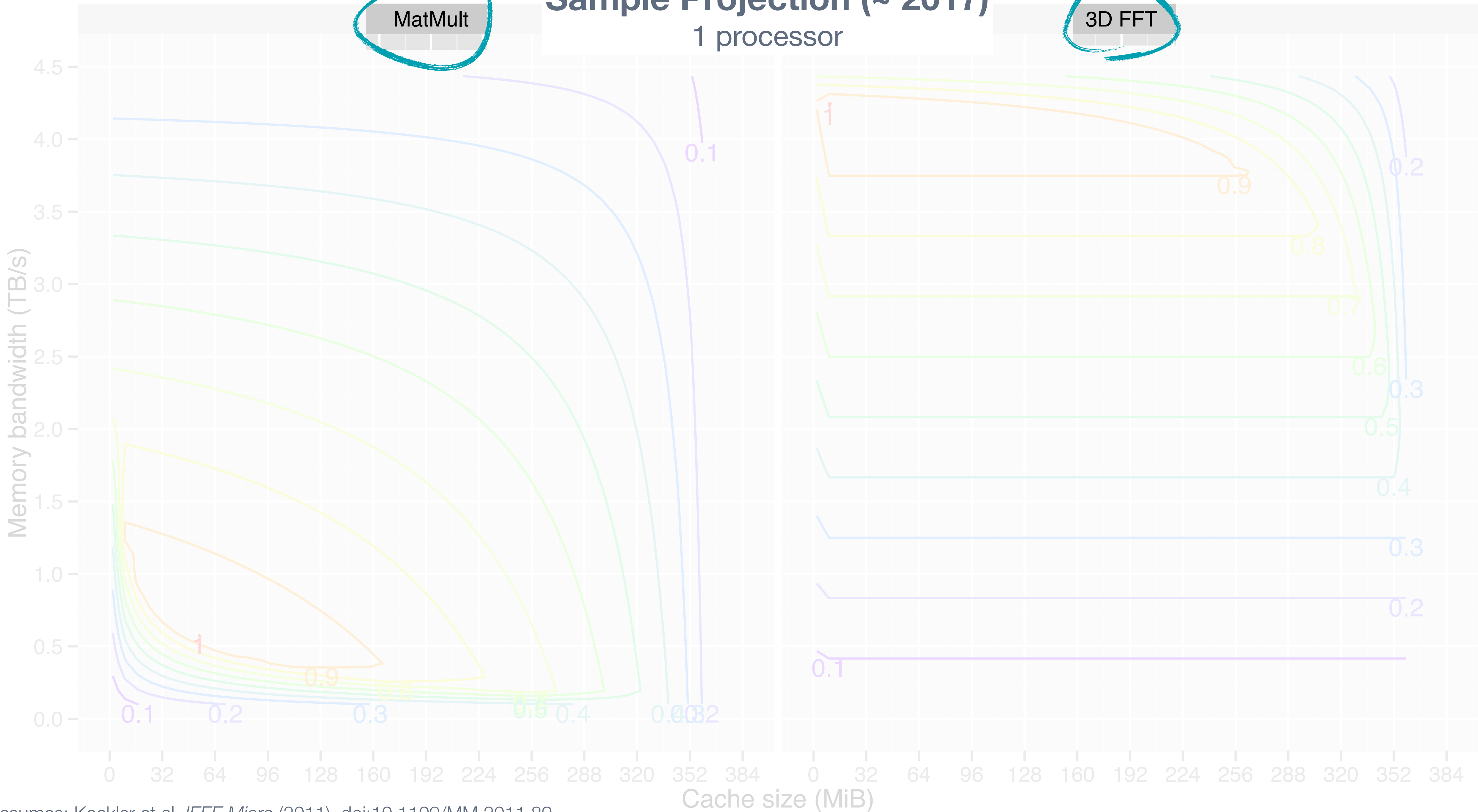


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

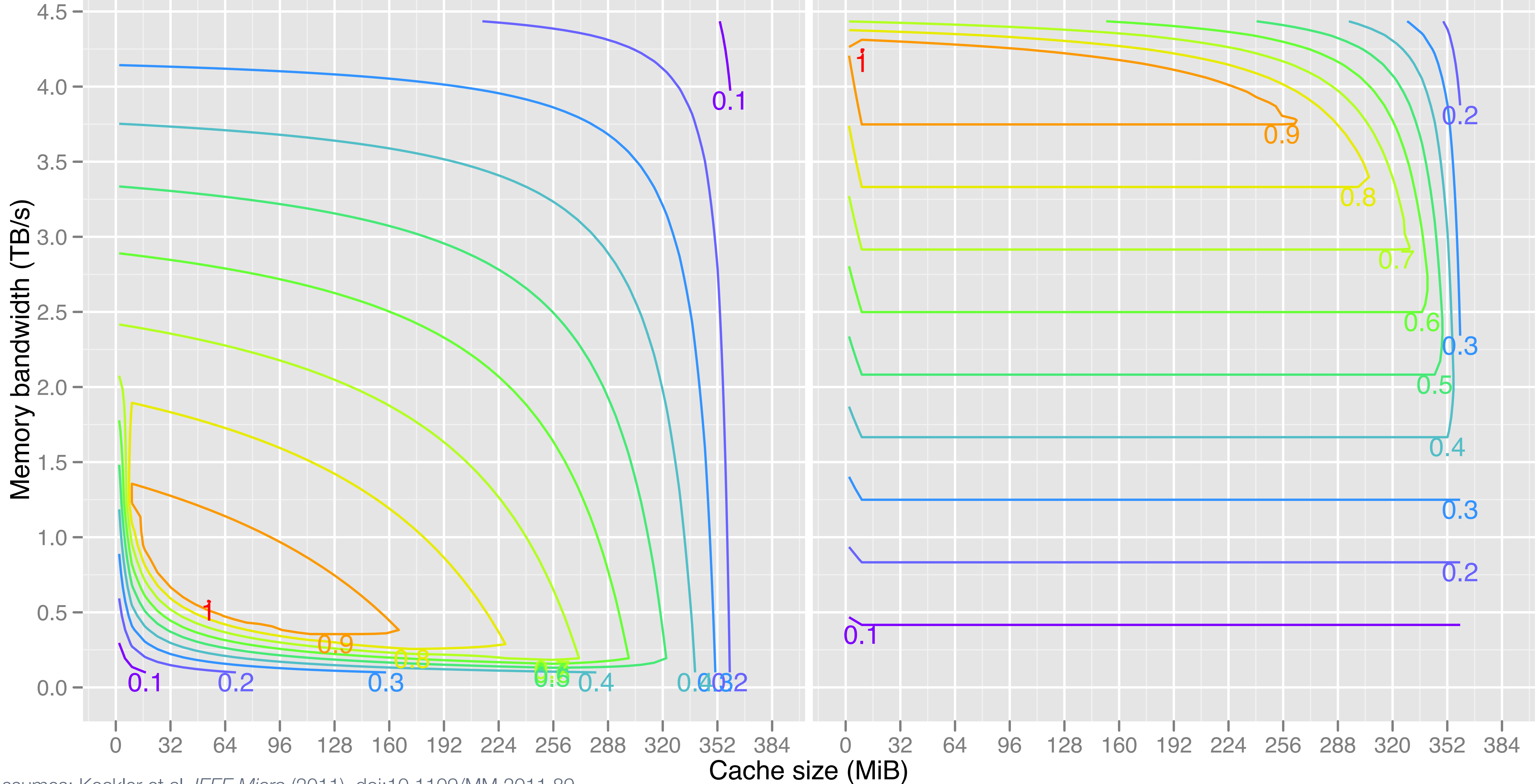


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

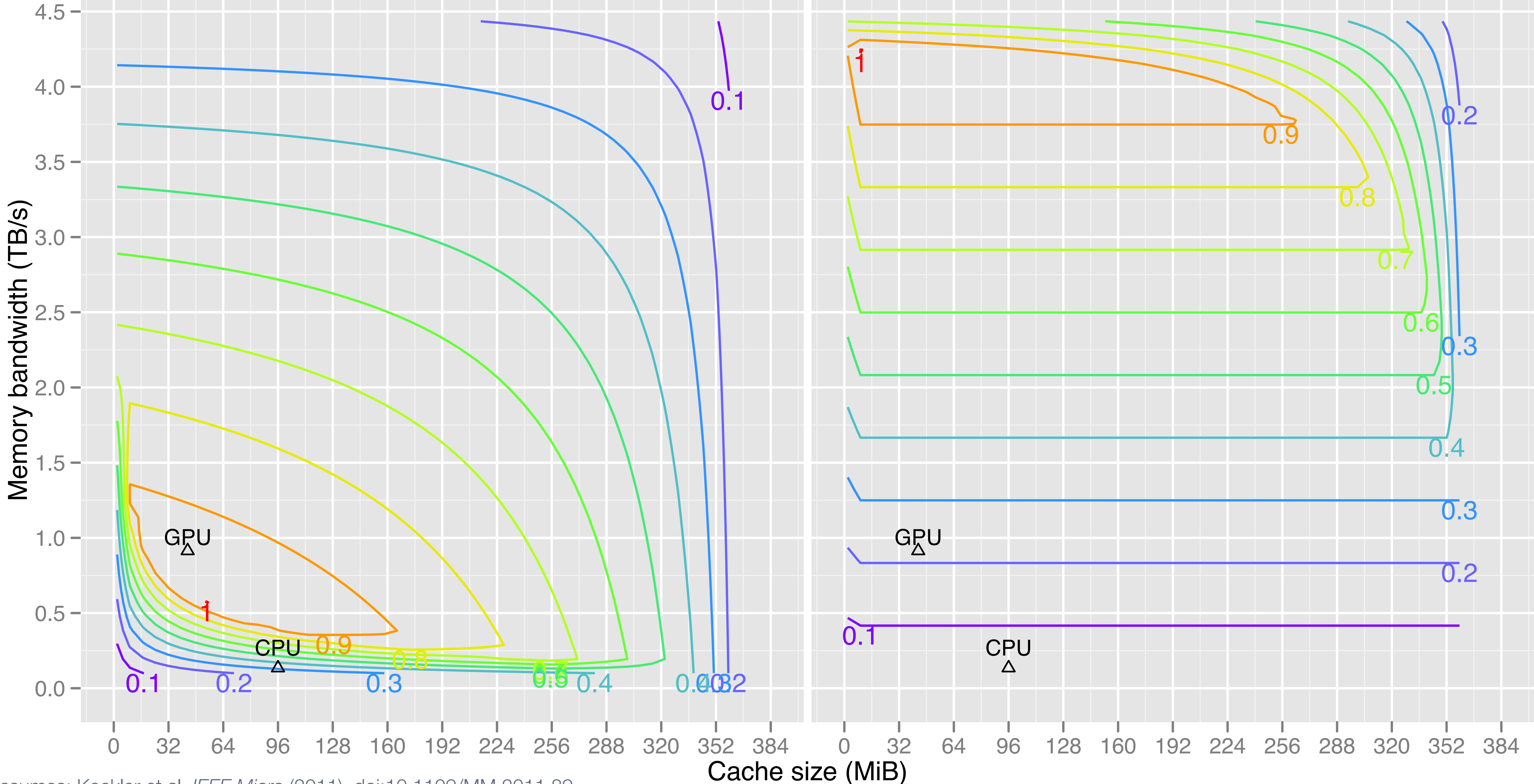


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

