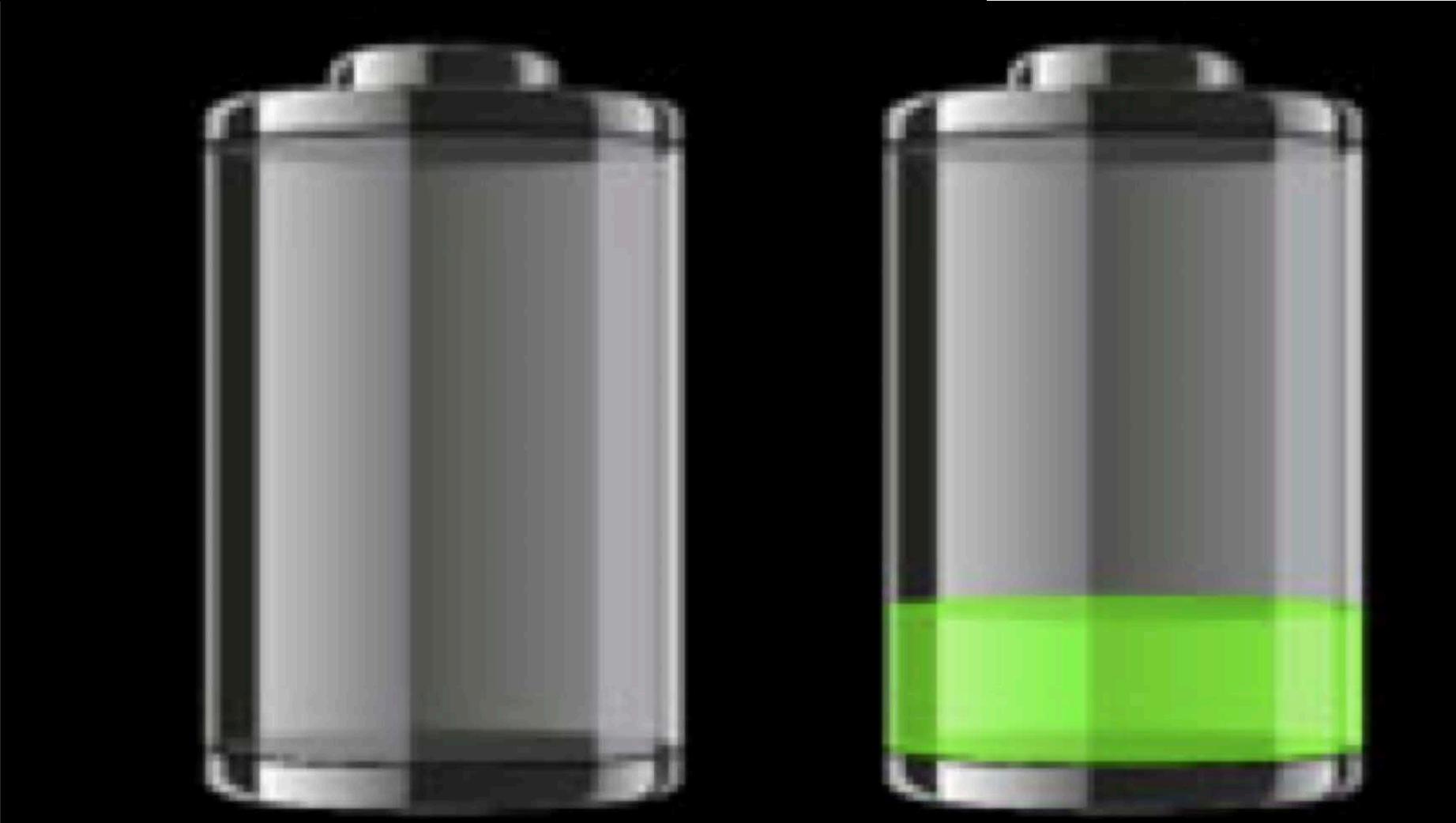


Rich Vuduc, Georgia Tech

What is an energy-efficient algorithm?



Energy Power



Rich Vuduc (...) Zoom Facilit... Presentations

Rich Vuduc (Georgia...) Zoom Facilitator 1 Presentations

Meeting Room Jerry Bernholc

Jerry Bernholc

Anatomy of a “Value” Metric

Increase performance with Disaggregation And Bandwidth Steering!

Deliver bandwidth and resources to where it is needed By taking it from where it isn’t

Performance

Measured Watt

30% of datacenter power goes to network

So max savings by creating a perfectly efficient (0 pJ/bit) optical interconnect is ONLY 30%!

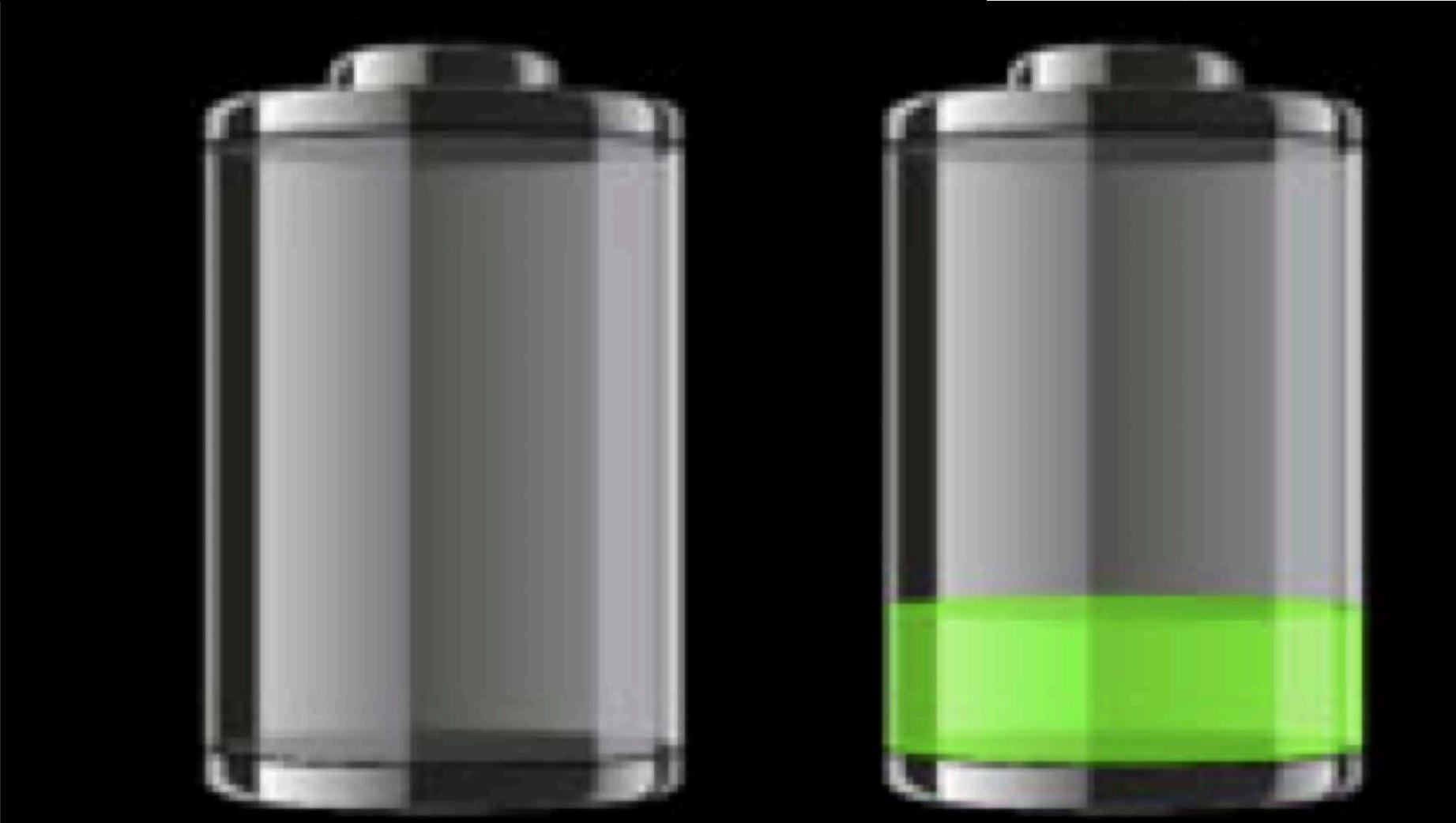
Exploit the unique capabilities to improve both numerator and denominator!

BERKELEY LAB



Rich Vuduc, Georgia Tech

Can you beat the “standard” model?



Energy Power

What is the “standard model?”

What can it tell me about building energy-efficient systems?

Recall: Sequential complexity analysis

$$\mathcal{O}(N^2) \longrightarrow \boxed{\mathcal{O}(N)}$$

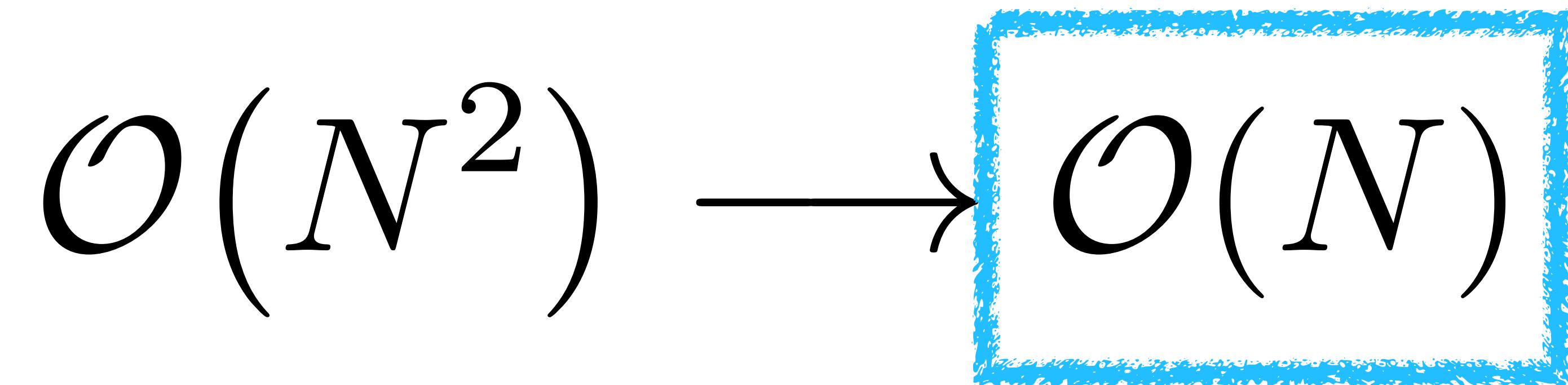
Count asymptotic operations & reduce them.

Recall: Sequential complexity analysis

$$\mathcal{O}(N^2) \longrightarrow \boxed{\mathcal{O}(N)}$$

What about energy?

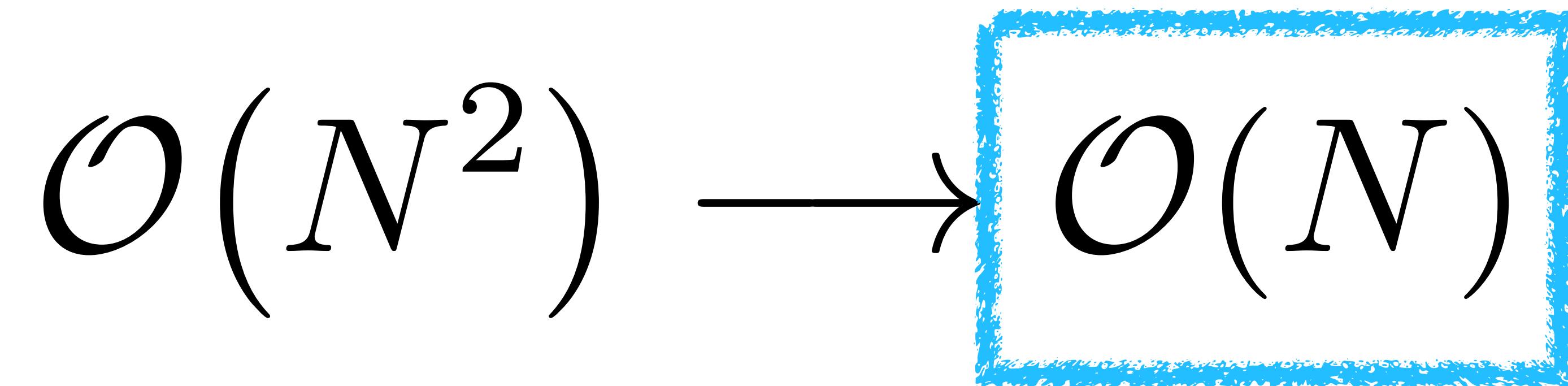
Recall: Sequential complexity analysis



Reduces **energy**, e.g., fewer ops, less storage

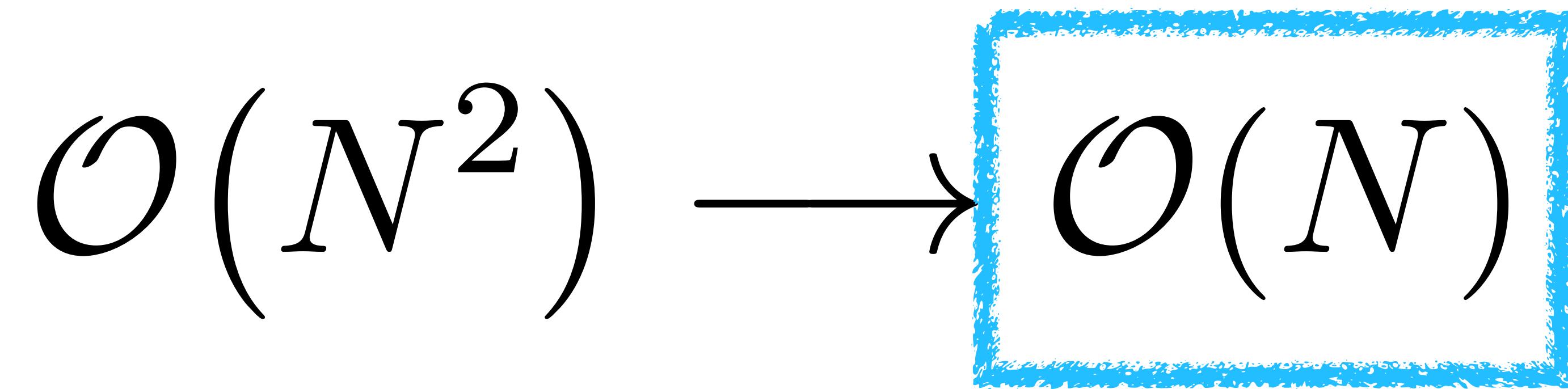
The “classical process” **is** energy reduction!

Recall: Sequential complexity analysis

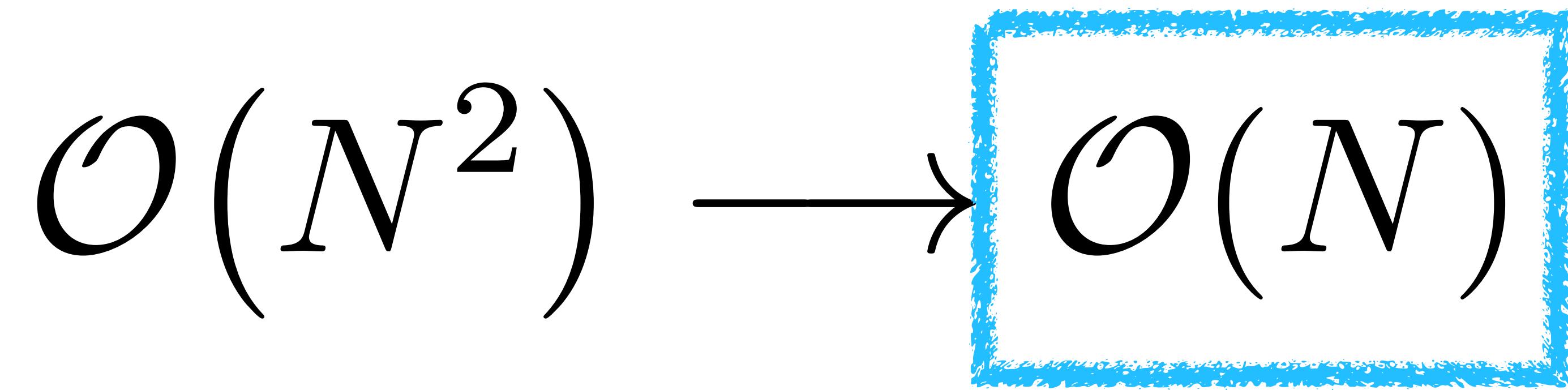


But if that's so, what does **explicit reasoning** about energy accomplish?

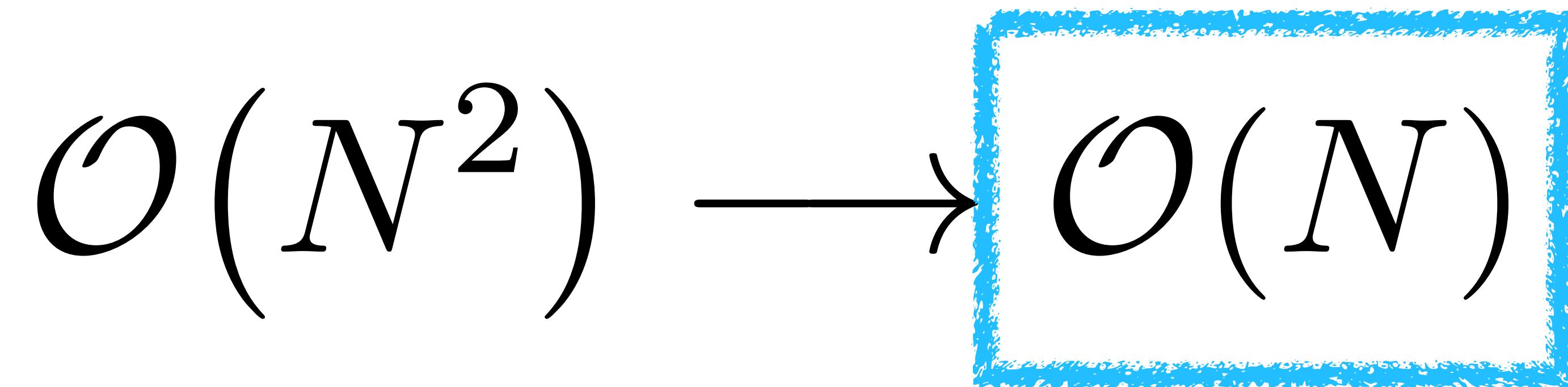
Recall: Sequential complexity analysis



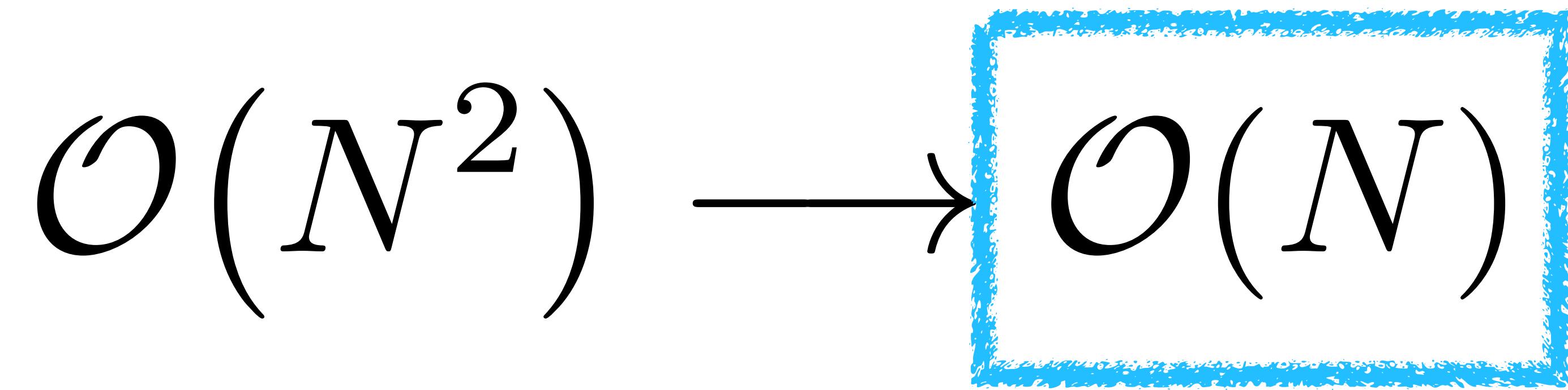
Answer: It **facilitates energy-centric codesign** and will **make or break** your new platform.



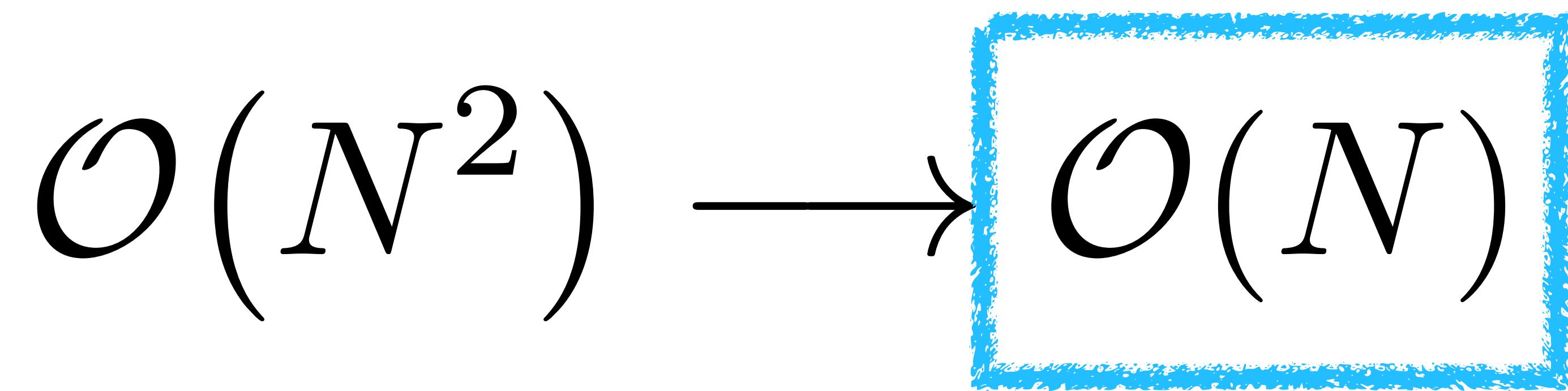
What is an “**operation**?”



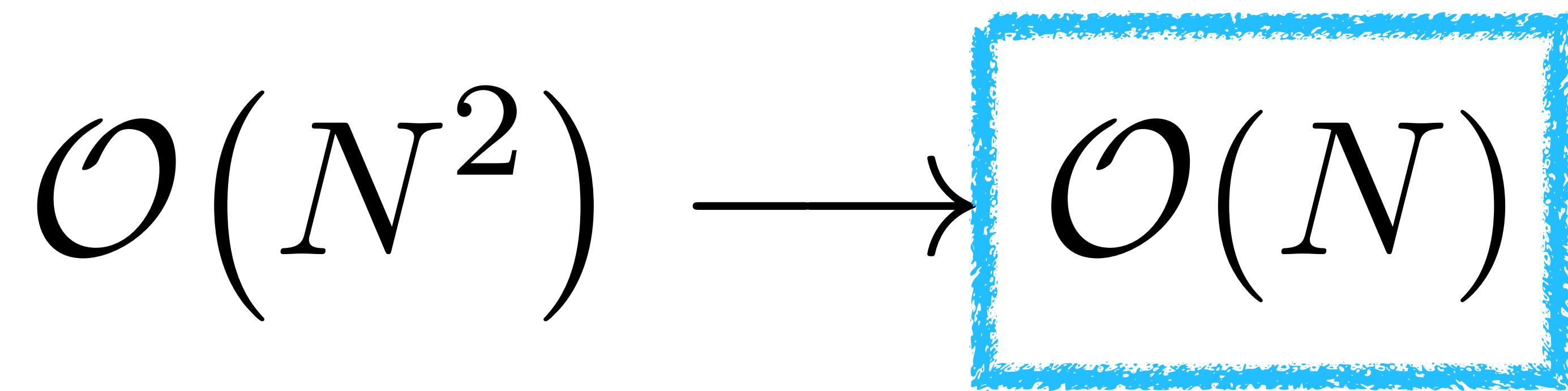
It is an “**interface**” that enabled >70 years of **productivity** in algorithms (and software).



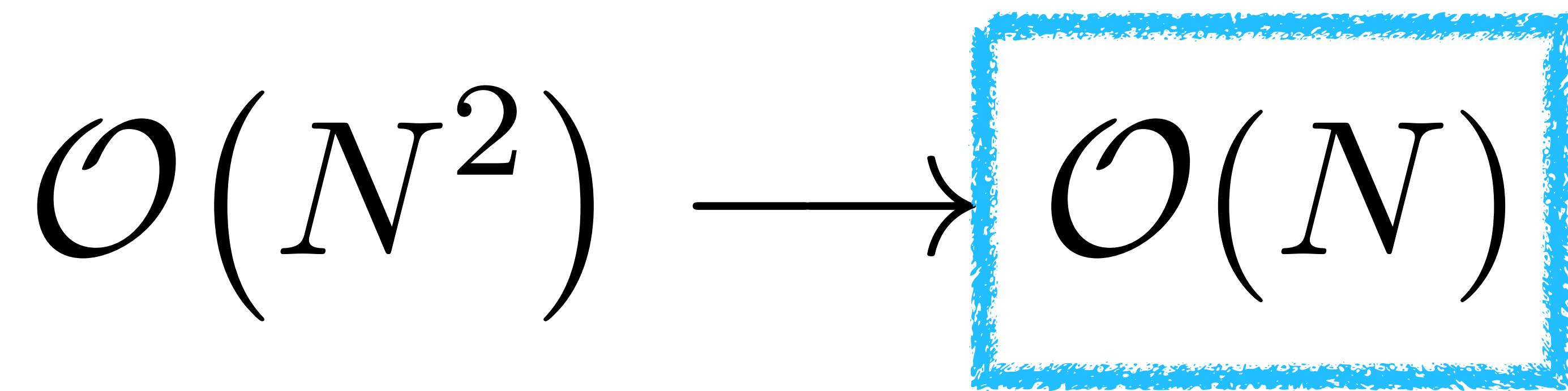
An implication? Your new computing substrate
abandons this interface at its **peril**.



What else does an “**operation**” do, in its role as an interface?

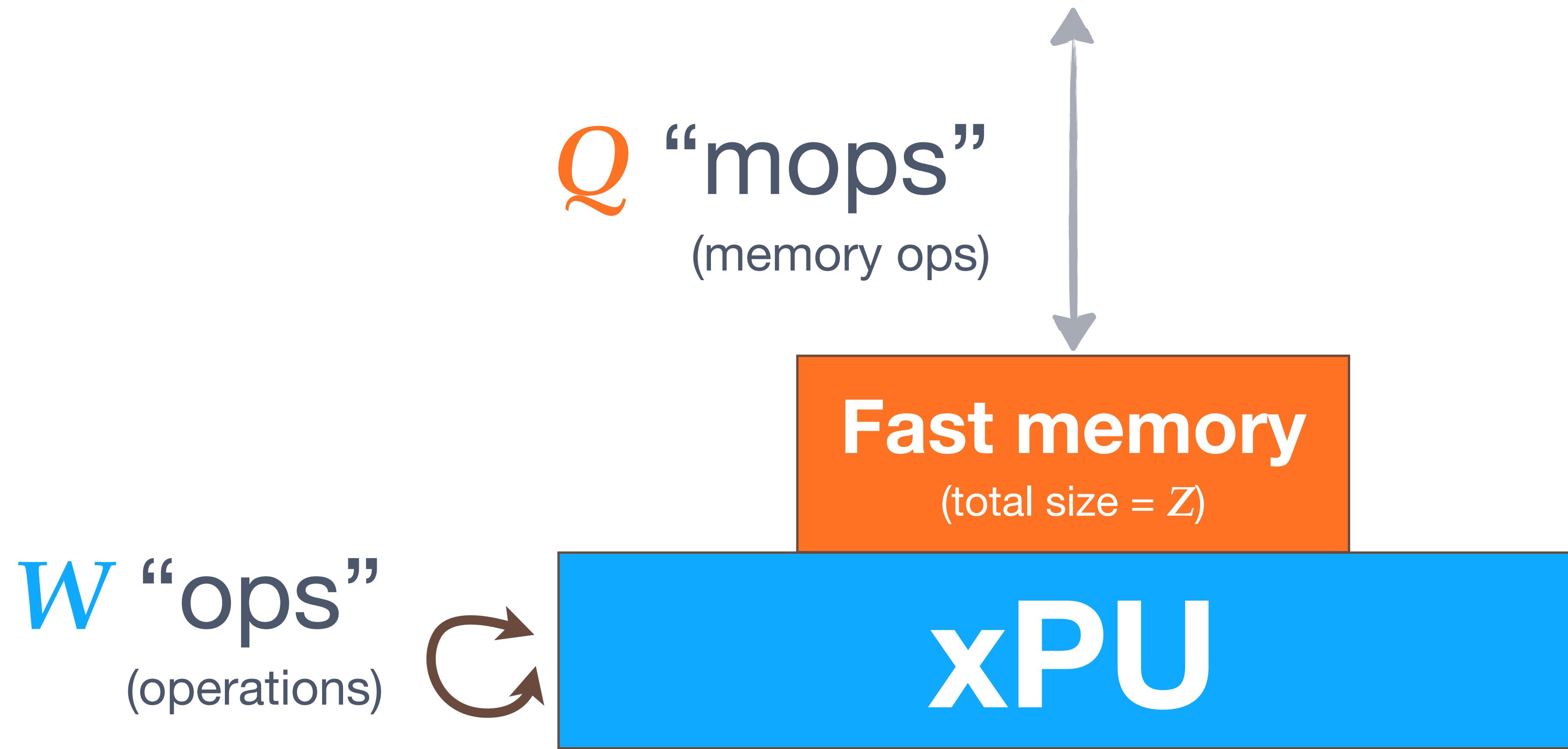


It defines a **unit of cost**.