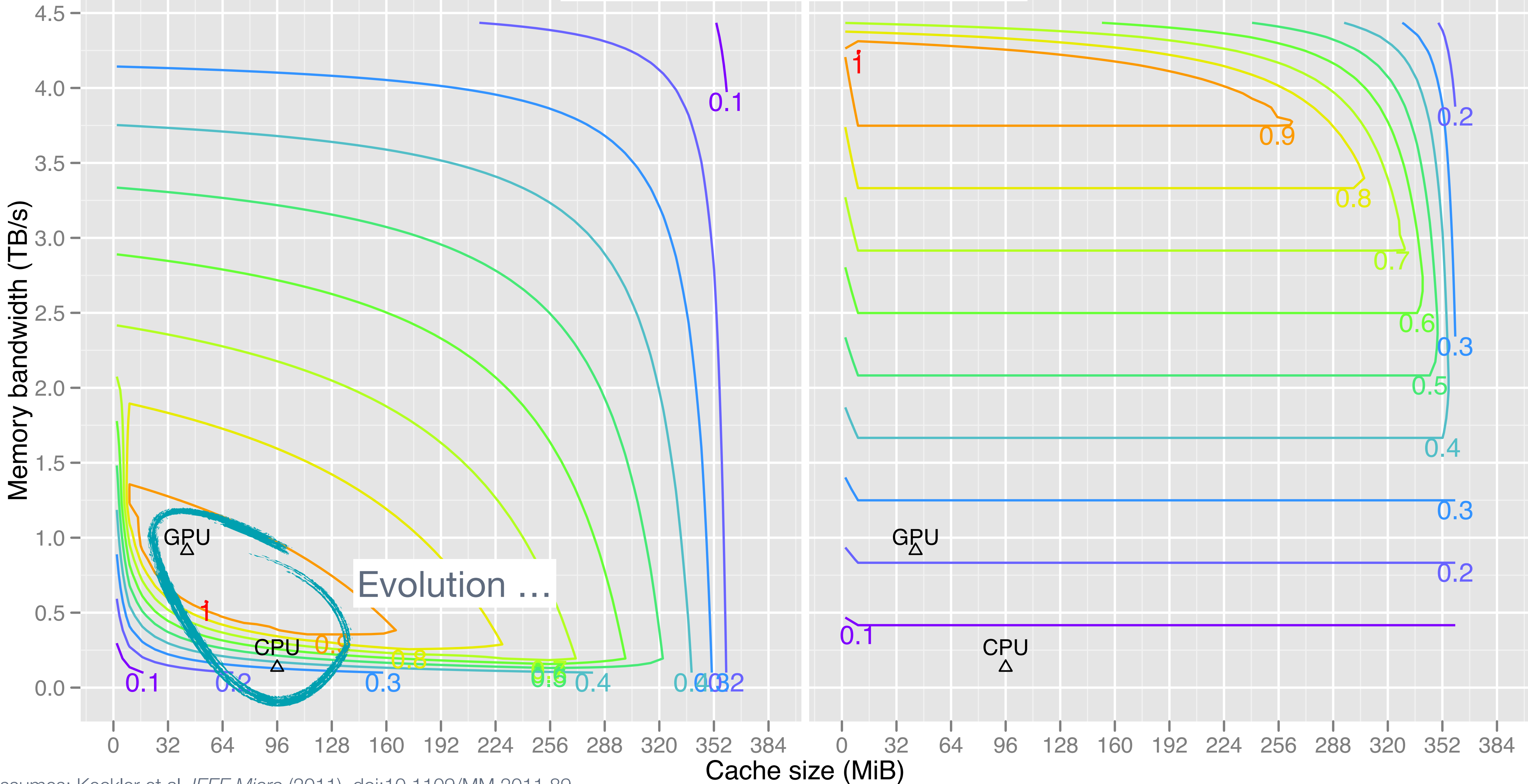


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

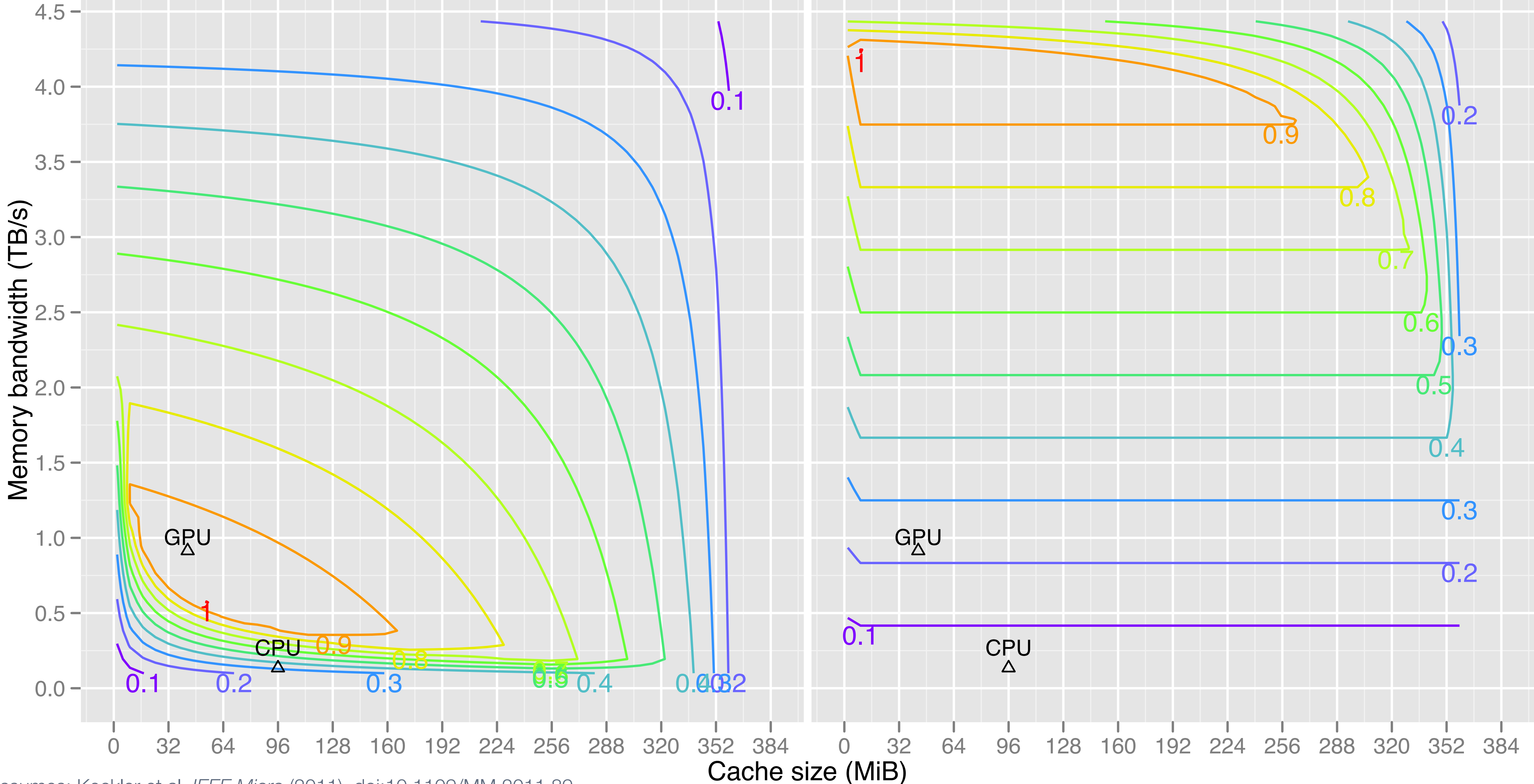


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT



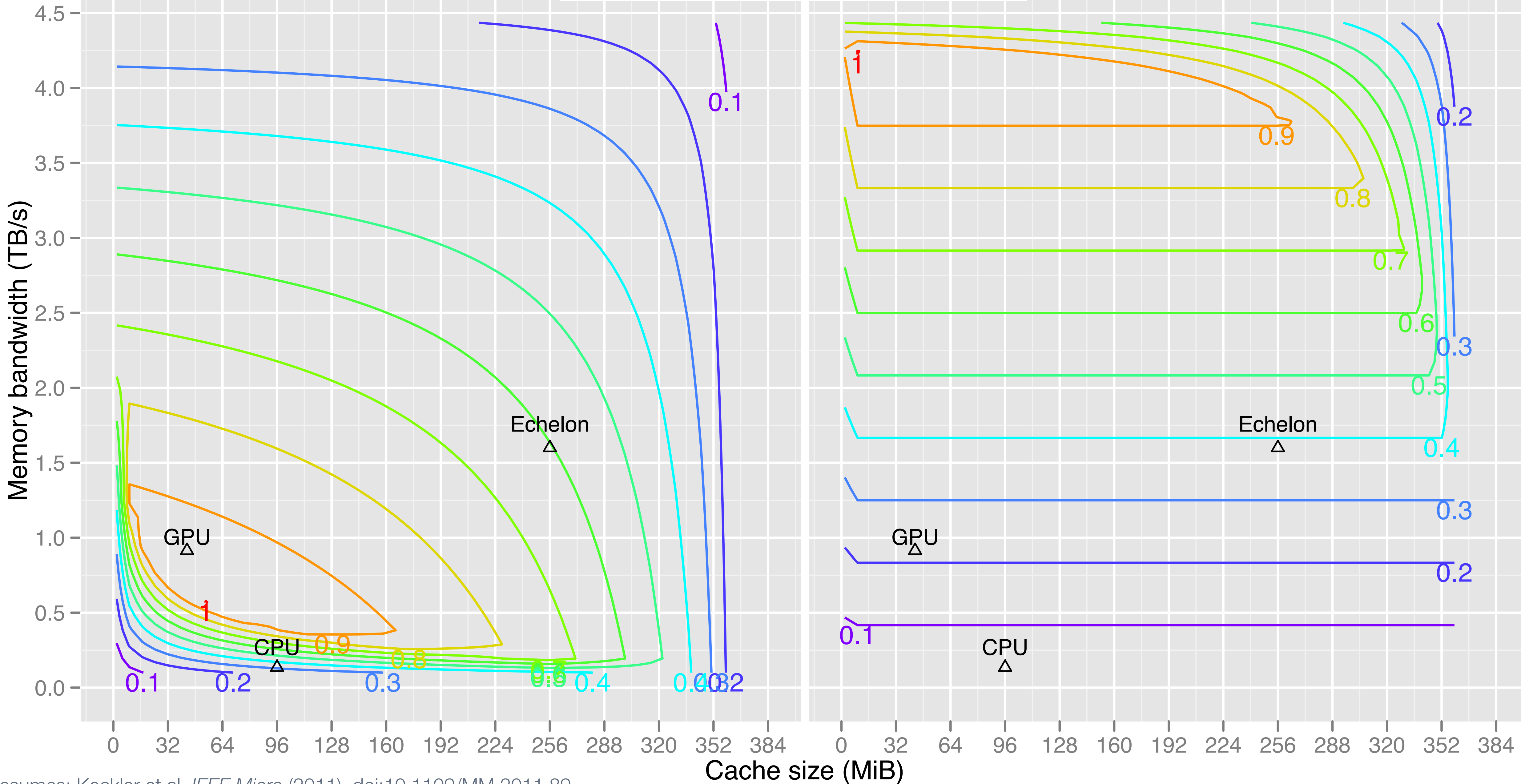
Assumes: Keckler et al. *IEEE Micro* (2011). doi:10.1109/MM.2011.89

Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

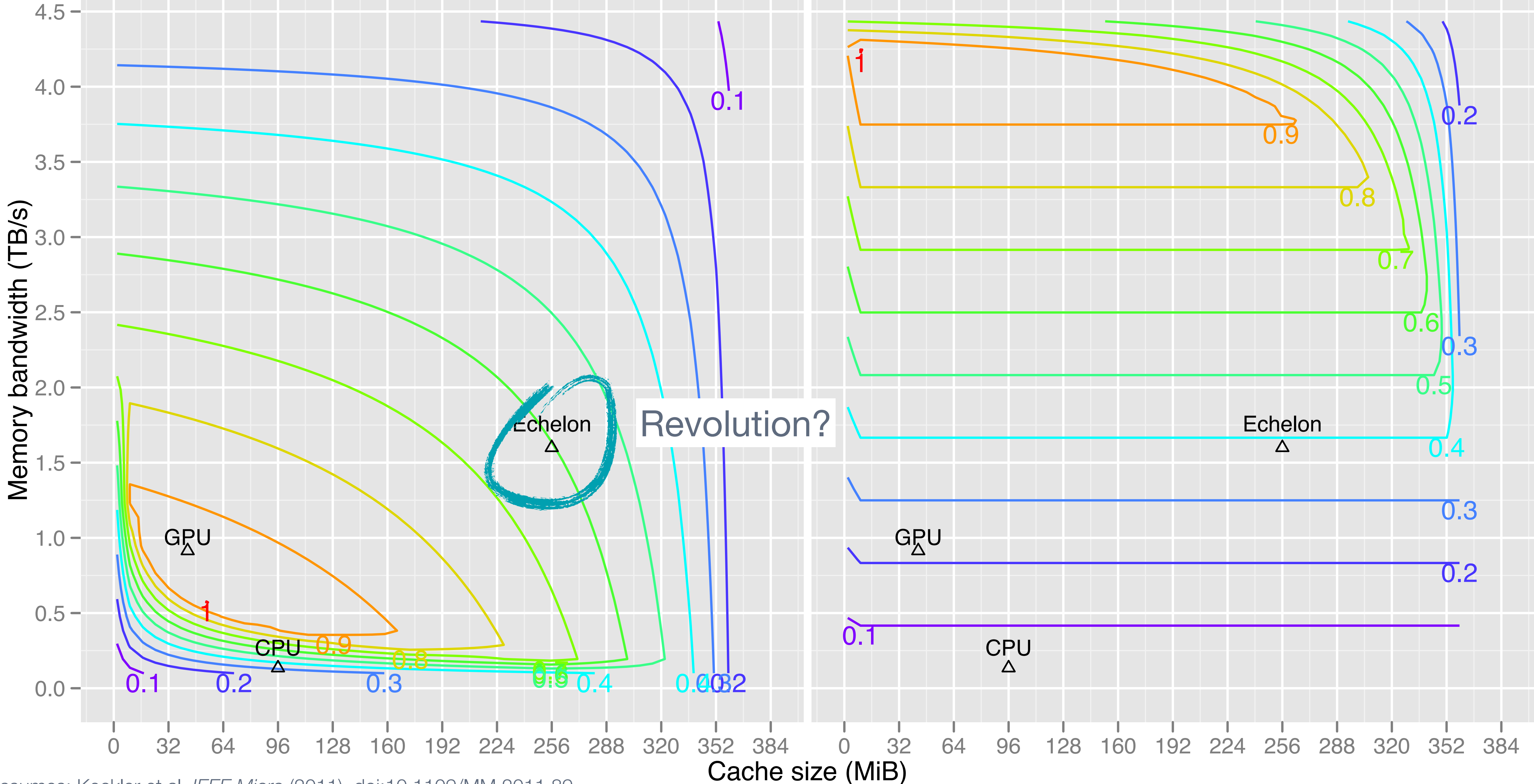


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

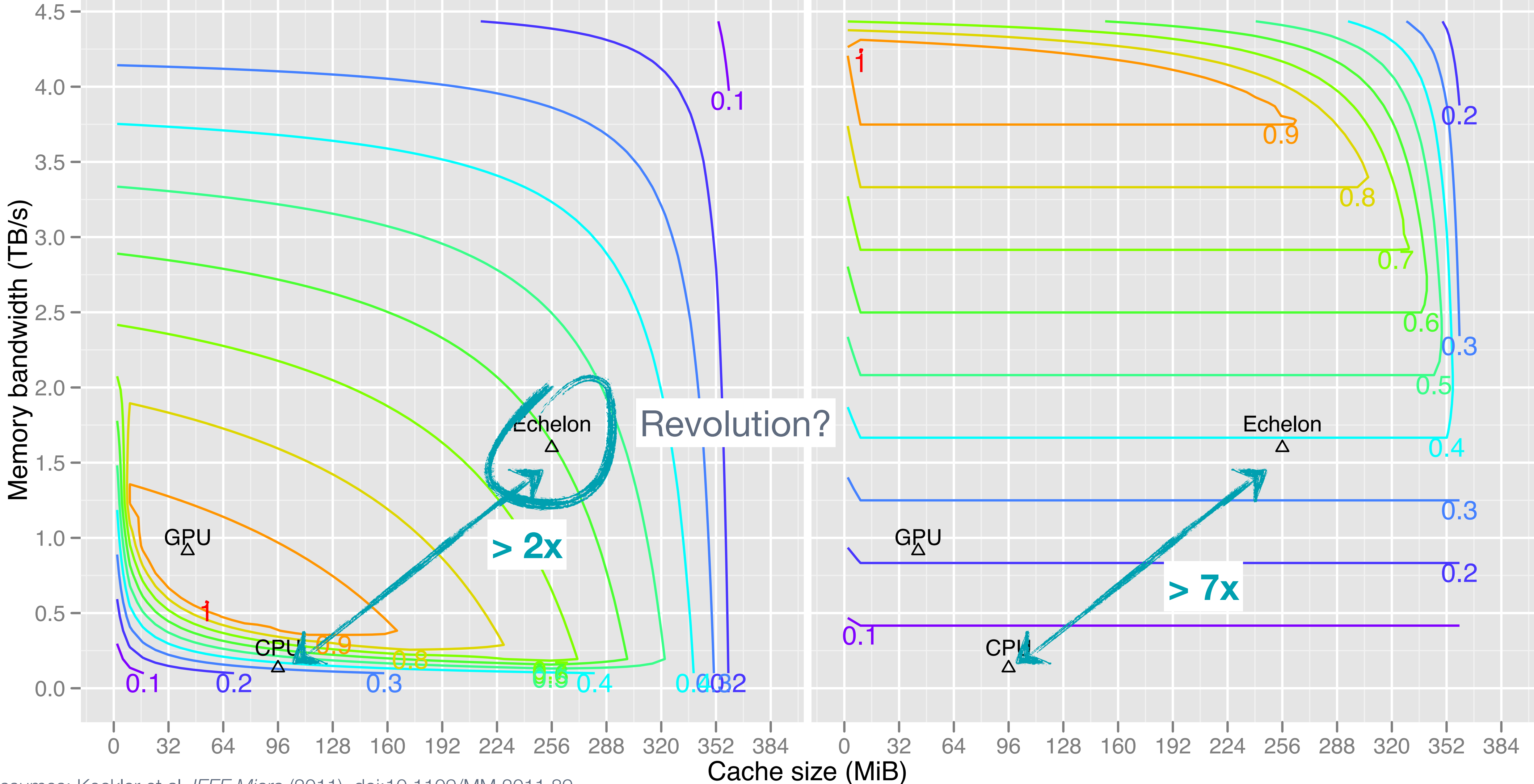


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

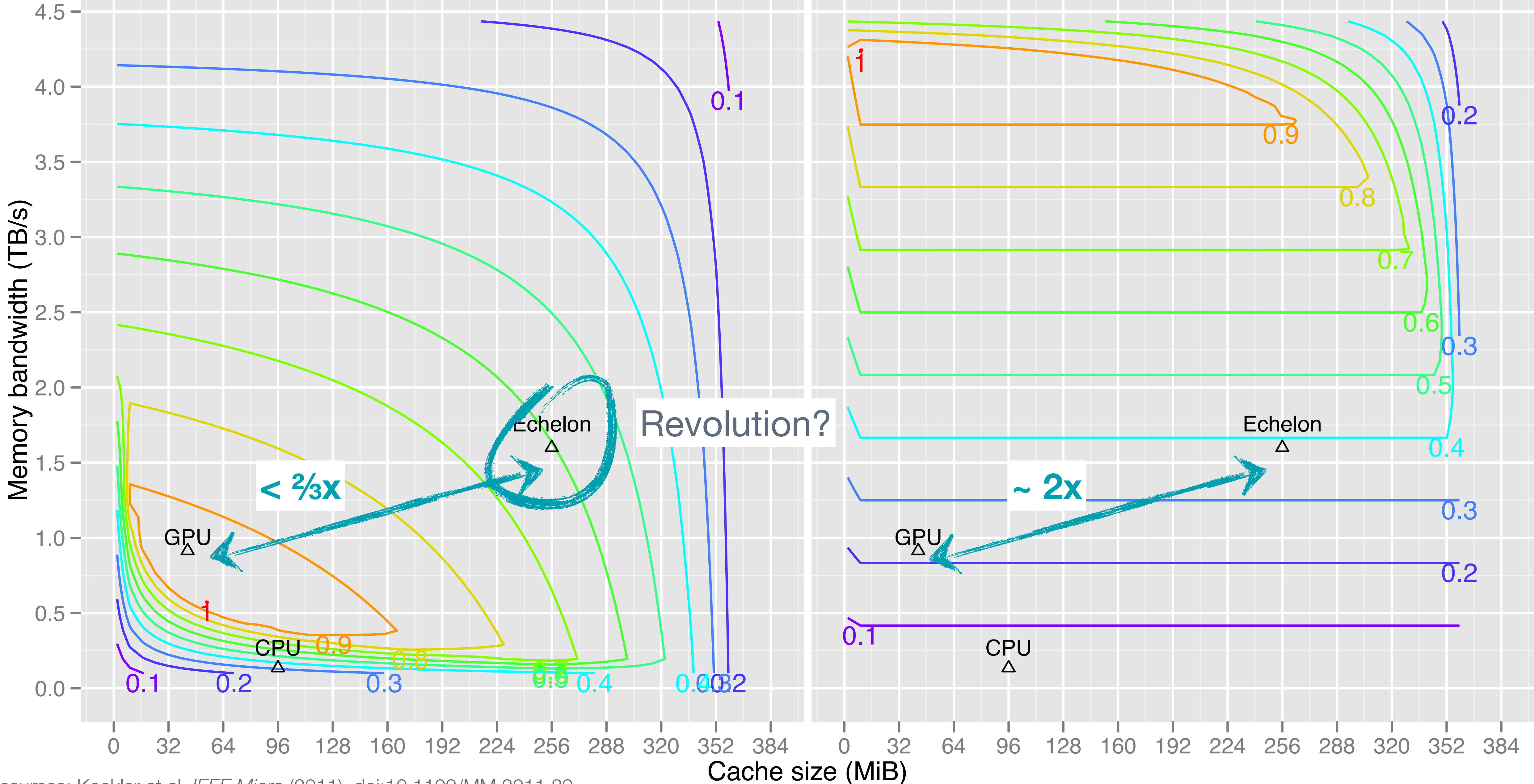


Sample Projection (~ 2017)

1 processor

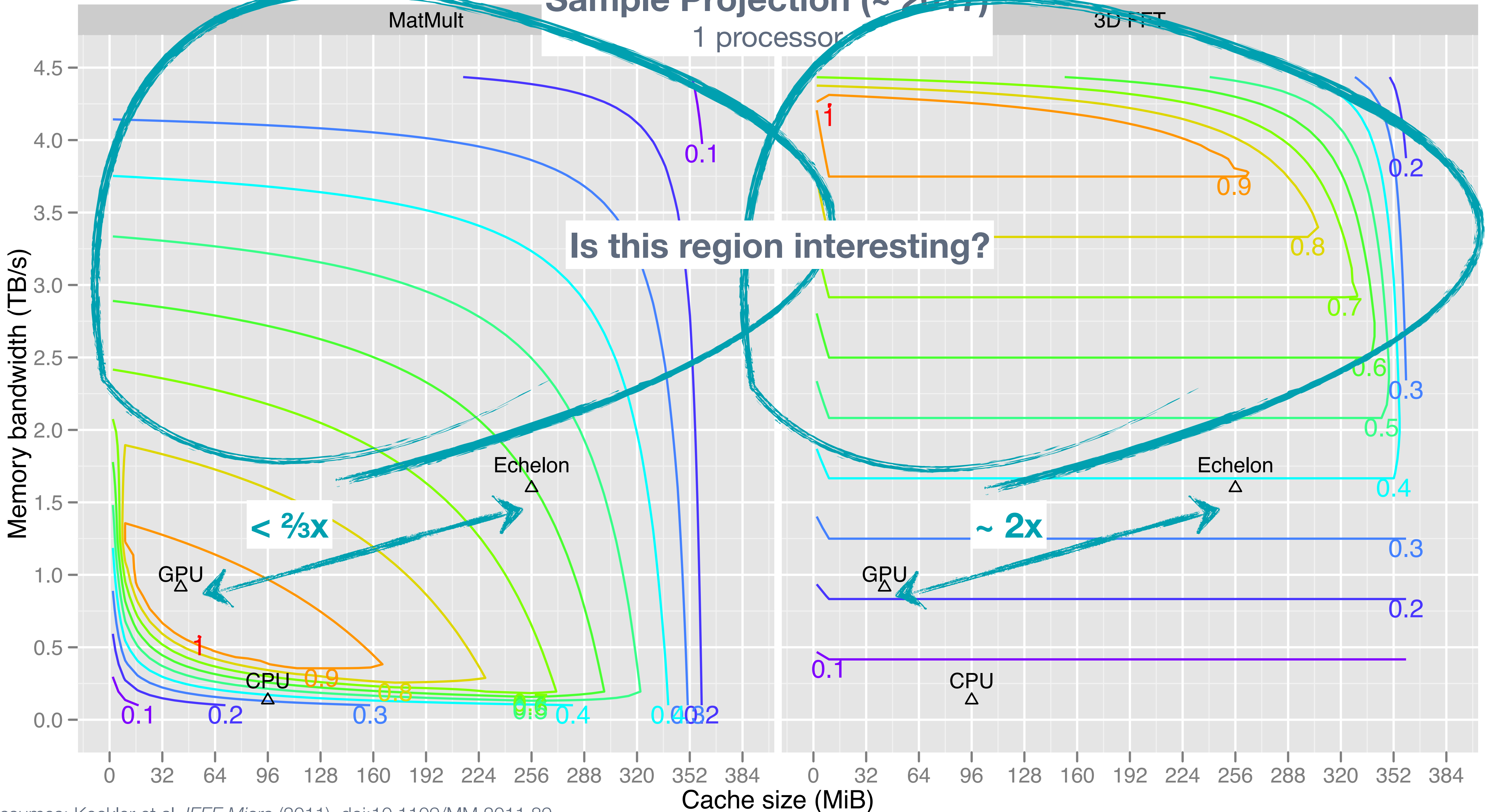
MatMult

3D FFT



Assumes: Keckler et al. *IEEE Micro* (2011). doi:10.1109/MM.2011.89

Sample Projection (~ 2017)