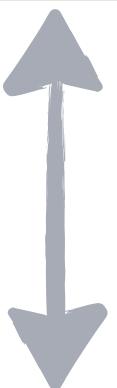


It defines a unit of cost. “Exciting” things happen when costs are **nonuniform**.

Slow memory



Slow

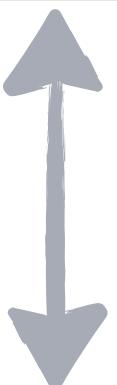


Fast

xPU

W ops
Q mops

Slow

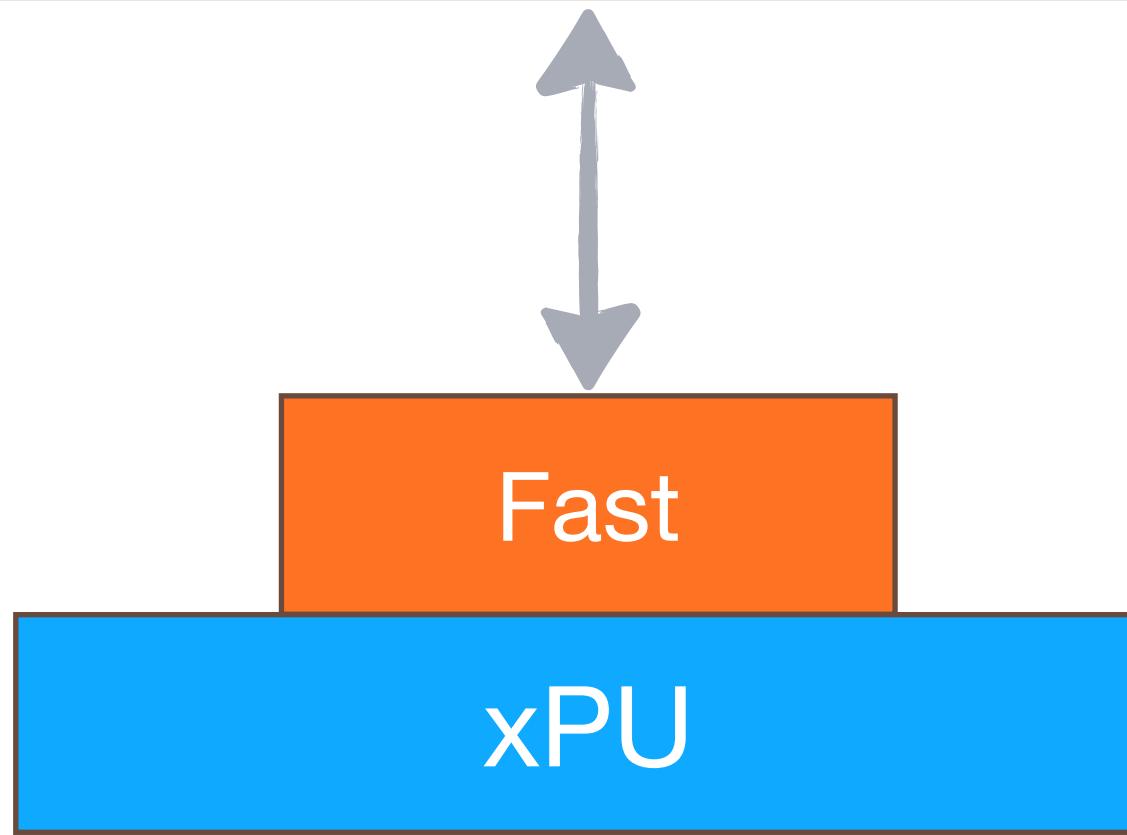


Fast

xPU

$\frac{W}{Q}$ op : mop
(e.g., “flop:Byte”)

Slow



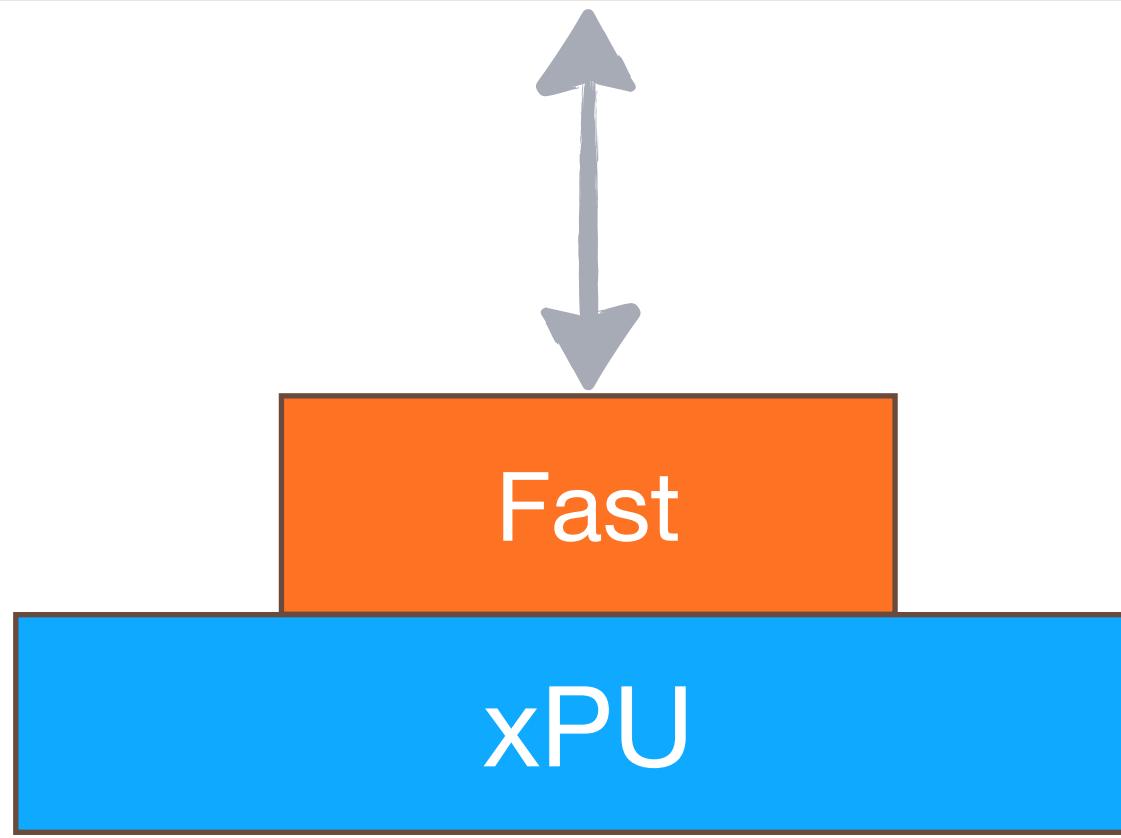
$W \tau_{\text{op}}$

$Q \tau_{\text{mop}}$

Time per op
(In truth,
inverse
throughput)

Time per mop

Slow

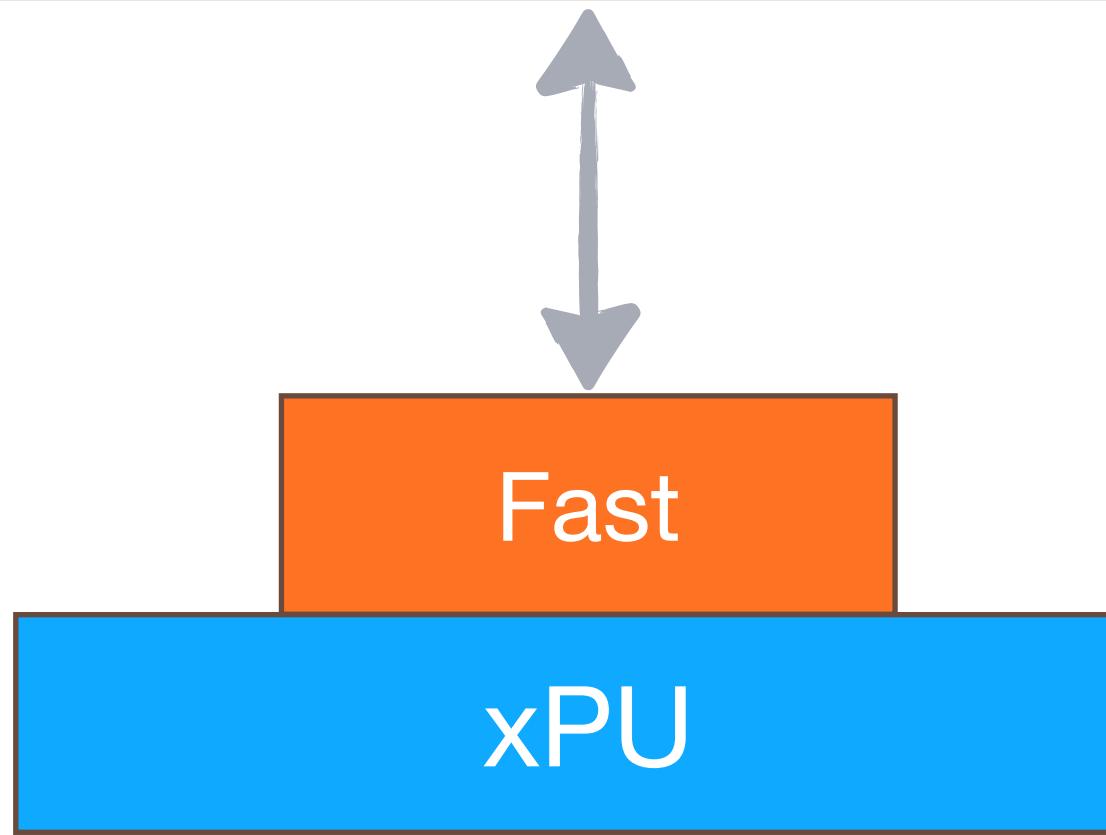


Machine balance
(inverse)

$$\frac{W\tau_{op}}{Q\tau_{mop}}$$

The equation is enclosed in a yellow rectangular frame with a hand-drawn texture. It consists of two terms separated by a horizontal line. The top term is $W\tau_{op}$ where W is blue and τ_{op} is blue. The bottom term is $Q\tau_{mop}$ where Q is orange and τ_{mop} is orange.

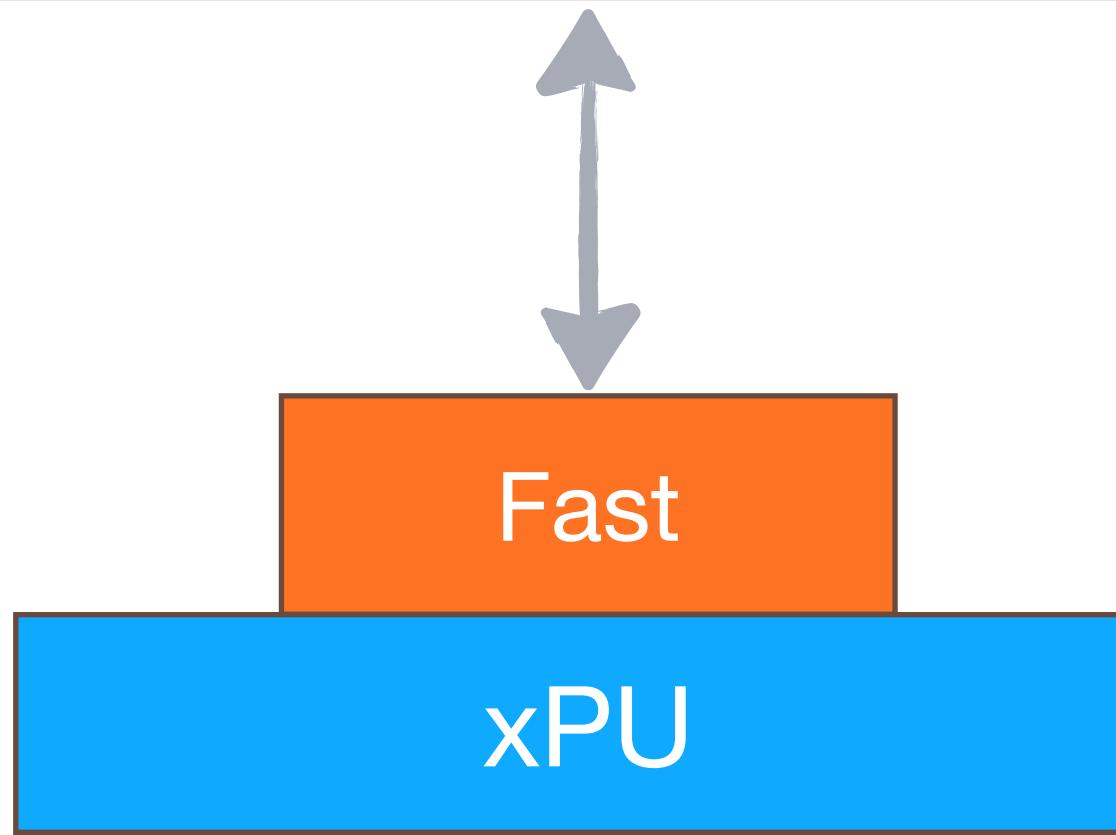
Slow



$\frac{W??}{Q??}$

What about **energy**?

Slow



$$\frac{W_{\epsilon_{op}}}{Q_{\epsilon_{mop}}} = \frac{\text{Energy per op}}{\text{Energy per mop}}$$

Energy per op

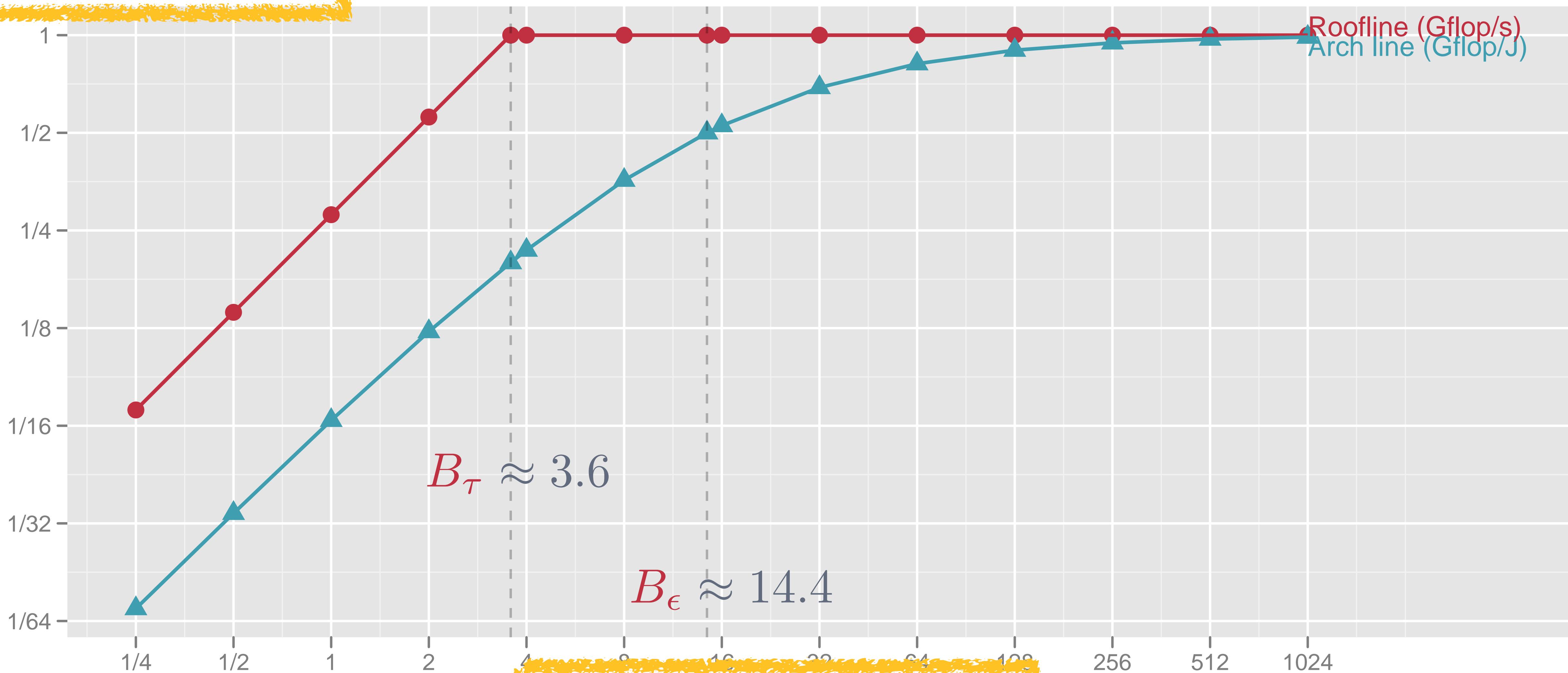
Energy per mop

Rooflines in time & (igloos in?) energy



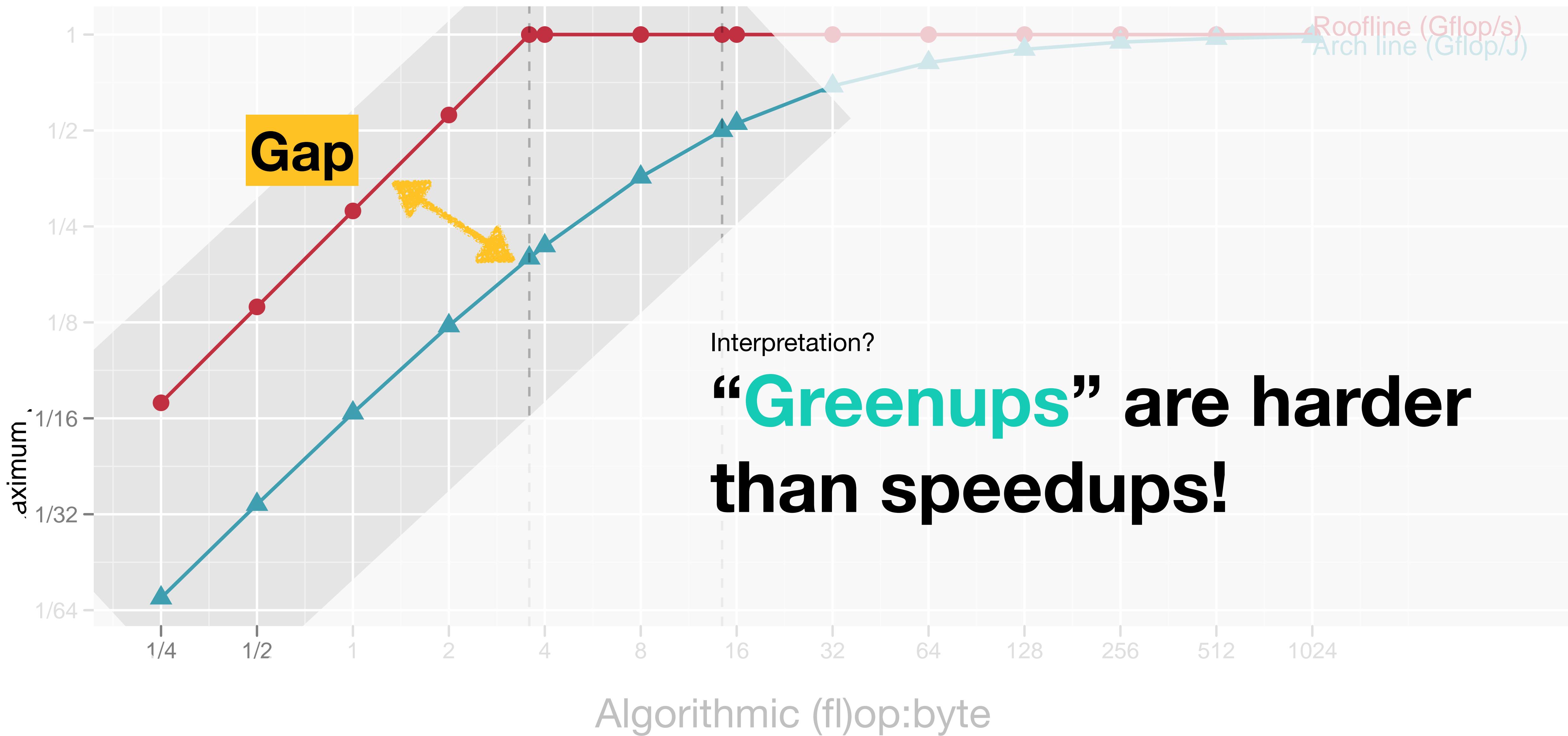
S. Williams, A. Waterman, D. Patterson. *Roofline: An insightful visual performance model for multicore architectures*. CACM, v52(4), 2009. doi:[10.1145/1498765.1498785](https://doi.org/10.1145/1498765.1498785)
J. Choi, D. Bedard, R. Fowler, R. Vuduc. *A roofline model of energy*. In IPDPS'13. doi:[10.1109/IPDPS.2013.77](https://doi.org/10.1109/IPDPS.2013.77)
J. Choi, M. Dukhan, X. Liu, R. Vuduc. *Algorithmic time, energy, and power on candidate HPC building blocks*. In IPDPS'14. doi:[10.1109/IPDPS.2014.54](https://doi.org/10.1109/IPDPS.2014.54)

Fraction of peak

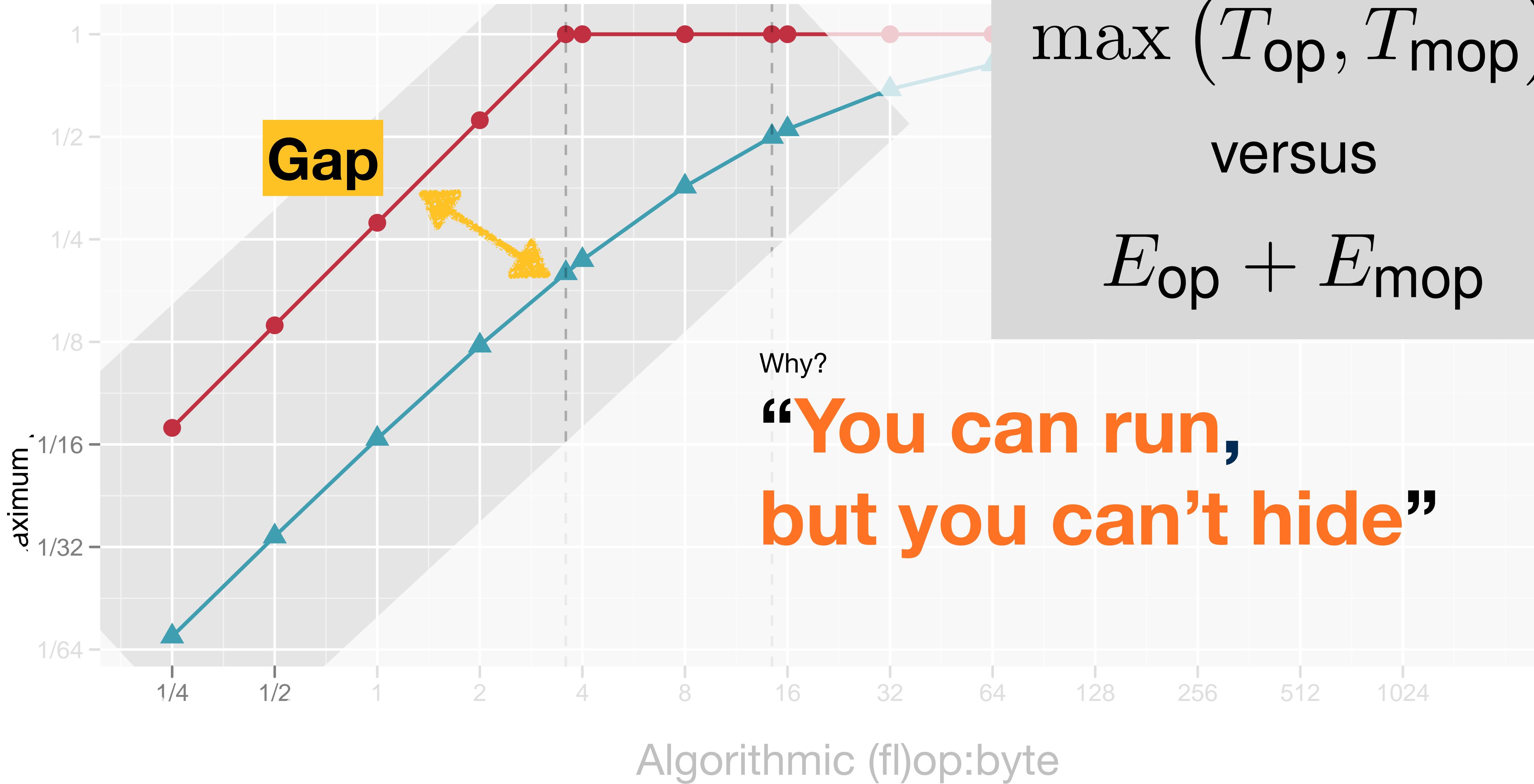


Algorithmic (fl)op:byte

Fraction of peak



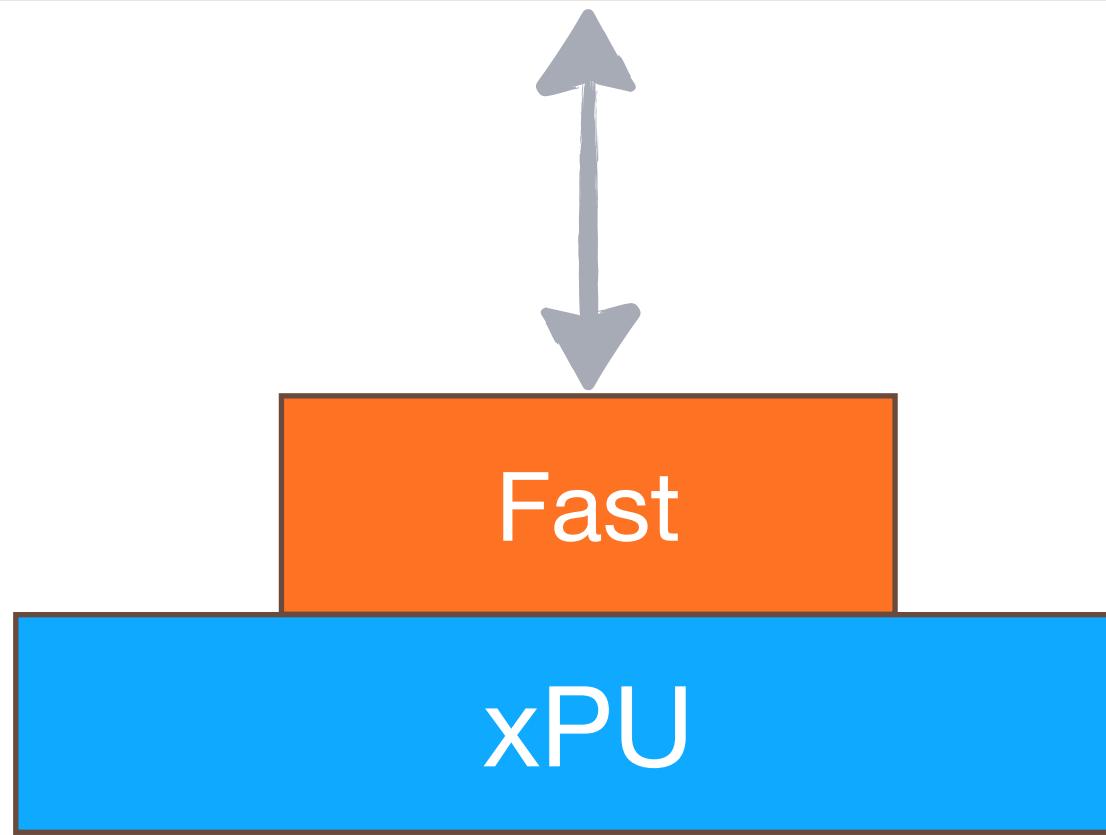
Fraction of peak



Why?