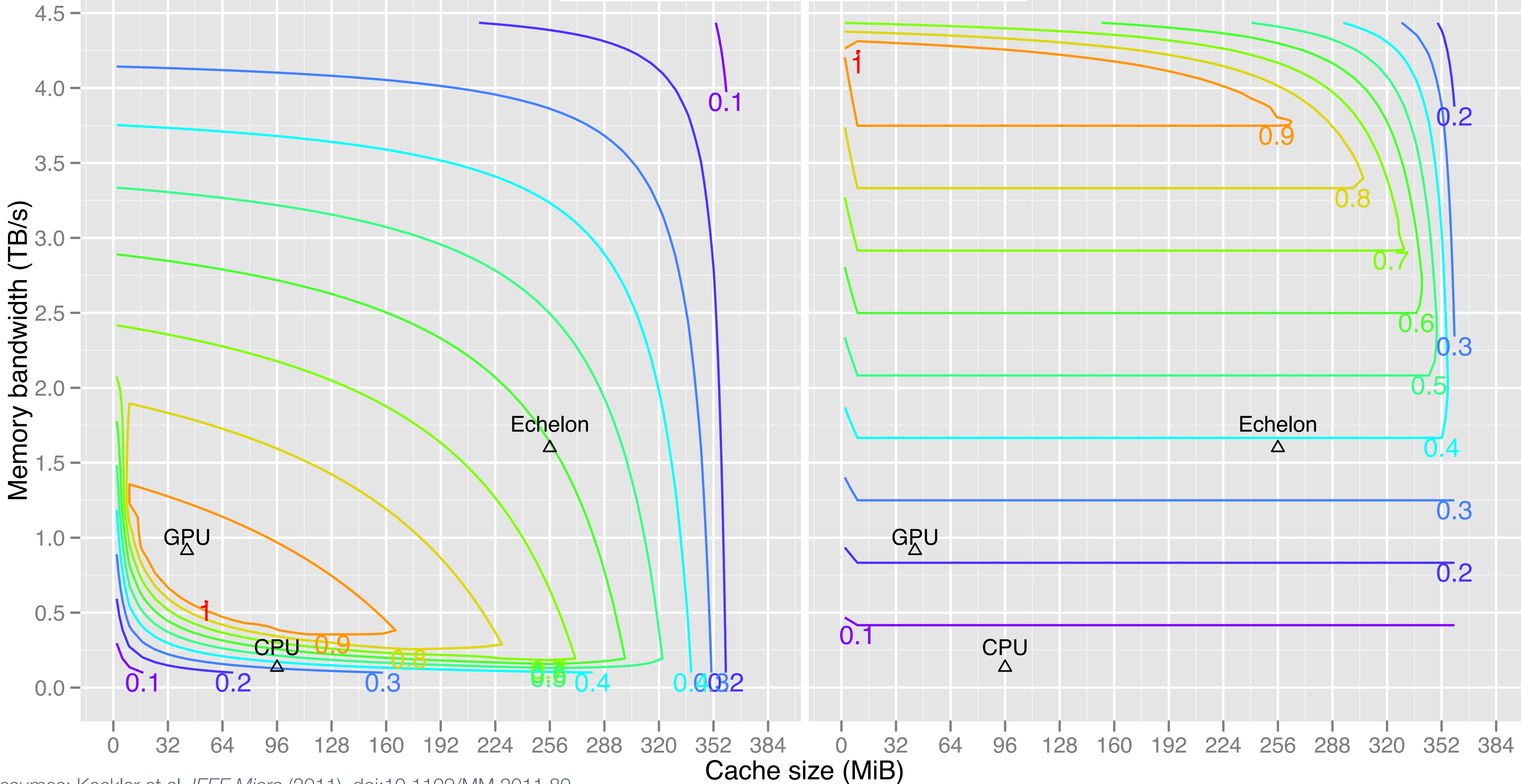
Assumes: Keckler et al. *IEEE Micro* (2011). doi:10.1109/MM.2011.89

Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

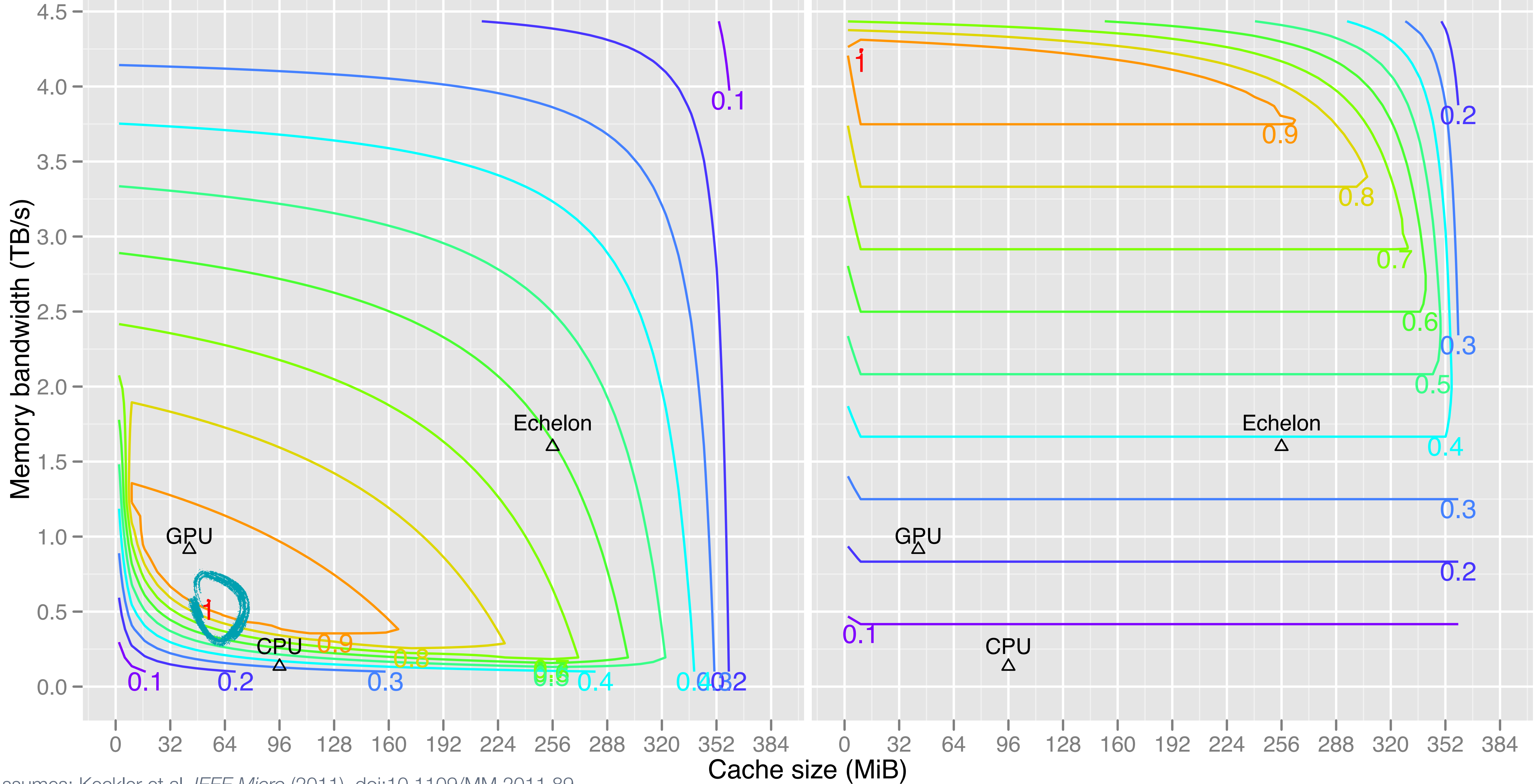


Sample Projection (~ 2017)

1 processor

MatMult

3D FFT

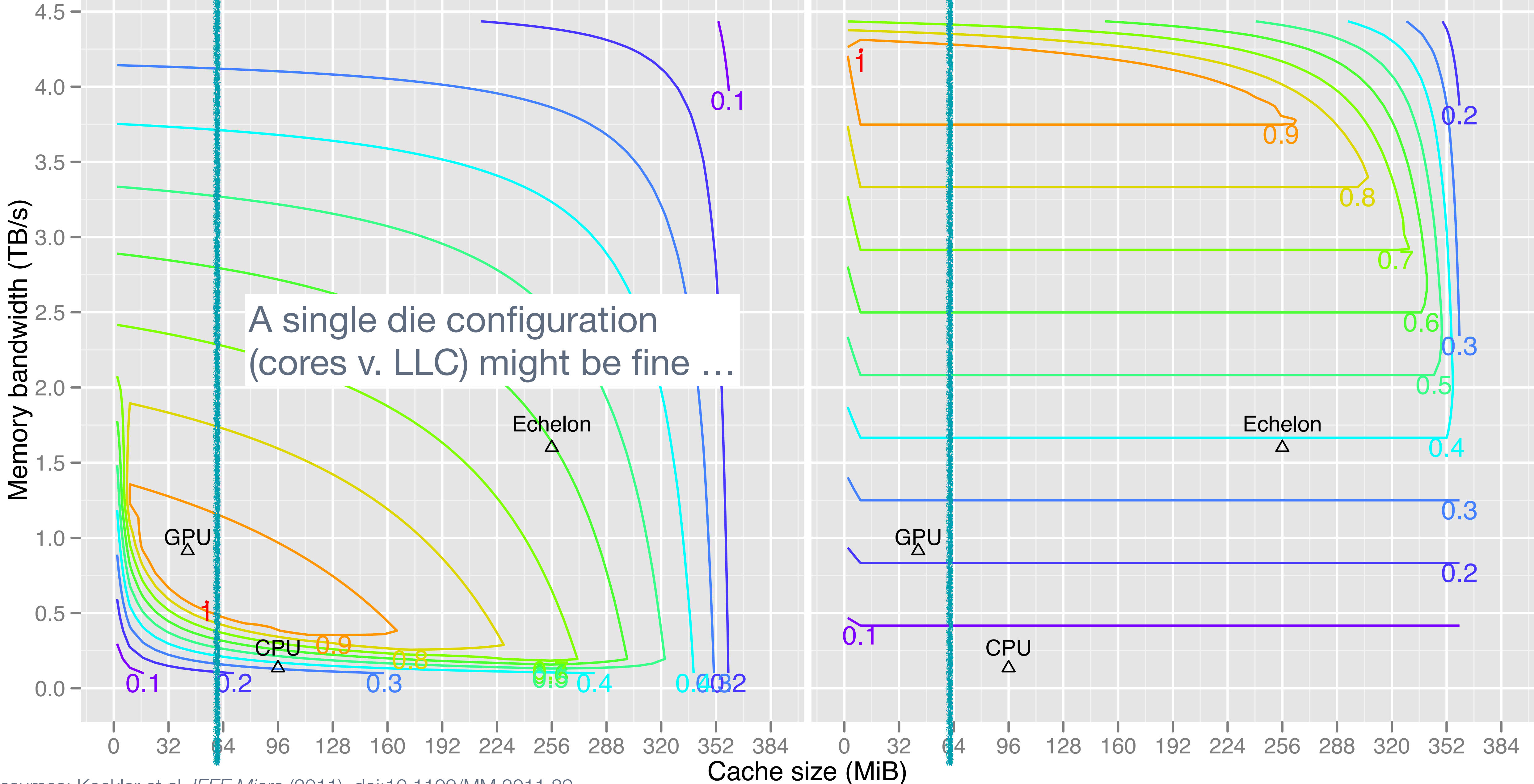


Sample Projection (~ 2017)

1 processor

MatMult

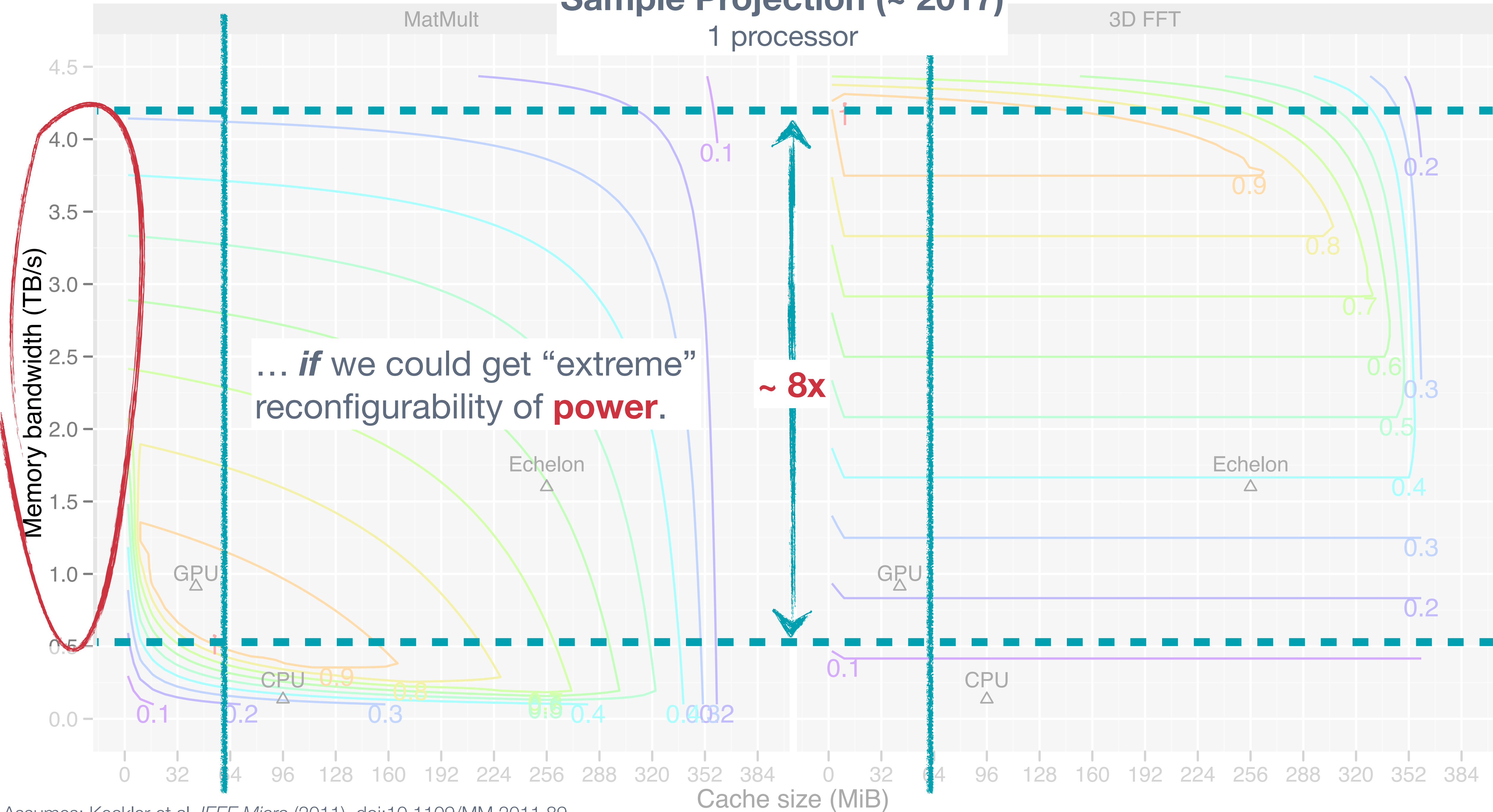
3D FFT



A single die configuration (cores v. LLC) might be fine ...