**“You can run,
but you can’t hide”**

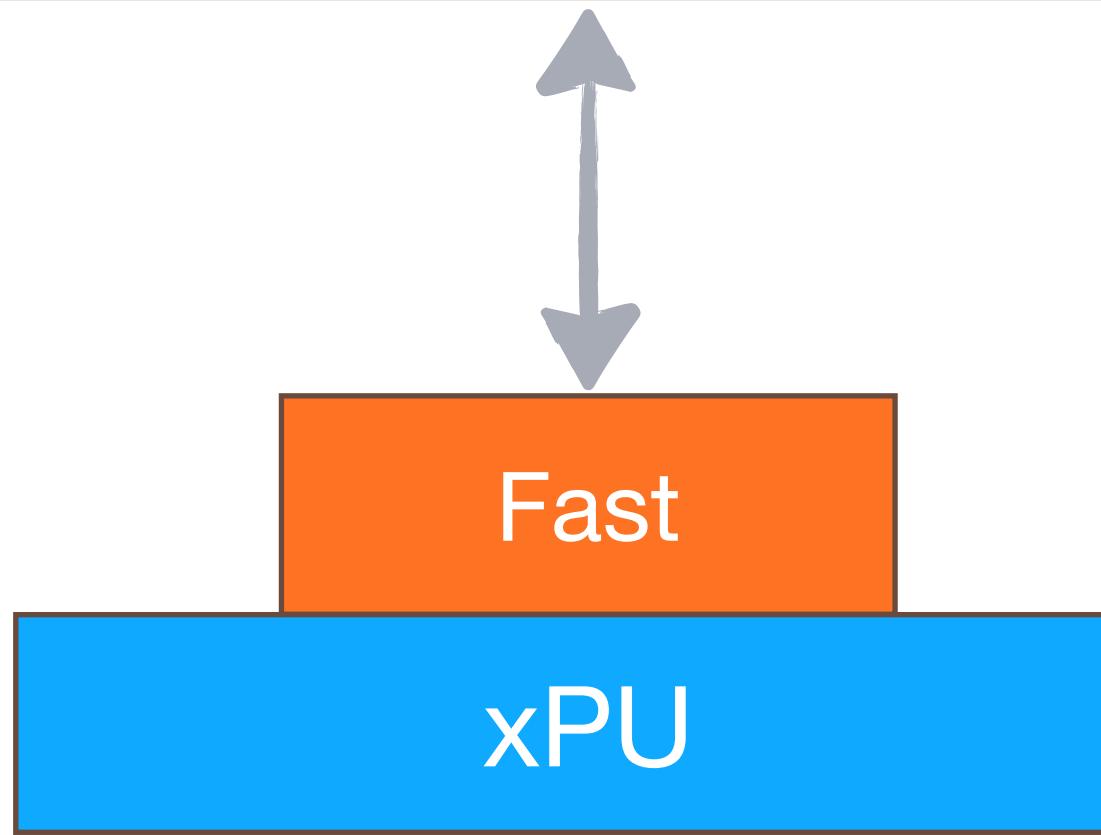
Slow



W??
— —
Q??

What about power?

Slow



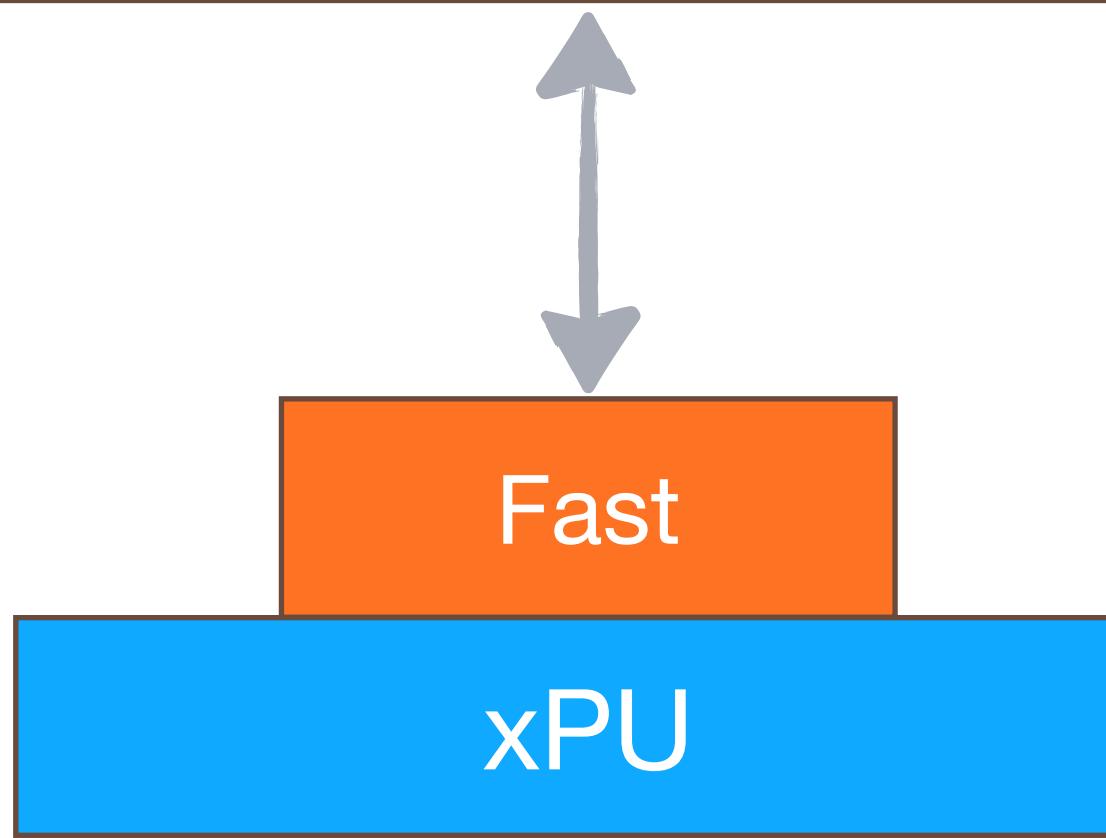
Power per op

$$\frac{W \epsilon_{\text{op}} / \tau_{\text{op}}}{Q \epsilon_{\text{mop}} / \tau_{\text{mop}}}$$

What about power?

Power per mop

Slow



Real-life complications:

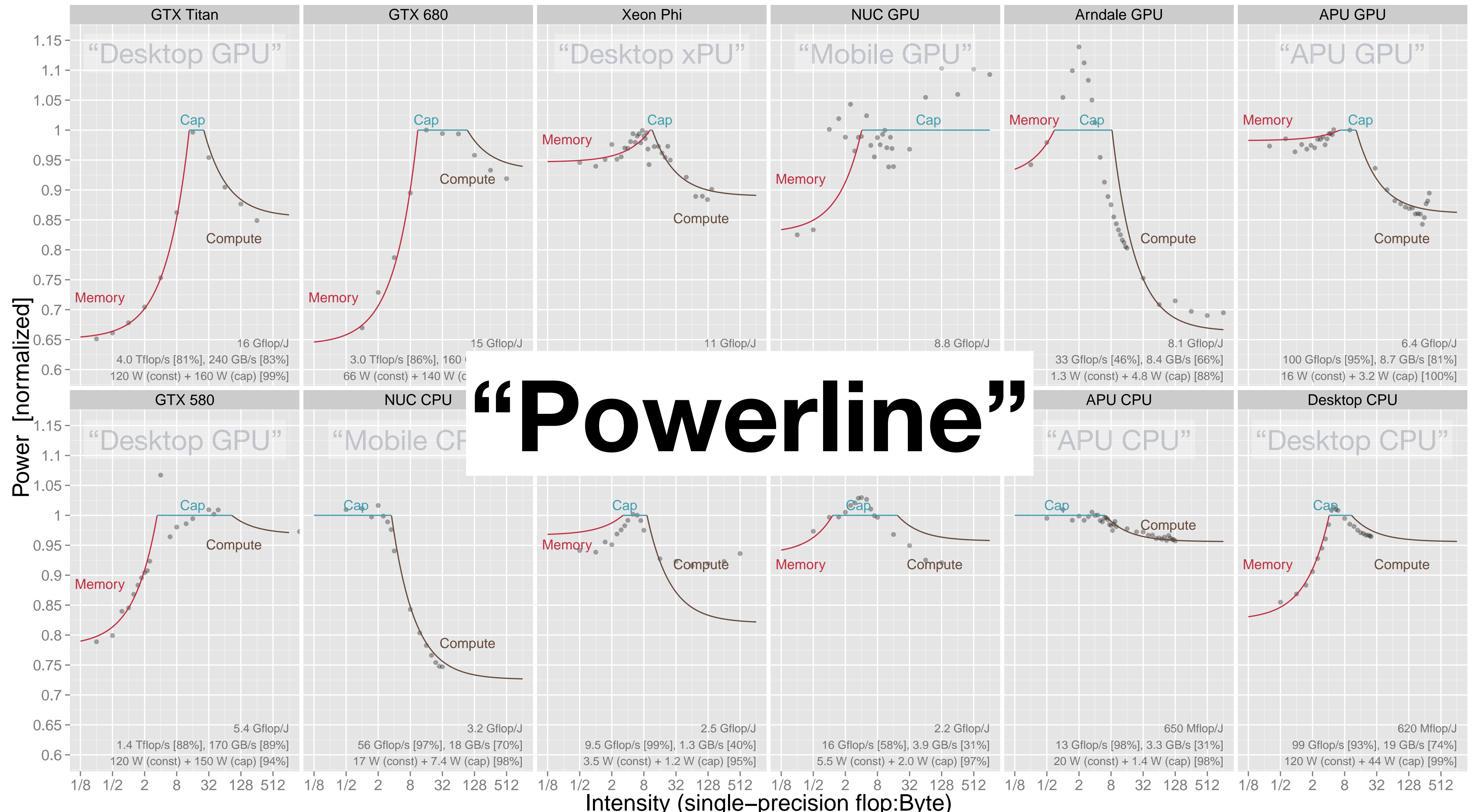
**Baseline power
& power caps**

Power per op

$$\frac{W \epsilon_{op}/\tau_{op}}{Q \epsilon_{mop}/\tau_{mop}}$$

Power per op

Power per mop



GTX Titan

GTX 6

Xeon II

NUC G

Arndale GPU

APU GPU

“Desktop GPU”

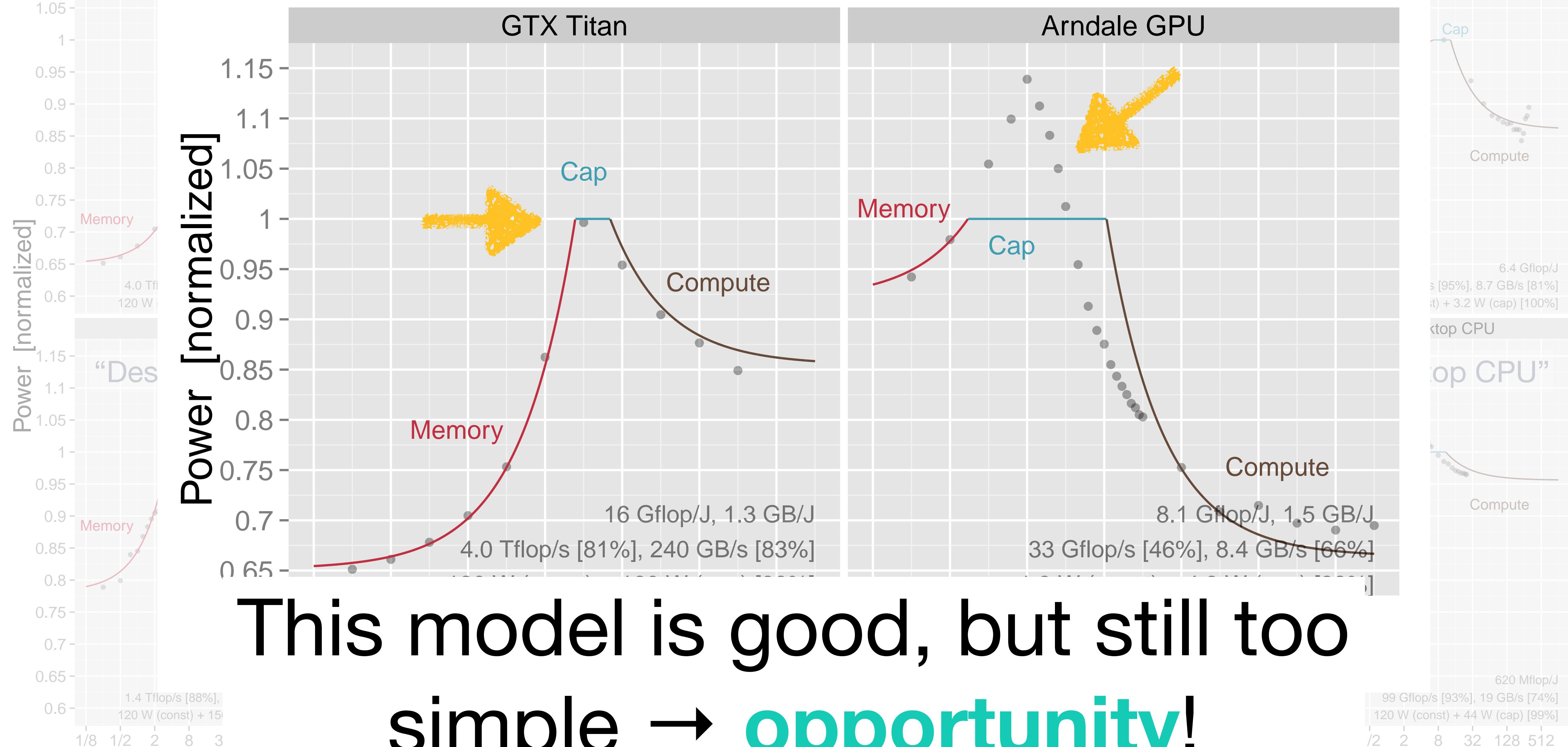
“Desktop GPU” “D

Desktop xPU™

“Mobile GPU”

“Mobile GPU”

“APU GPU”



GTX Titan

GTX 680

Xeon Phi

NUC GPU

Arndale GPU

APU GPU

“Desktop GPU”

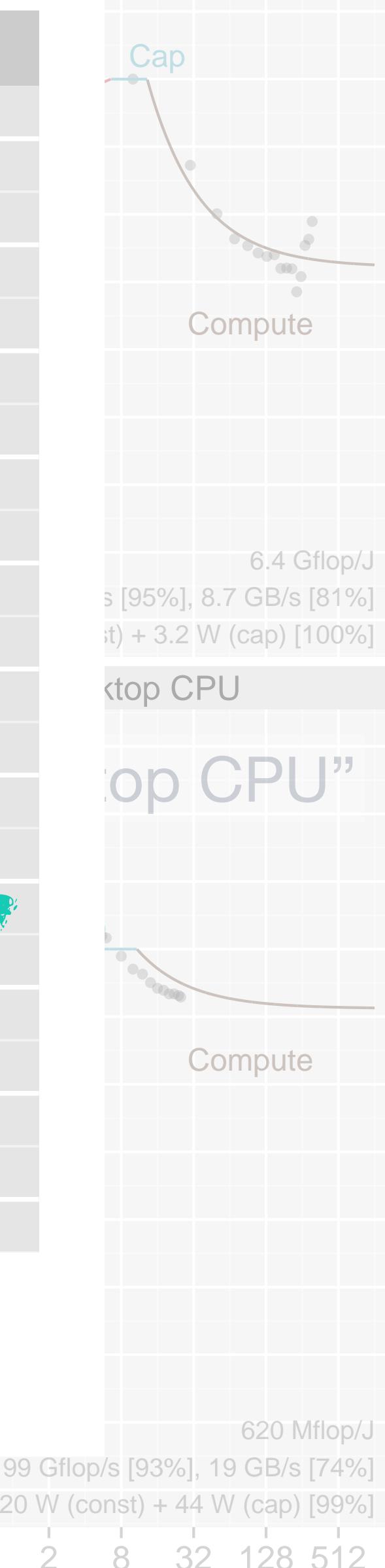
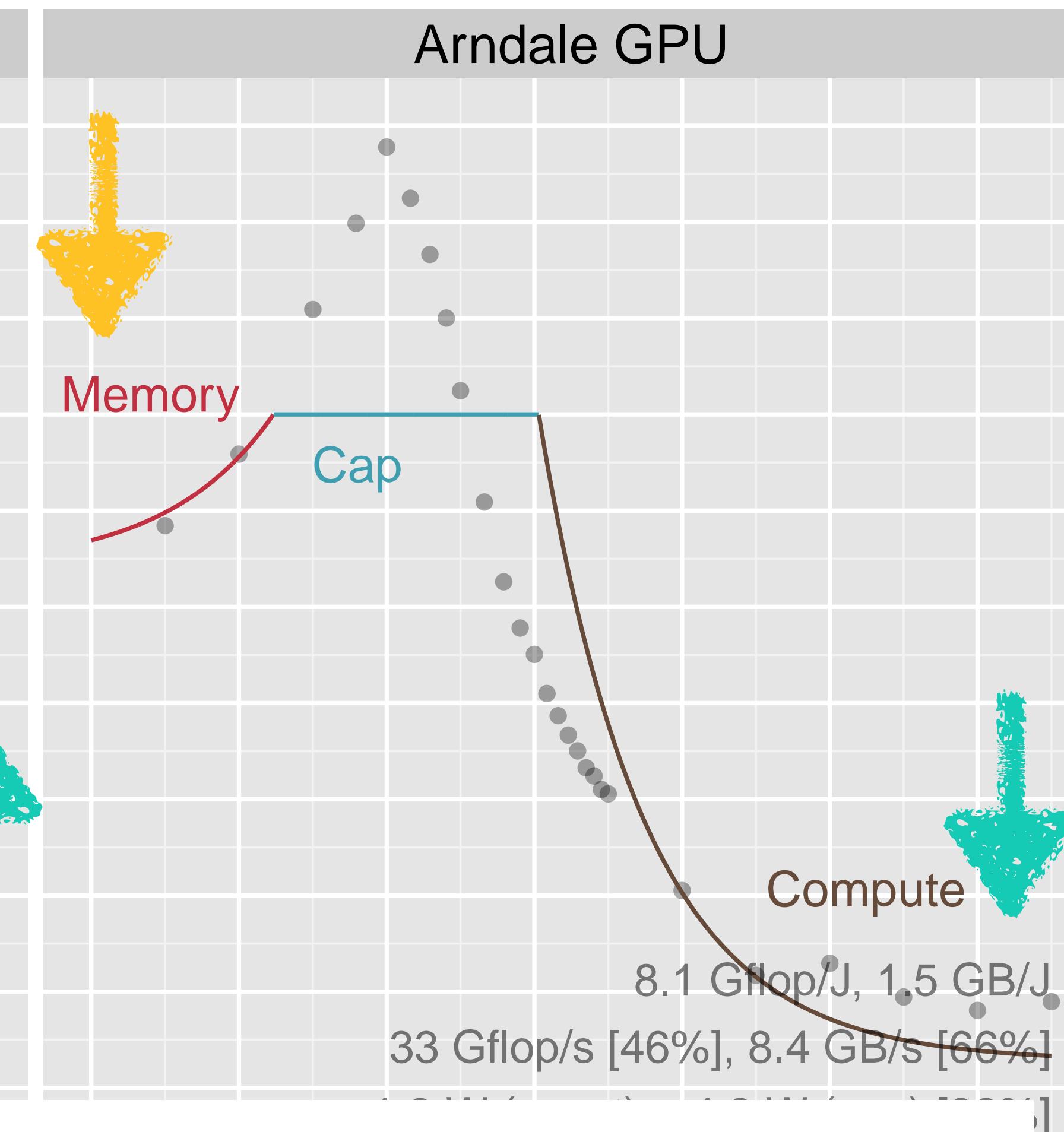
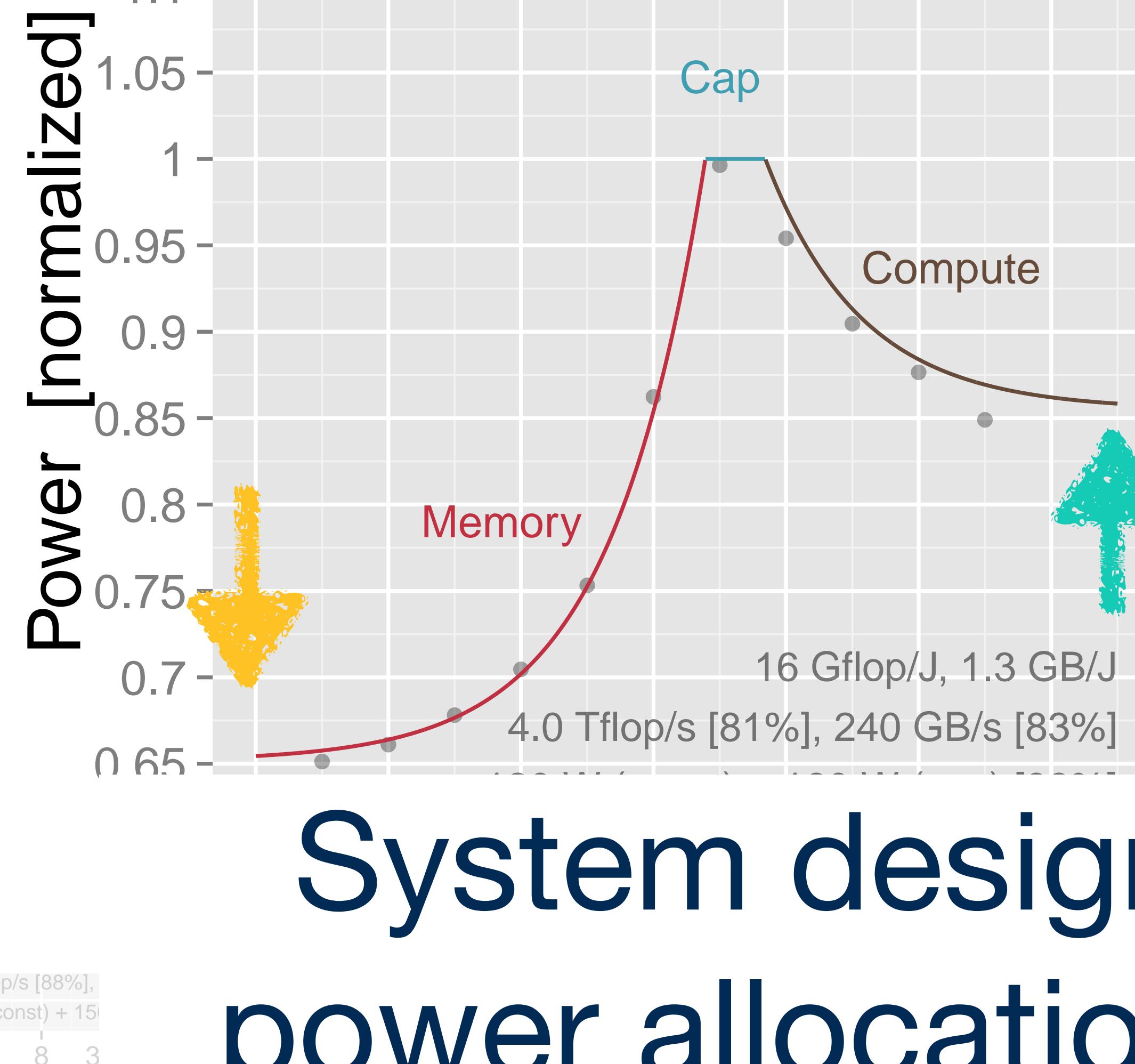
“Desktop GPU”

“Desktop xPU”

“Mobile GPU”

“Mobile GPU”

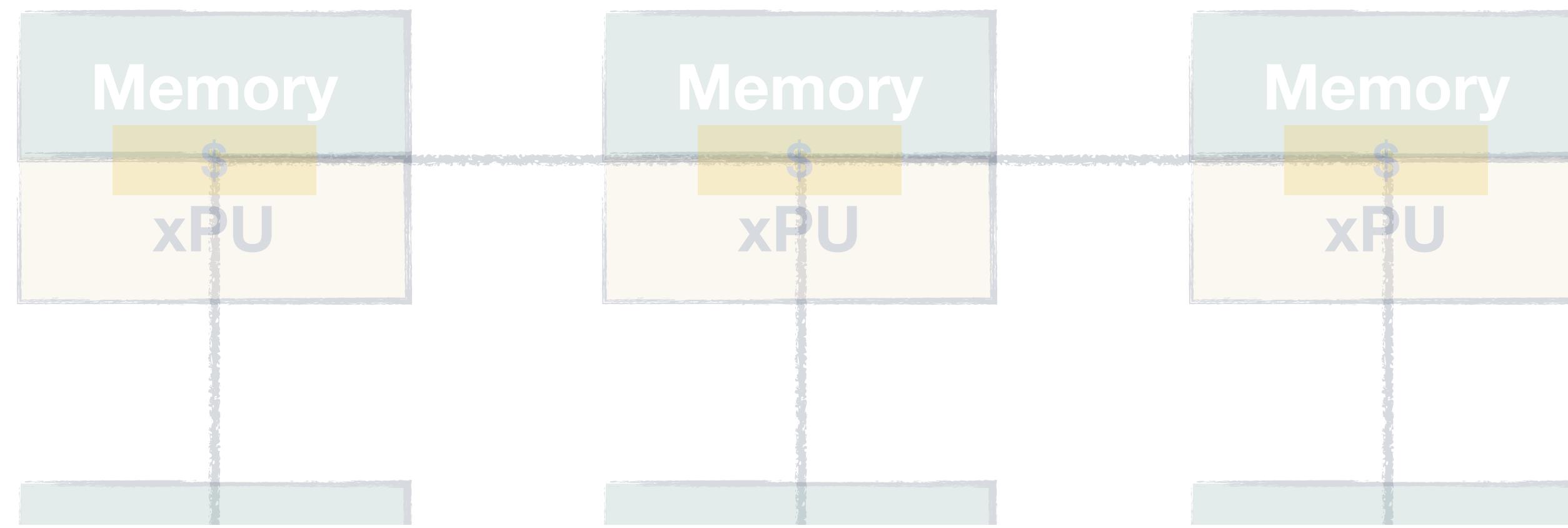
“APU GPU”



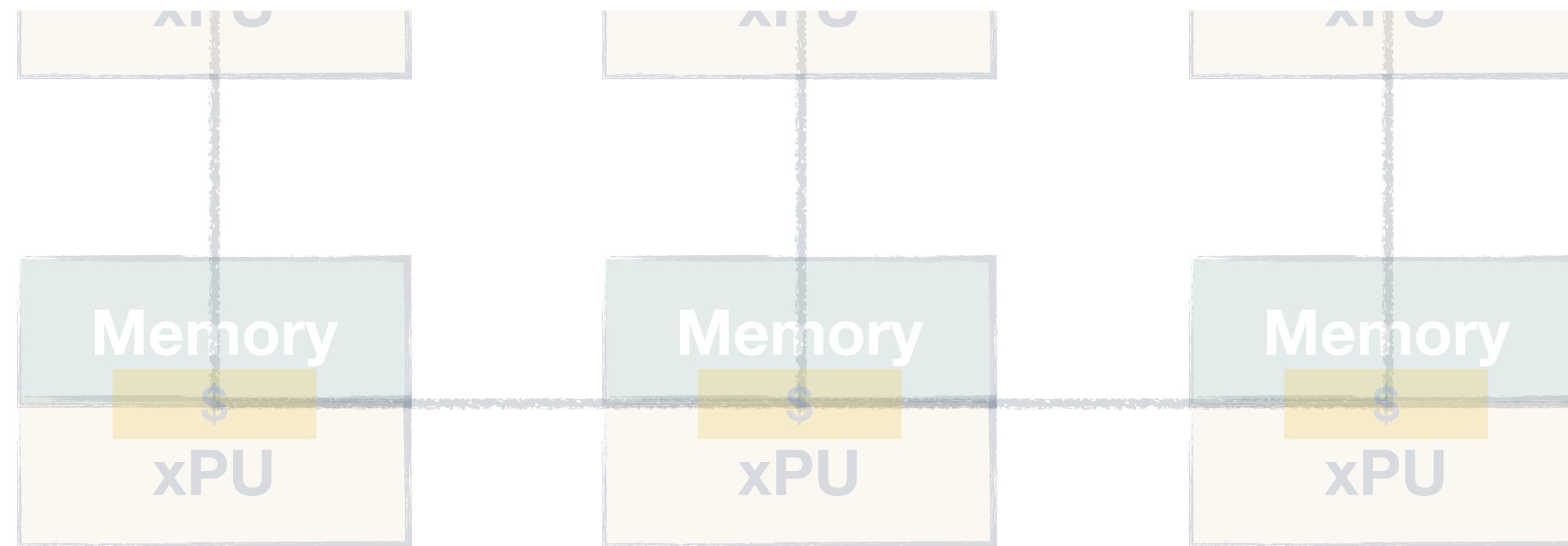
System designers choose power allocations → +Oppy!