Sample Projection (~ 2017)

1 processor



Assumes: Keckler et al. *IEEE Micro* (2011). doi:10.1109/MM.2011.89

We can extend this analysis to a “full” system,
in which we consider, on-chip networks,
multiple nodes, network topology, ...

Solutions to the constrained optimization problem

Matrix multiply — 8 EF/s peak (1M “nodes”)

3D FFT — 230 PF/s peak (4k “nodes”)

Solutions to the constrained optimization problem

Matrix multiply — 8 EF/s peak (1M “nodes”)

Relative to notional Echelon:

~ **5x faster** for MM

~ **0.9x as fast** for 3D FFT

3D FFT — 230 PF/s peak (4k “nodes”)

Relative to notional Echelon:

~ **28x faster** for 3D FFTs

~ **0.14x as fast** for MM

Solutions to the constrained optimization problem

Matrix multiply — 8 EF/s peak (1M “nodes”)

3D FFT — 230 PF/s peak (4k “nodes”)

Relative to notional Echelon:

~ **5x faster** for MM

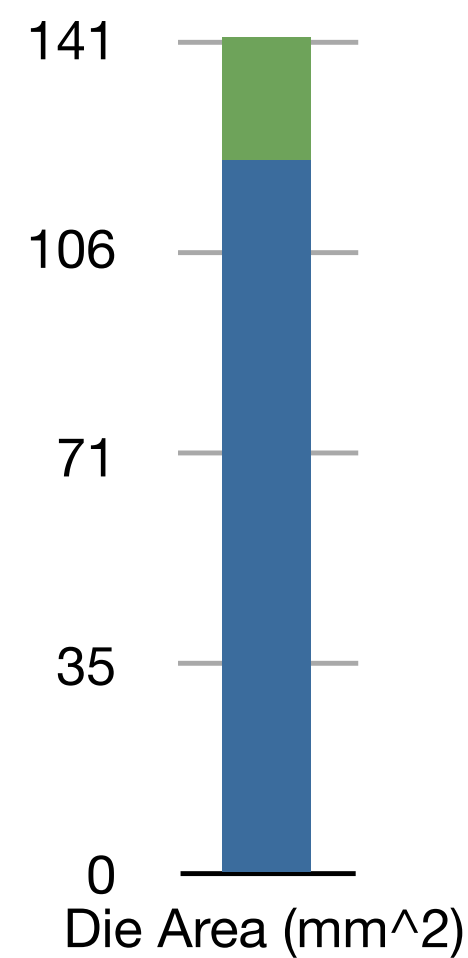
~ **0.9x as fast** for 3D FFT

Relative to notional Echelon:

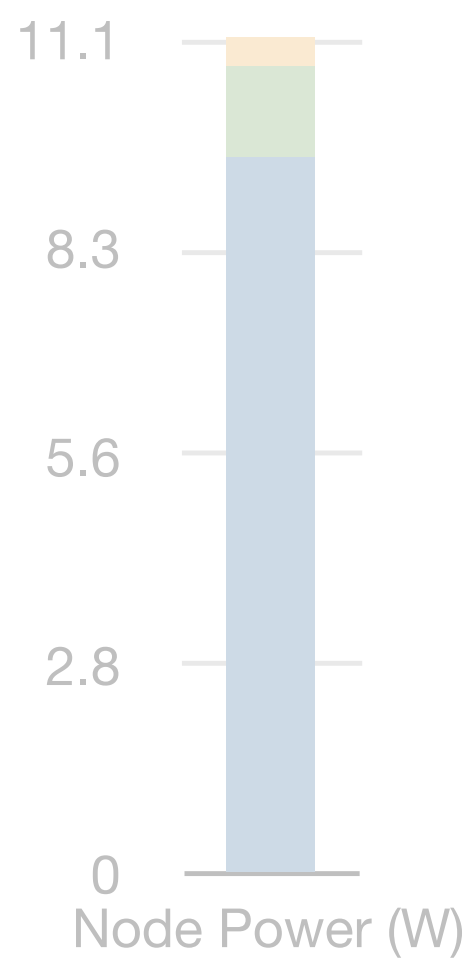
~ **28x faster** for 3D FFTs

~ **0.14x as fast** for MM

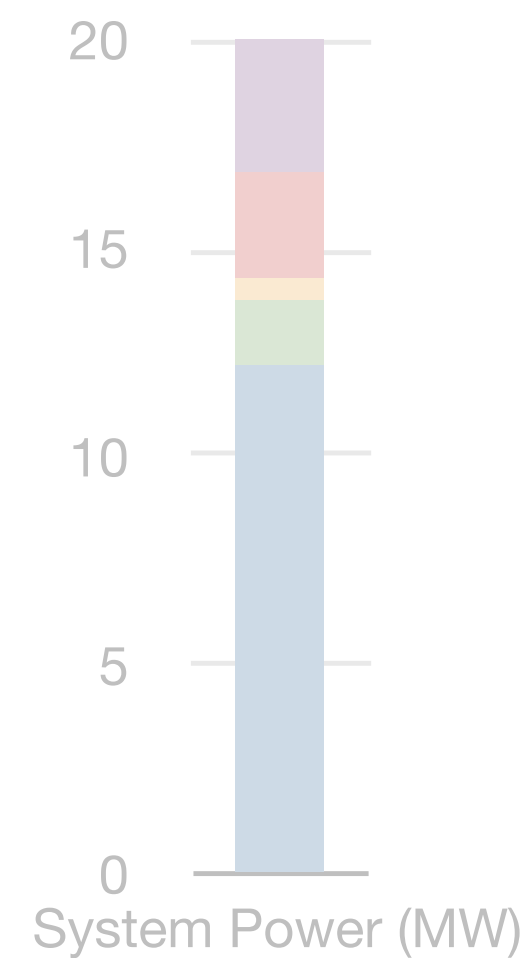
Transistor Budget



Node Power Budget



System Power

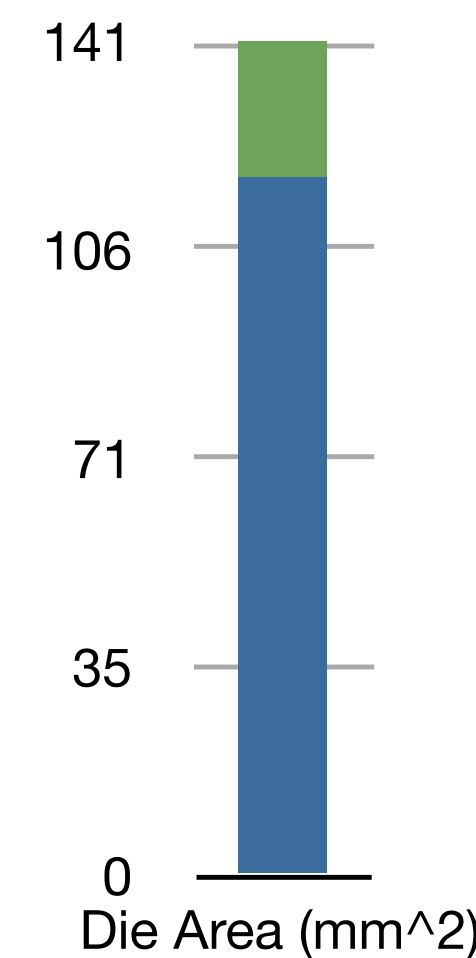


Cache Area
Core Area

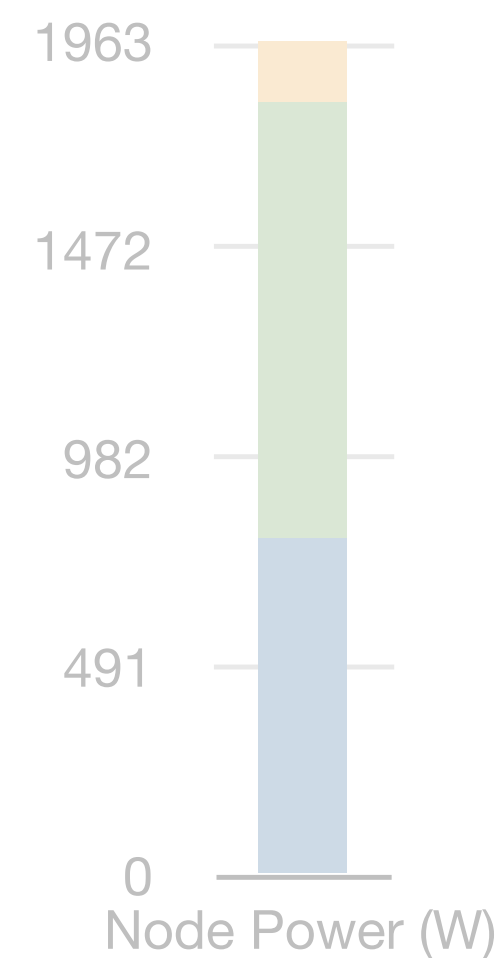
Noc Bandwidth
Mem Bandwidth
Cores

Network
Node Overhead
On-Chip Network
Memory Bandwidth
Computation

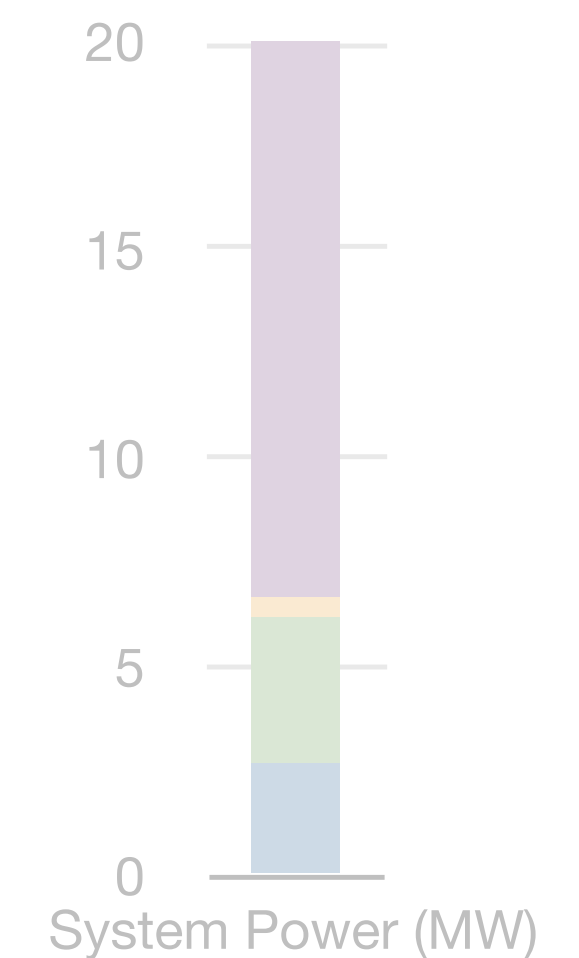
Transistor Budget



Node Power Budget



System Power



Cache Area
Core Area

Noc Bandwidth
Mem Bandwidth
Cores

Network
Node Overhead
On-Chip Network
Memory Bandwidth
Computation

Plausibility issue: Unconstrained power density

Solutions to the constrained optimization problem

Matrix multiply — 8 EF/s peak (1M “nodes”)

3D FFT — 230 PF/s peak (4k “nodes”)

Relative to notional Echelon:

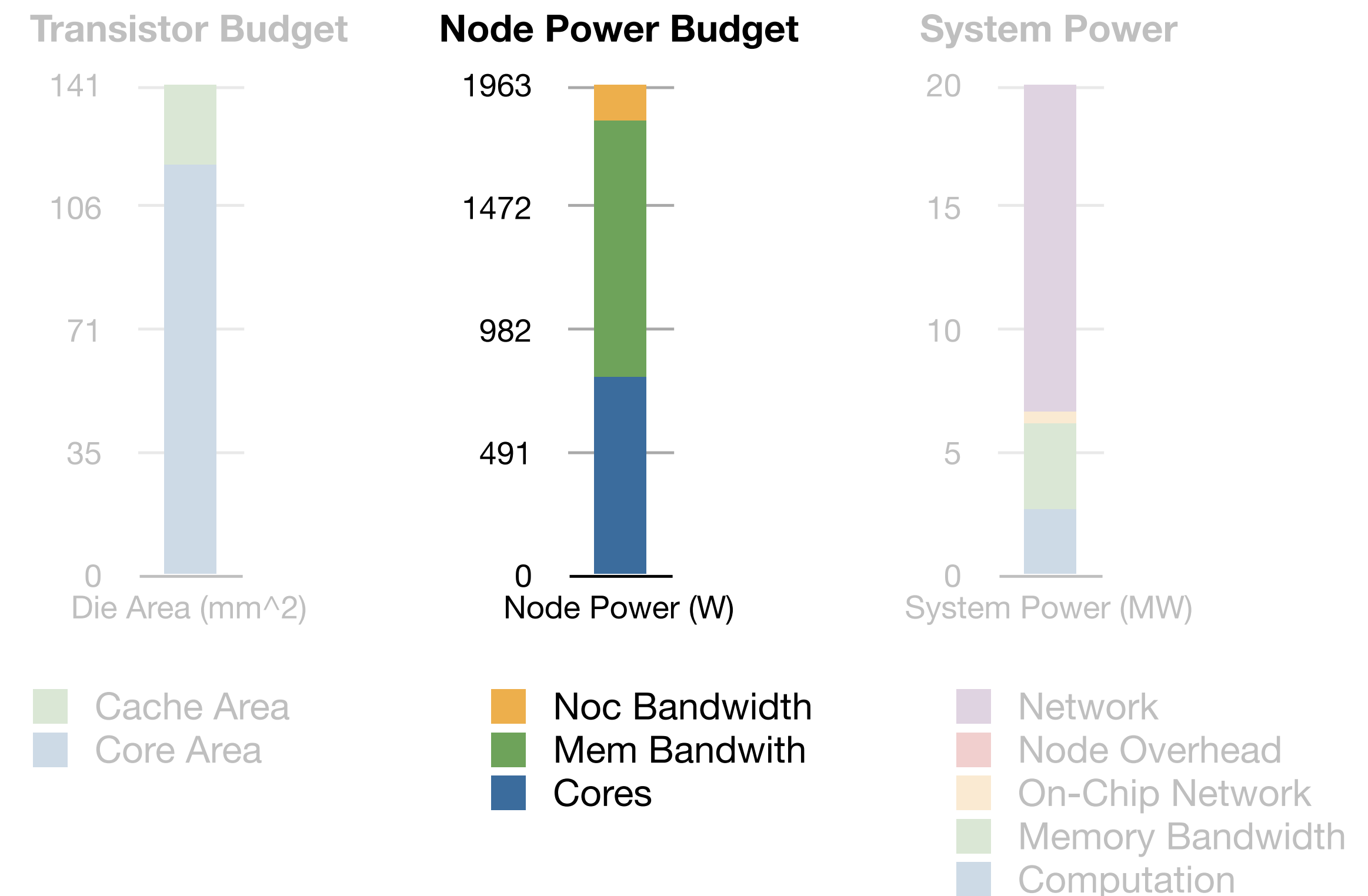
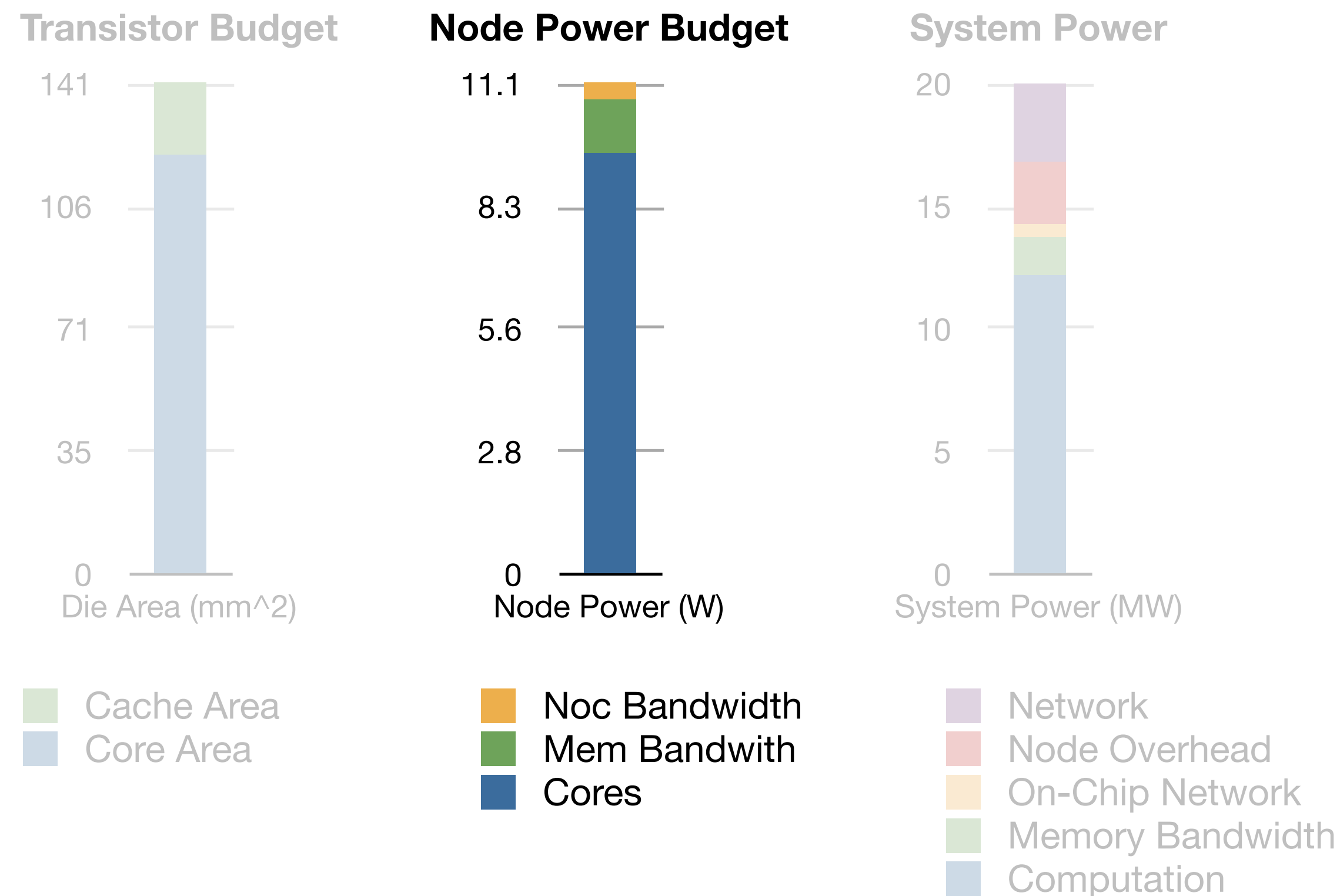
~ **5x faster** for MM

~ **0.9x as fast** for 3D FFT

Relative to notional Echelon:

~ **28x faster** for 3D FFTs

~ **0.14x as fast** for MM



Plausibility issue: Unconstrained power density

Solutions to the constrained optimization problem

Matrix multiply — 8 EF/s peak (1M “nodes”)

3D FFT — 230 PF/s peak (4k “nodes”)

Relative to notional Echelon:

~ **5x faster** for MM

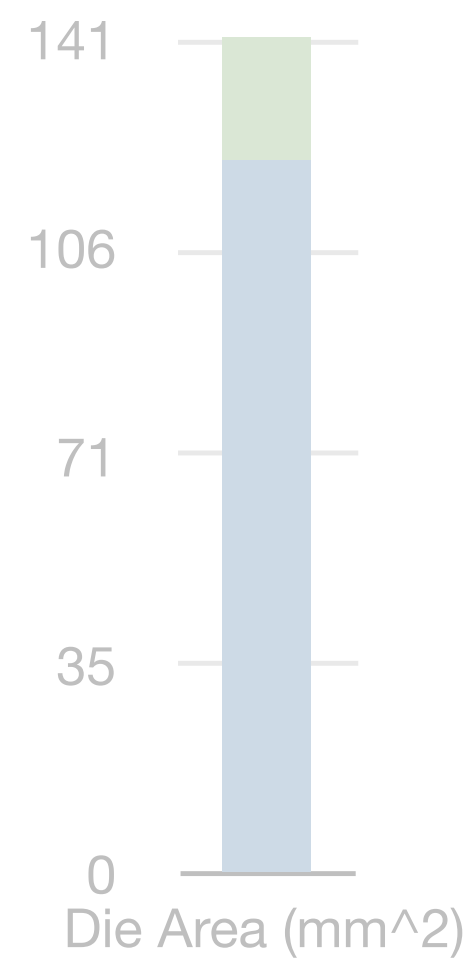
~ **0.9x as fast** for 3D FFT

Relative to notional Echelon:

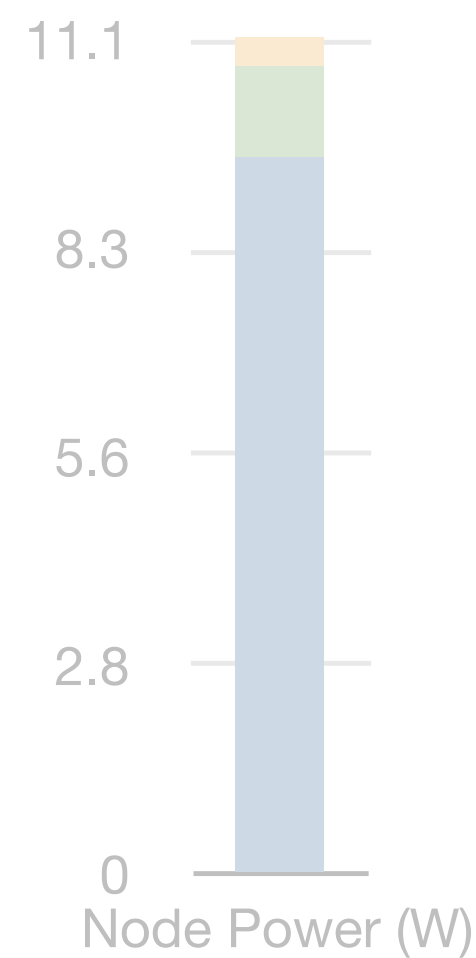
~ **28x faster** for 3D FFTs

~ **0.14x as fast** for MM

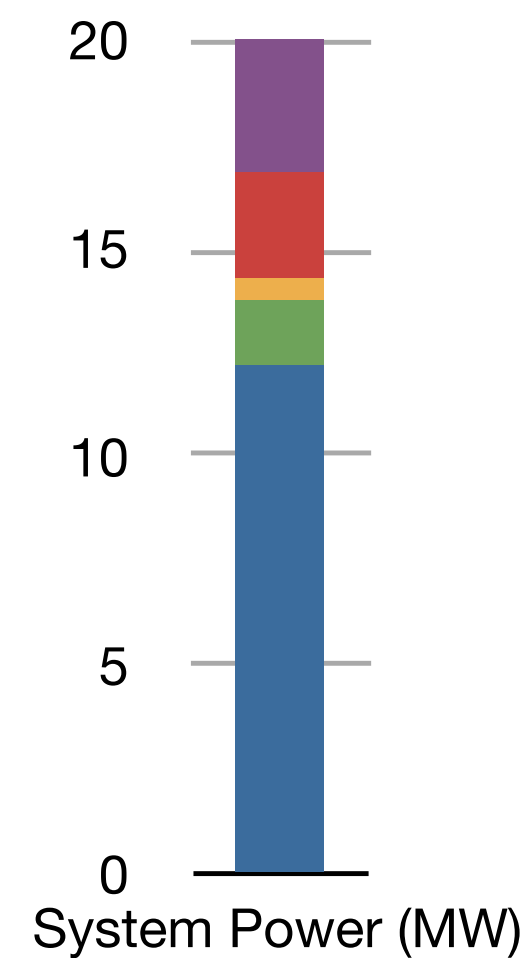
Transistor Budget



Node Power Budget



System Power

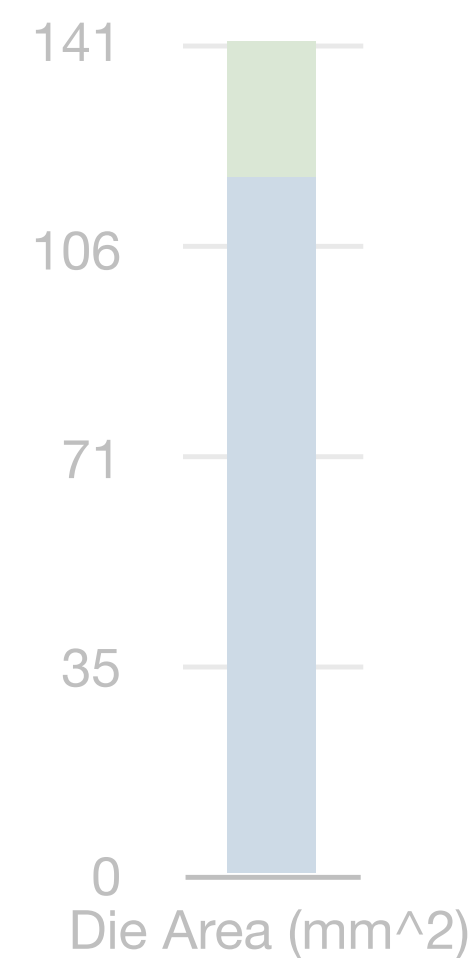


Cache Area
Core Area

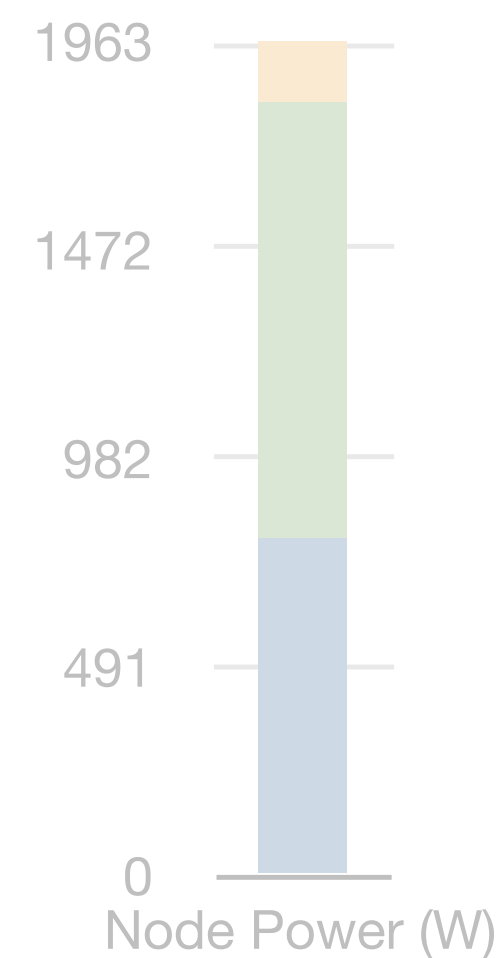
Noc Bandwidth
Mem Bandwidth
Cores

Network
Node Overhead
On-Chip Network
Memory Bandwidth
Computation

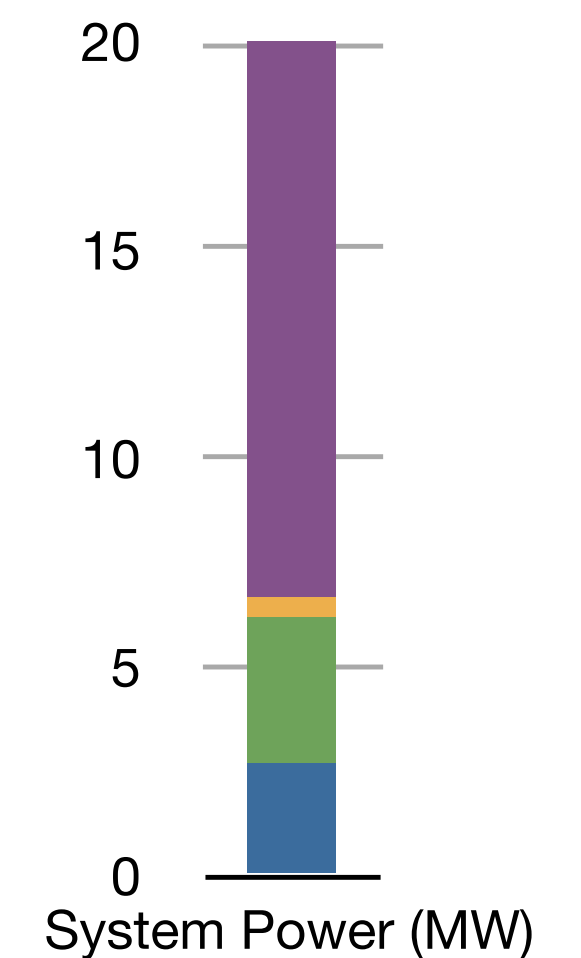
Transistor Budget



Node Power Budget



System Power



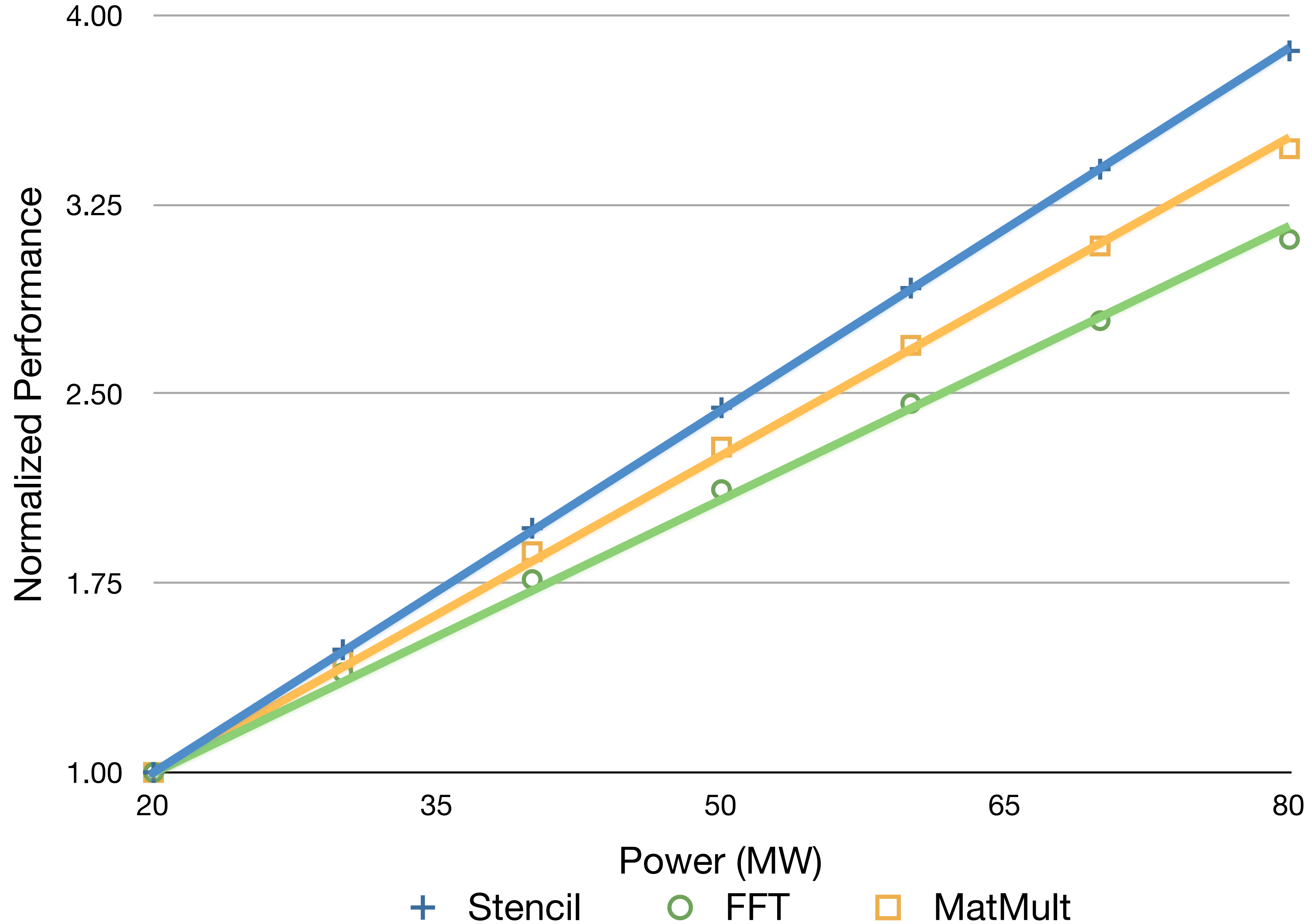
Cache Area
Core Area

Noc Bandwidth
Mem Bandwidth
Cores

Network
Node Overhead
On-Chip Network
Memory Bandwidth
Computation

Plausibility issue: Unconstrained power density

Performance as a Function of System Power



Summary:

Are we looking in the right place in the space of “exascale-able” designs?

Can starting from algorithmic first-principles and minimizing time, with power and area as the principal constraints, identify other regions of interest?

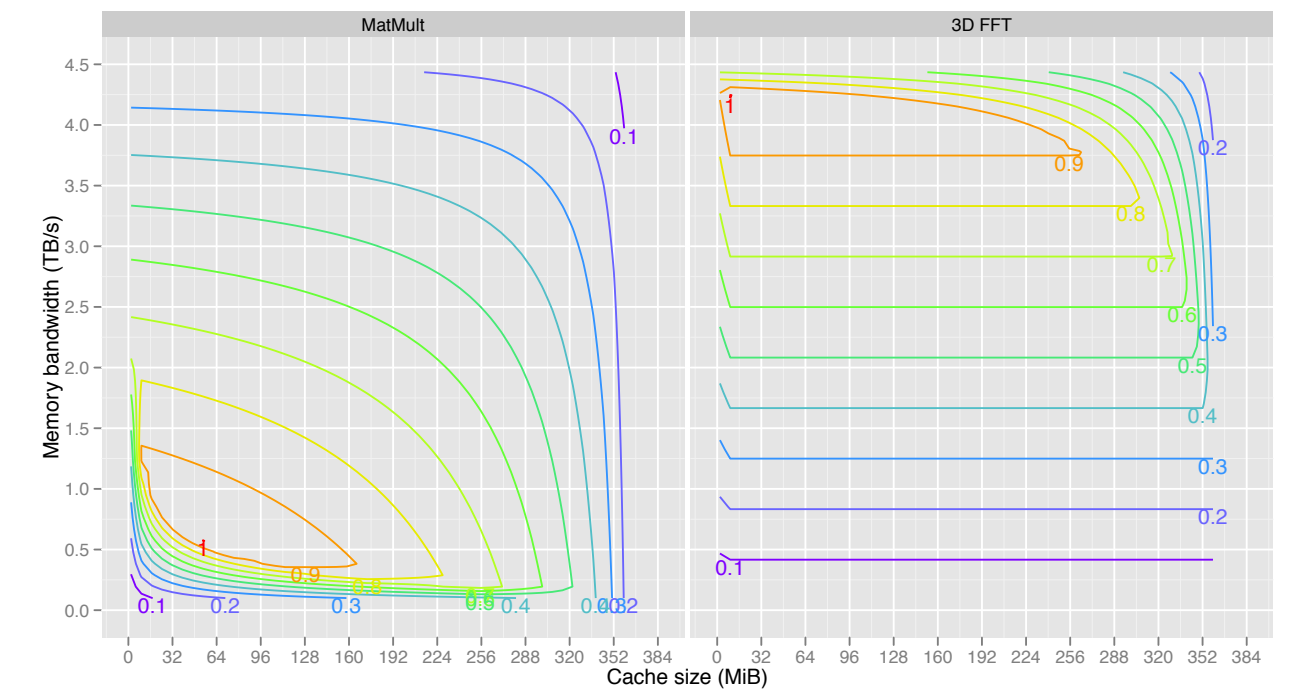
Questions:

What workloads and algorithms matter most?