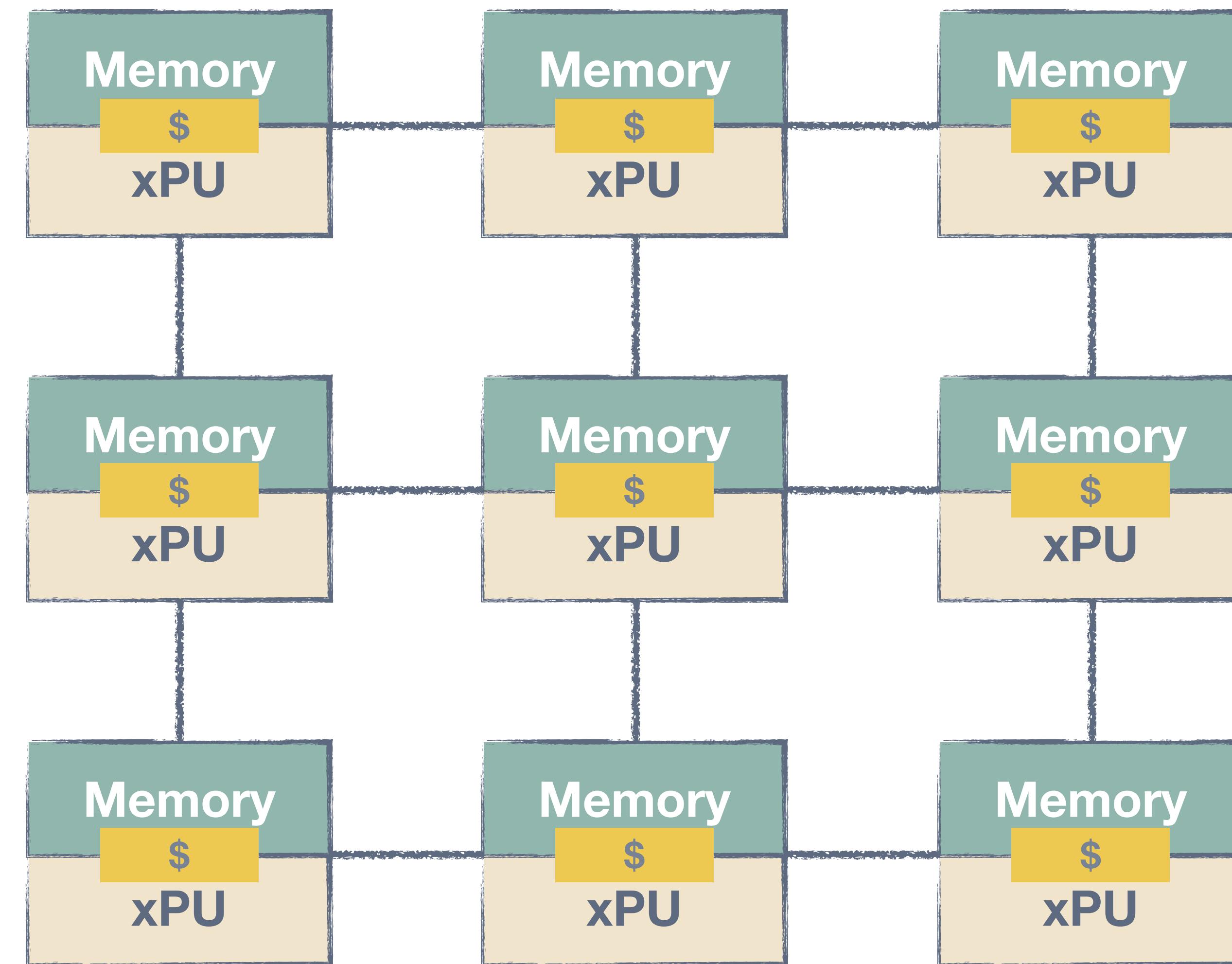
What is the “standard model?”

What can it tell me about building energy-efficient systems?



What if you could start from scratch?





Problem:

Given an algorithm
and a fixed power
& transistor budget,
pick the cores, caches, topology,
& all speeds and feeds
to minimize execution time.



Problem:

Given an algorithm
and a fixed power
& transistor budget,
pick the cores, caches, topology,
& all speeds and feeds
to minimize execution time.



Problem:

Given an algorithm
and a fixed power
& transistor budget,
**pick the cores, caches, topology,
& all speeds and feeds**
to minimize execution time.



Problem:

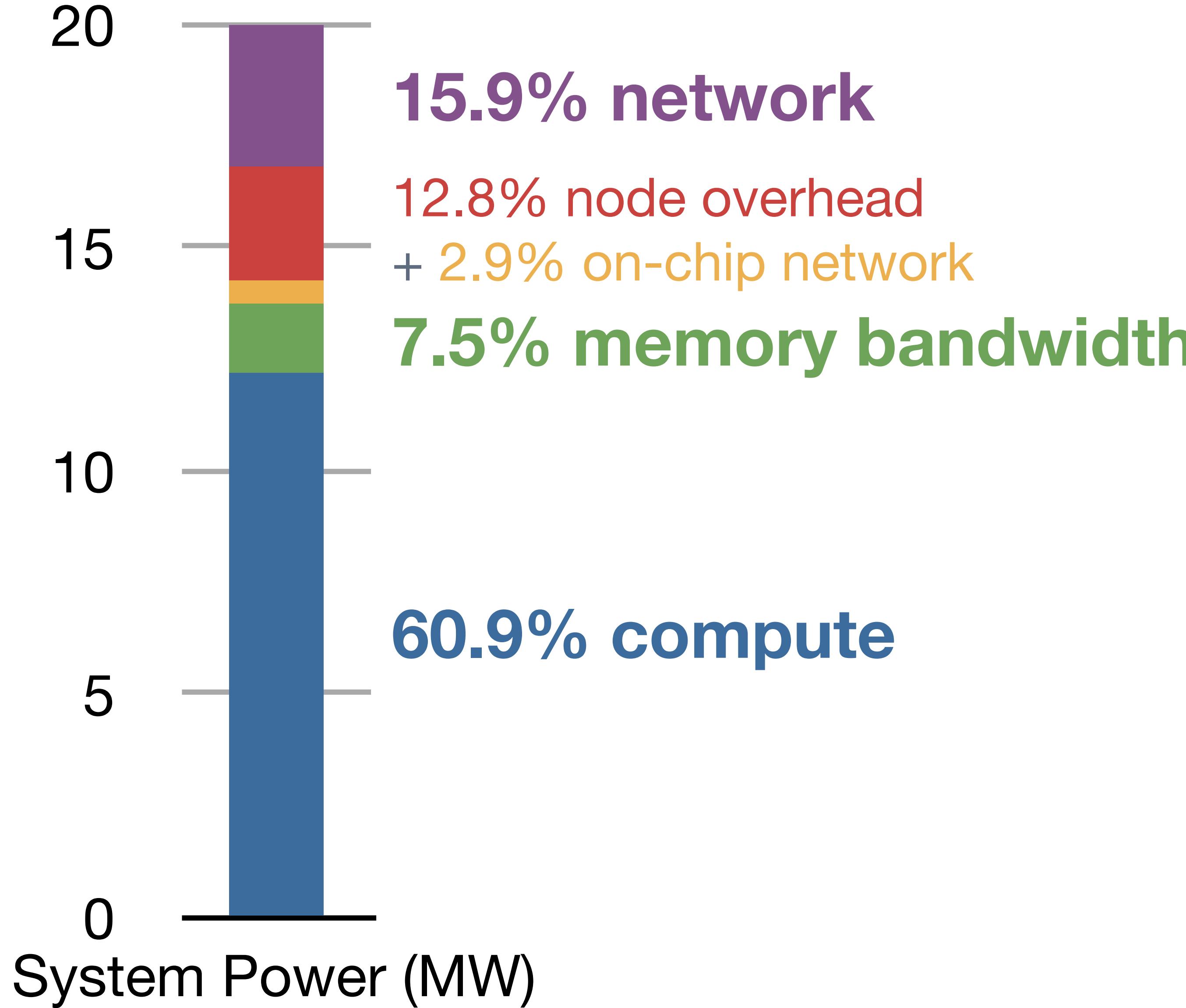
Given an algorithm
and a fixed power
& transistor budget,
pick the cores, caches, topology,
& all speeds and feeds
to minimize execution time.



Power allocation for an “optimal” matrix multiply machine?

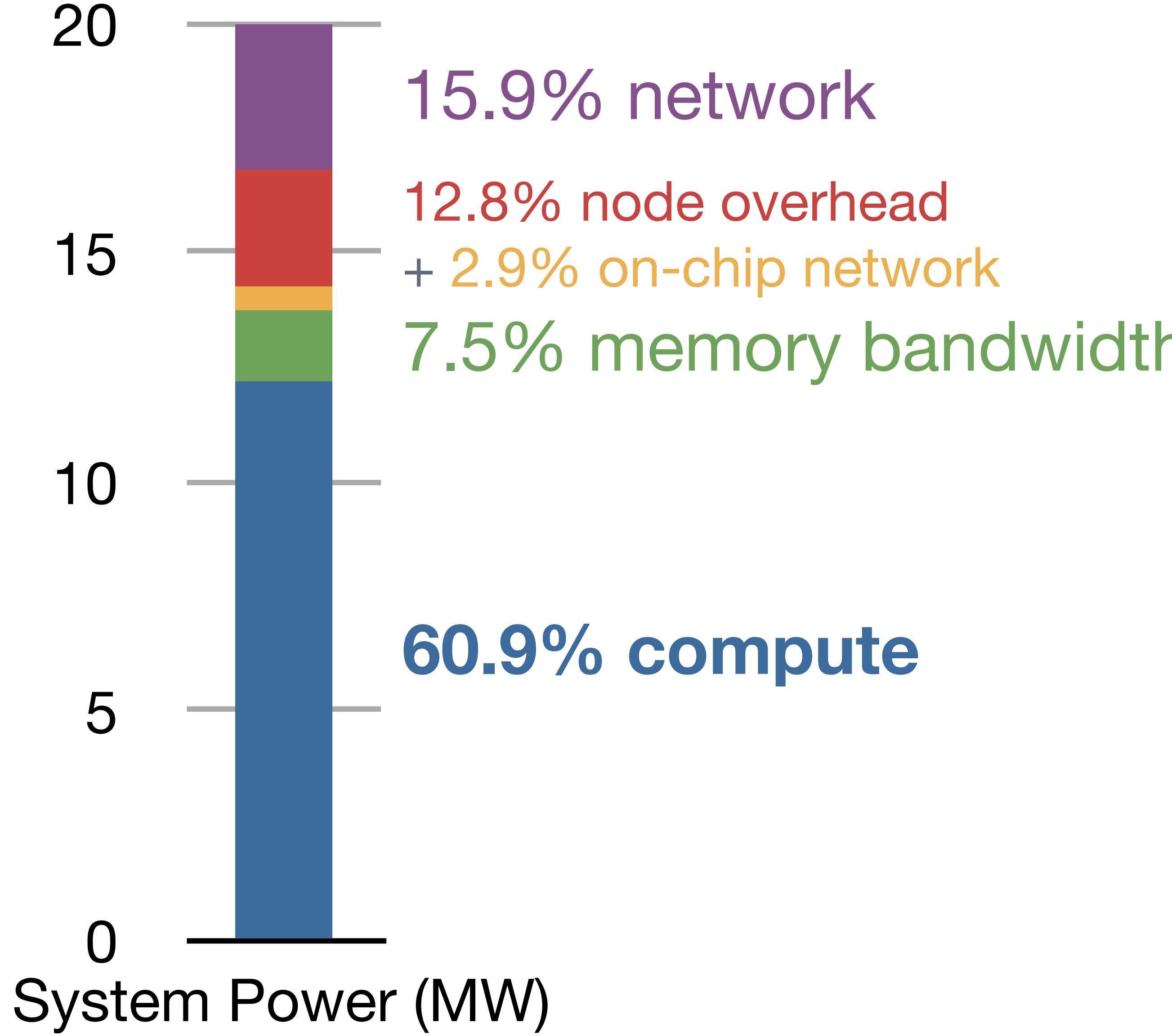
Power allocation for an “optimal” matrix multiply machine

41



Power allocation for an “optimal” matrix multiply machine

42



ORNL Summit (13-14 MW):
67.0% GPU compute
14.9% CPU compute

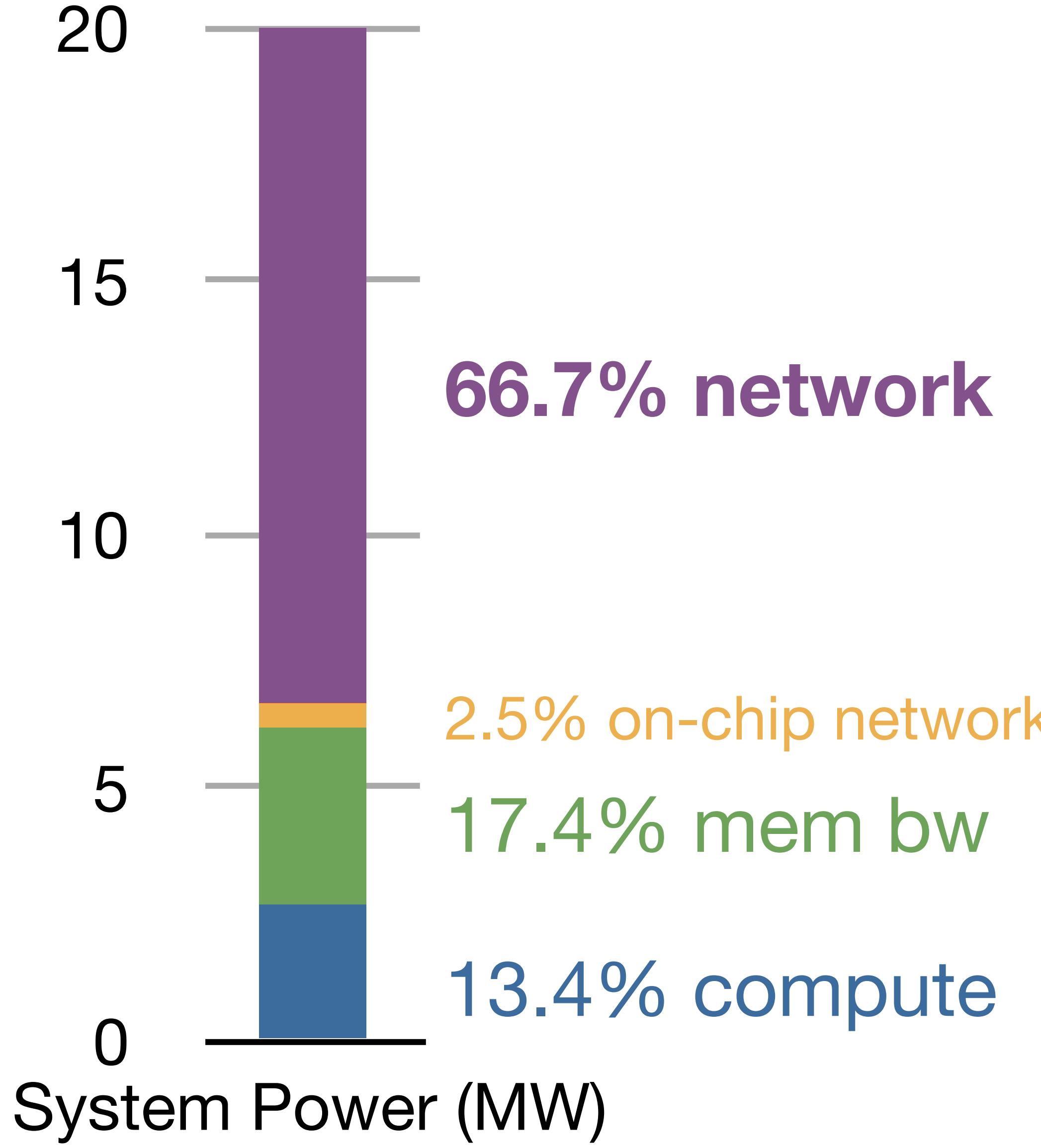
4.8% memory
5.3% network + disk
8% node overhead

P.S.: $R_{max} / R_{peak} \sim 75\%$

Power allocation for an “optimal” 3D FFT machine?

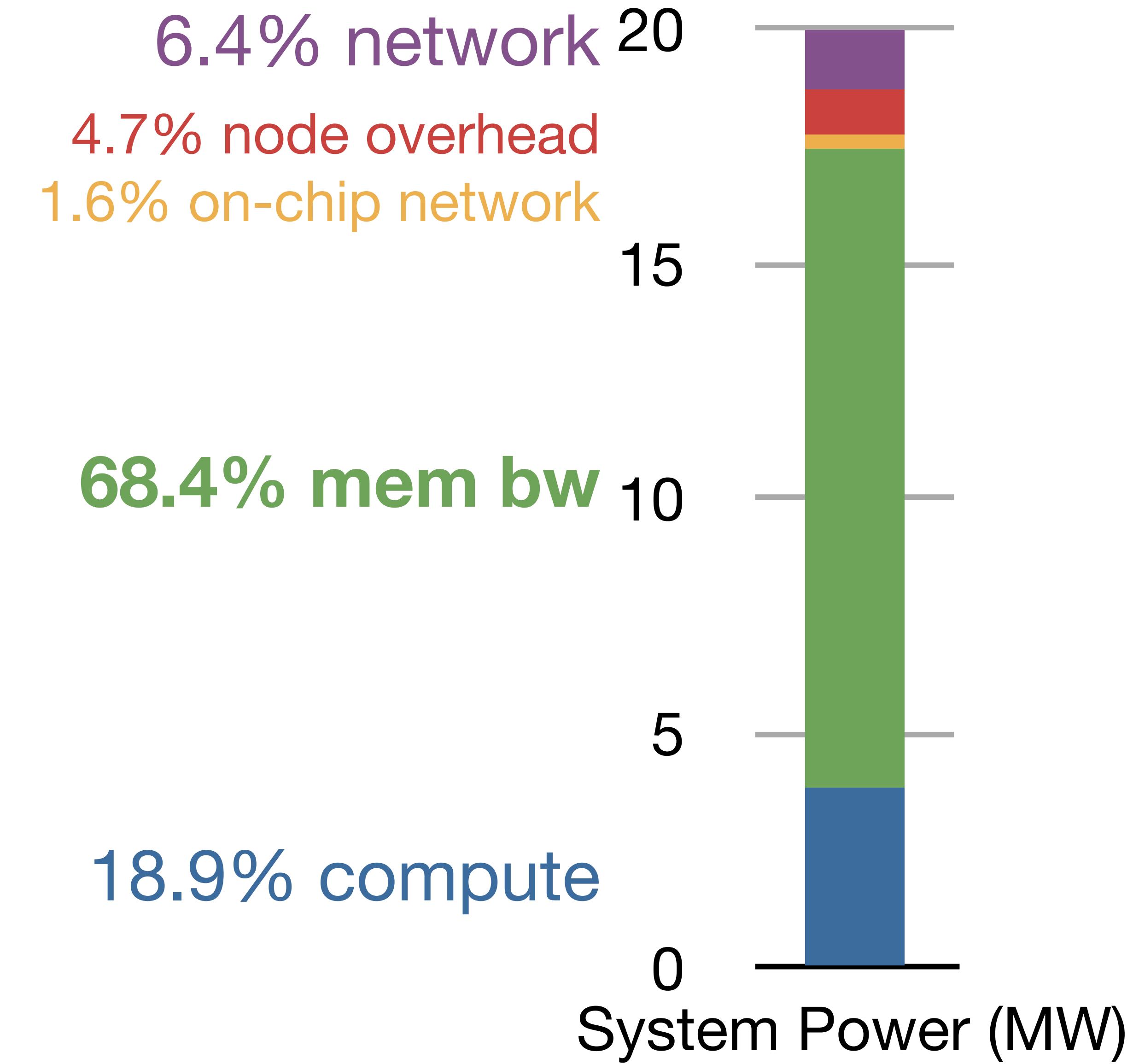
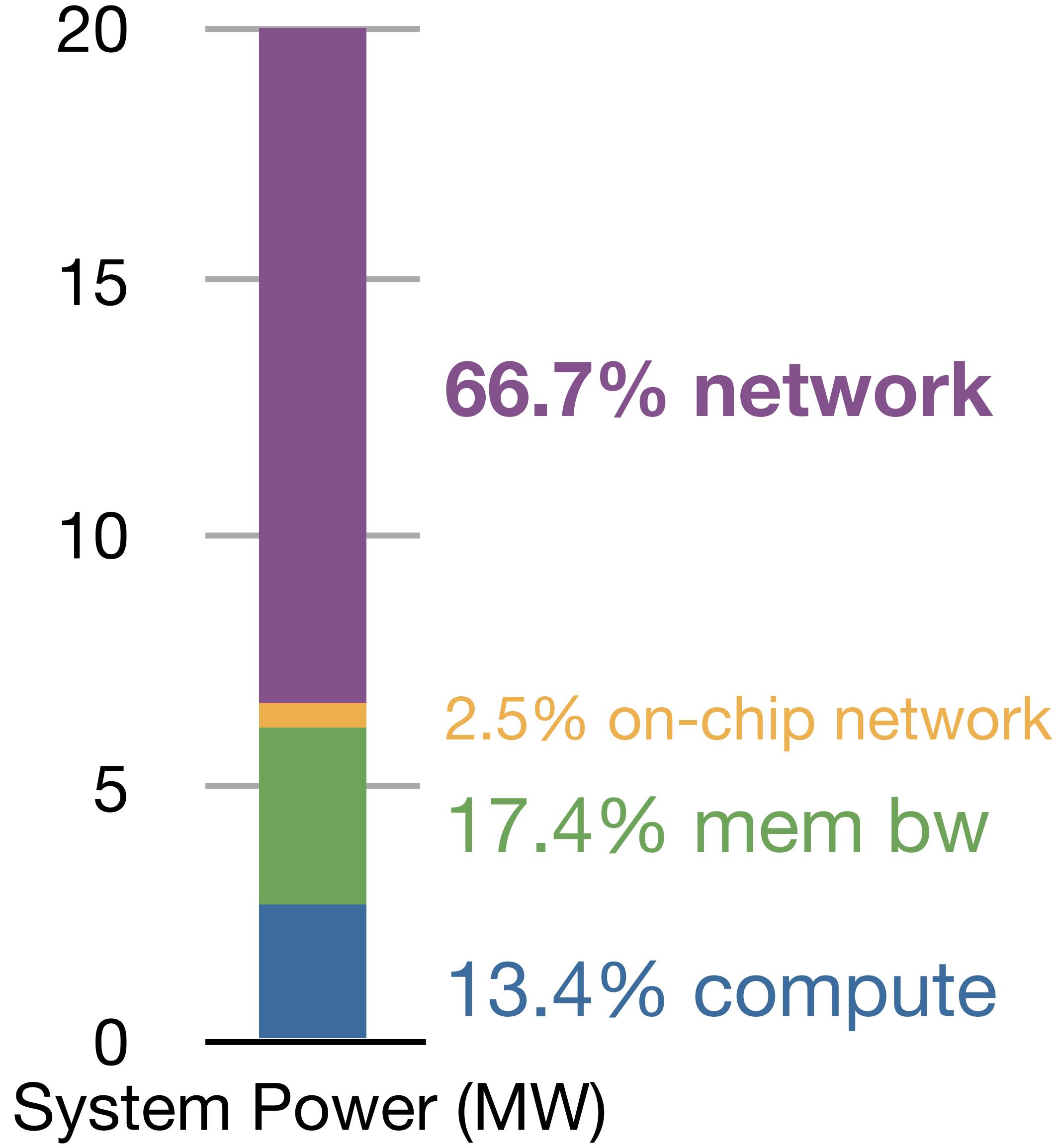
Power allocation for an “optimal” 3D FFT machine

44



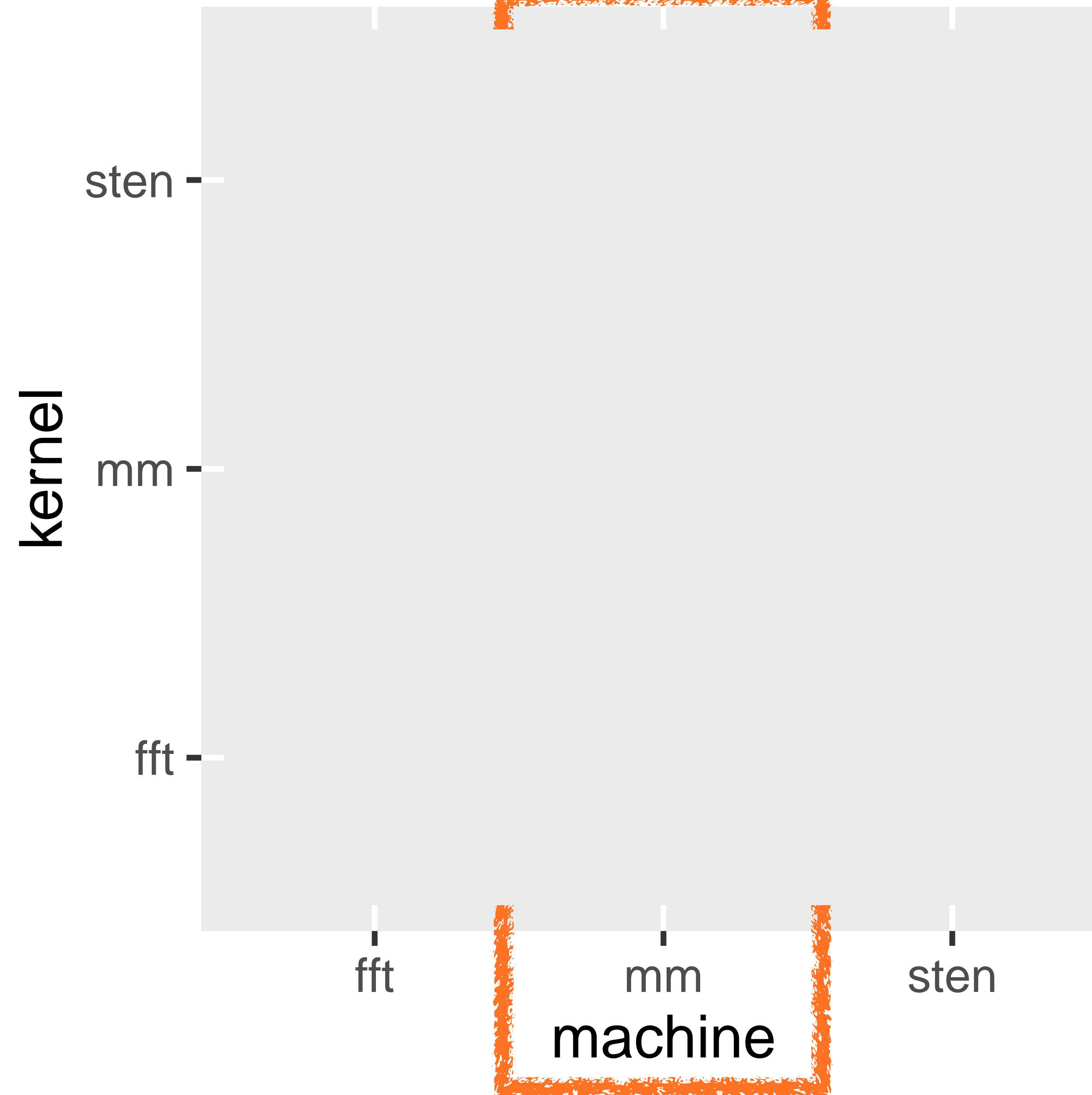
3D FFT vs. “Stencil” machines

45



Relative time (slowdown)

46



Can you beat the standard model?

- **What is your computational model?**
- Can I teach it to a freshman?
("Conte" question)
- Is it productive?
- Is it "easy" to provide provably guarantees for your model?
- **Can your model easily accommodate a variety of physical costs?**

These slides: bit.ly/4elibtM