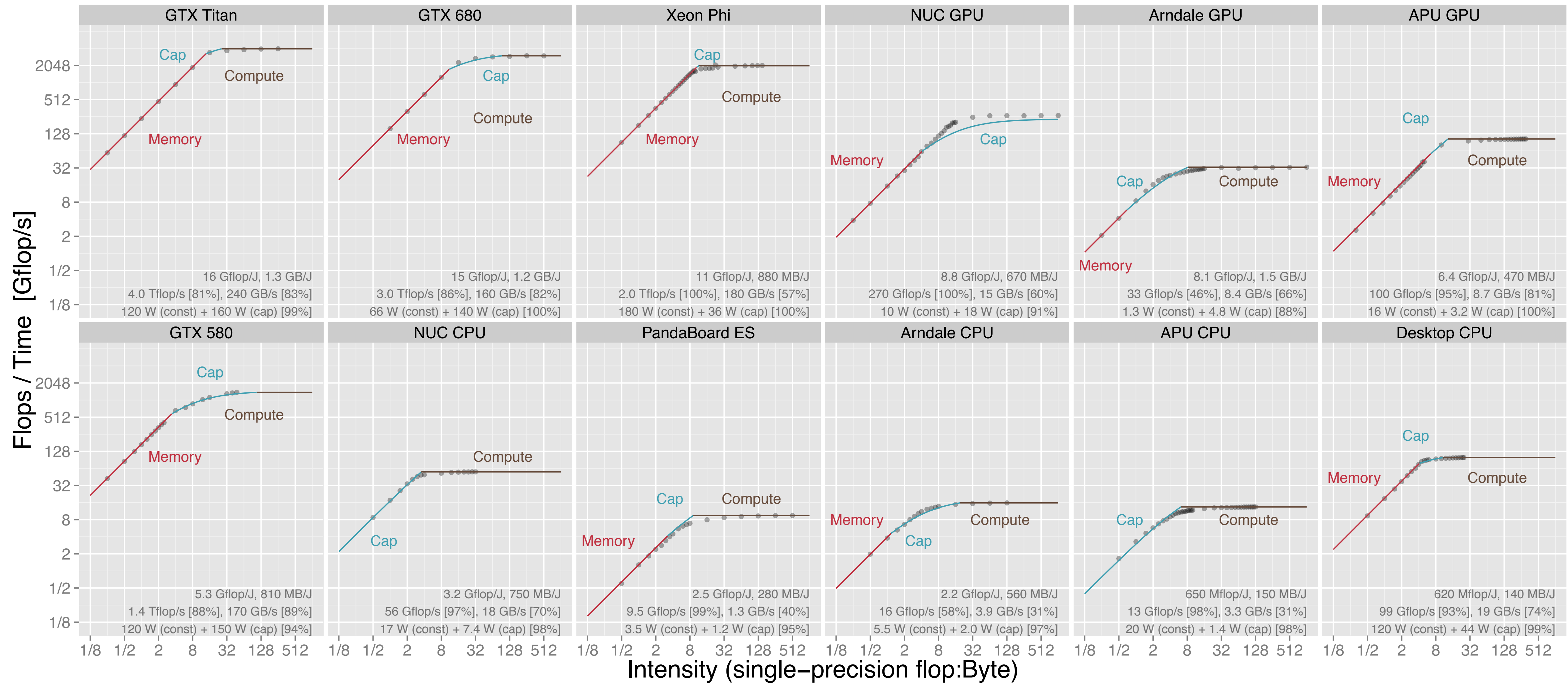
What algorithmic trade-offs can we analyze?

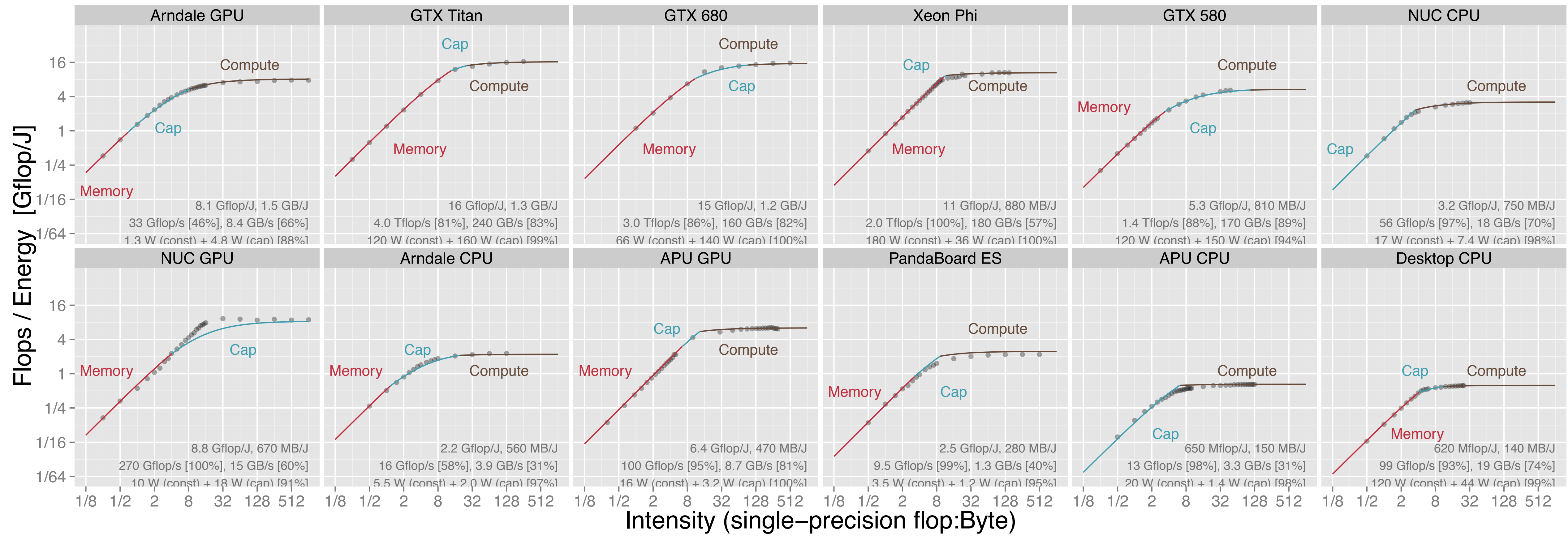
What architectural designs are feasible?

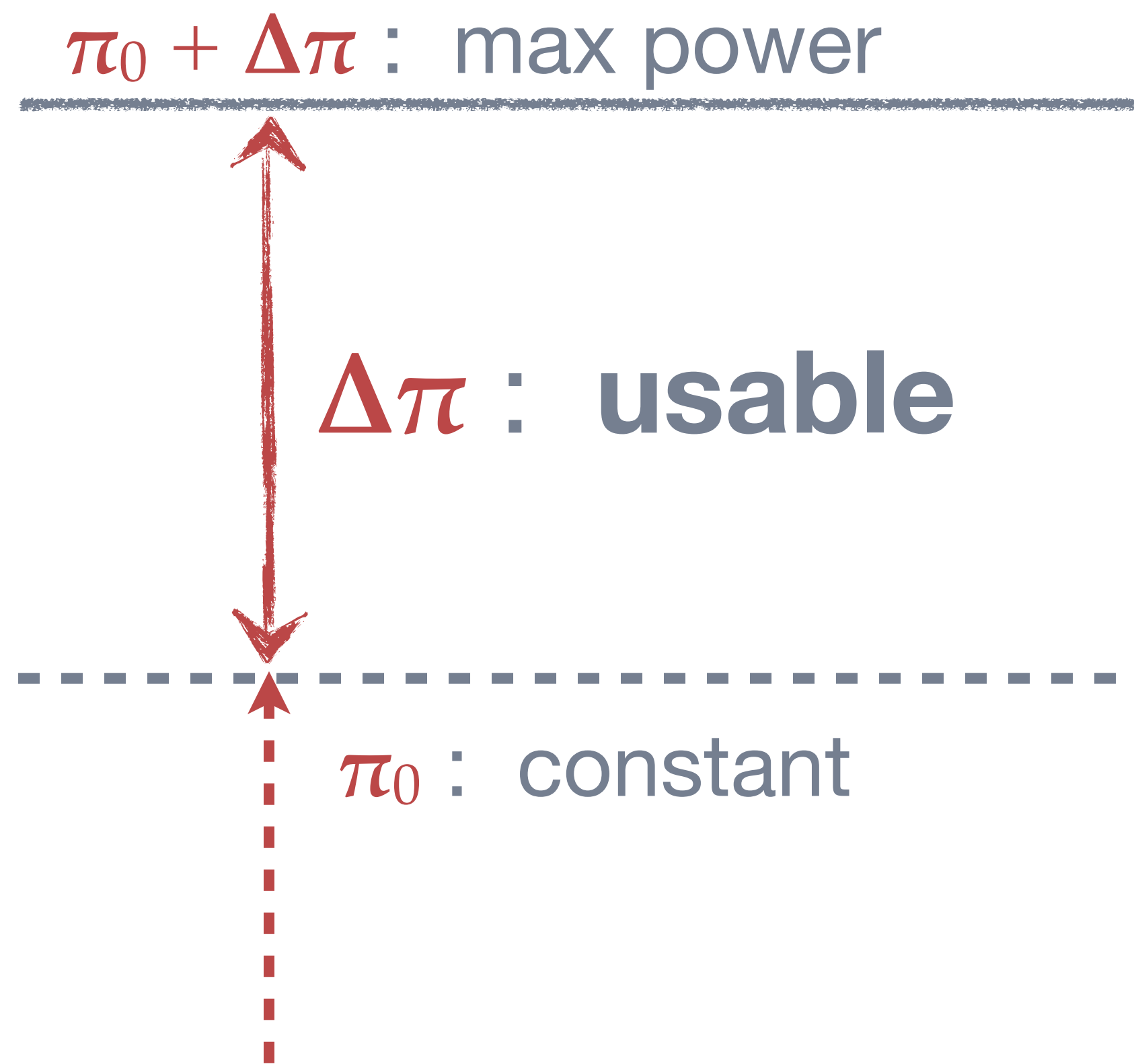
Can we estimate the impact of disruptive technologies?



Other junk







Peak power per flop (or mop)

$$\pi_{\text{flop}} \equiv \frac{\epsilon_{\text{flop}}}{\tau_{\text{flop}}}$$

$$\pi_{\text{mem}} \equiv \frac{\epsilon_{\text{mem}}}{\tau_{\text{mem}}}$$

↓

$$\Delta\pi \geq \pi_{\text{flop}} + \pi_{\text{mem}}$$

$$= \pi_{\text{flop}} \left(1 + \frac{B_\epsilon}{B_\tau} \right)$$

Caps imply throttling!

What power cap will obviate throttling?
 The *balance gap* dictates a sufficient (algorithm-independent) condition.



NVIDIA

LANL

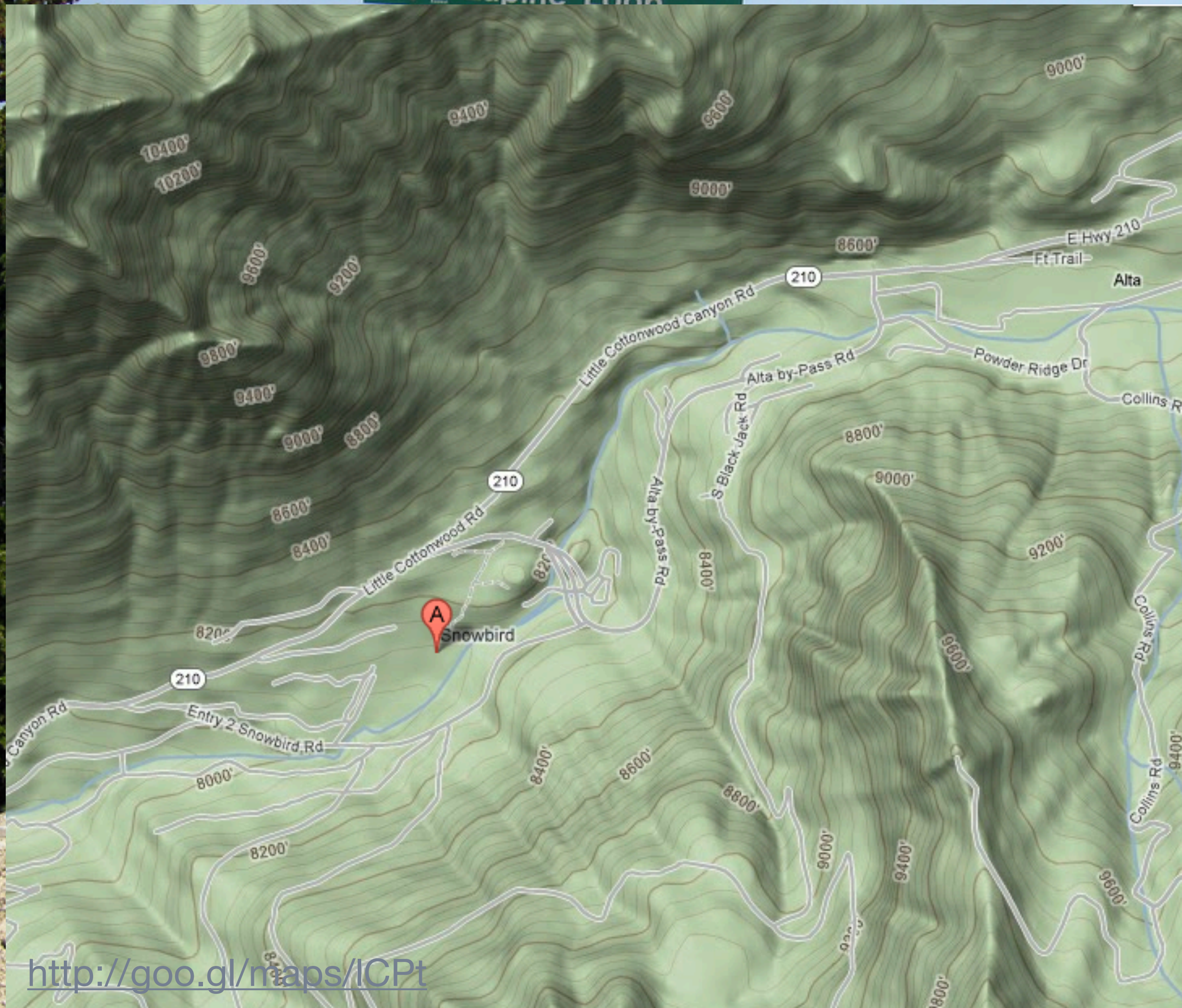


NVIDIA

LANL

Cray

← Lupine Loop



<http://goo.gl/maps/ICPt>