

Algorithmic time, energy, and power trade-offs in graph computations (?)

Sara Karamati · Jeff Young · Jee Choi [IBM Research] · [Richard \(Rich\) Vuduc](#)
· Oded Green [*Bader lab @ GT*] · Marat Dukhan · Anita Zakrzewska [*Bader lab @ GT*]

September 7, 2015 — 11th Int'l. Conf. Parallel Processing and Applied Mathematics (PPAM) — <http://ppam.pl>

**Kraków
bound!**



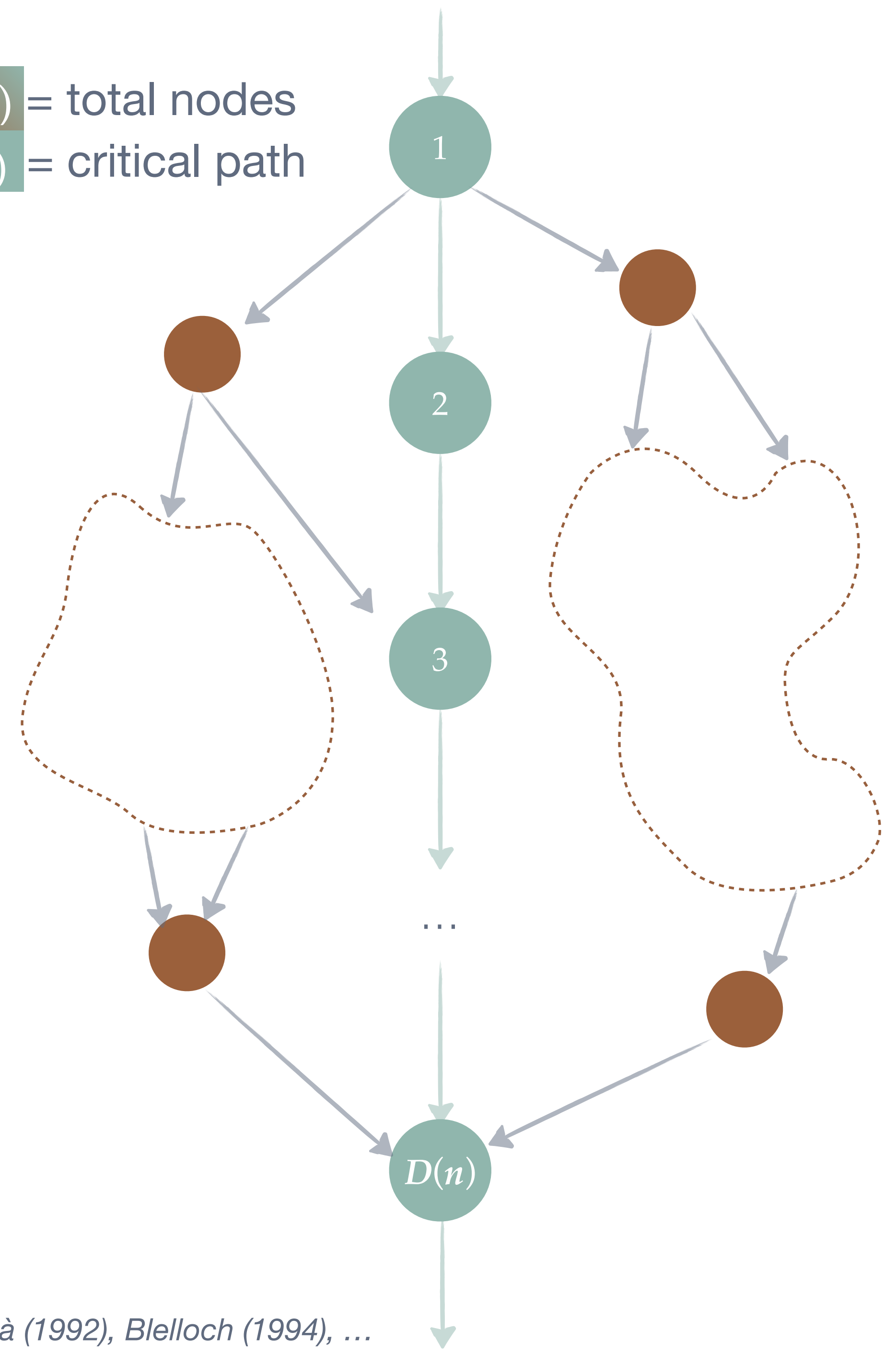
Do we need **new principles of algorithm design** when optimizing **energy** or **power** instead of **time** (or storage)?

Do we need **new principles of algorithm design** when optimizing **energy** or **power** instead of **time** (or storage)?

My answer: **No** — We have the main principles, but they are *even more important* when the metric changes, especially to energy instead of time.

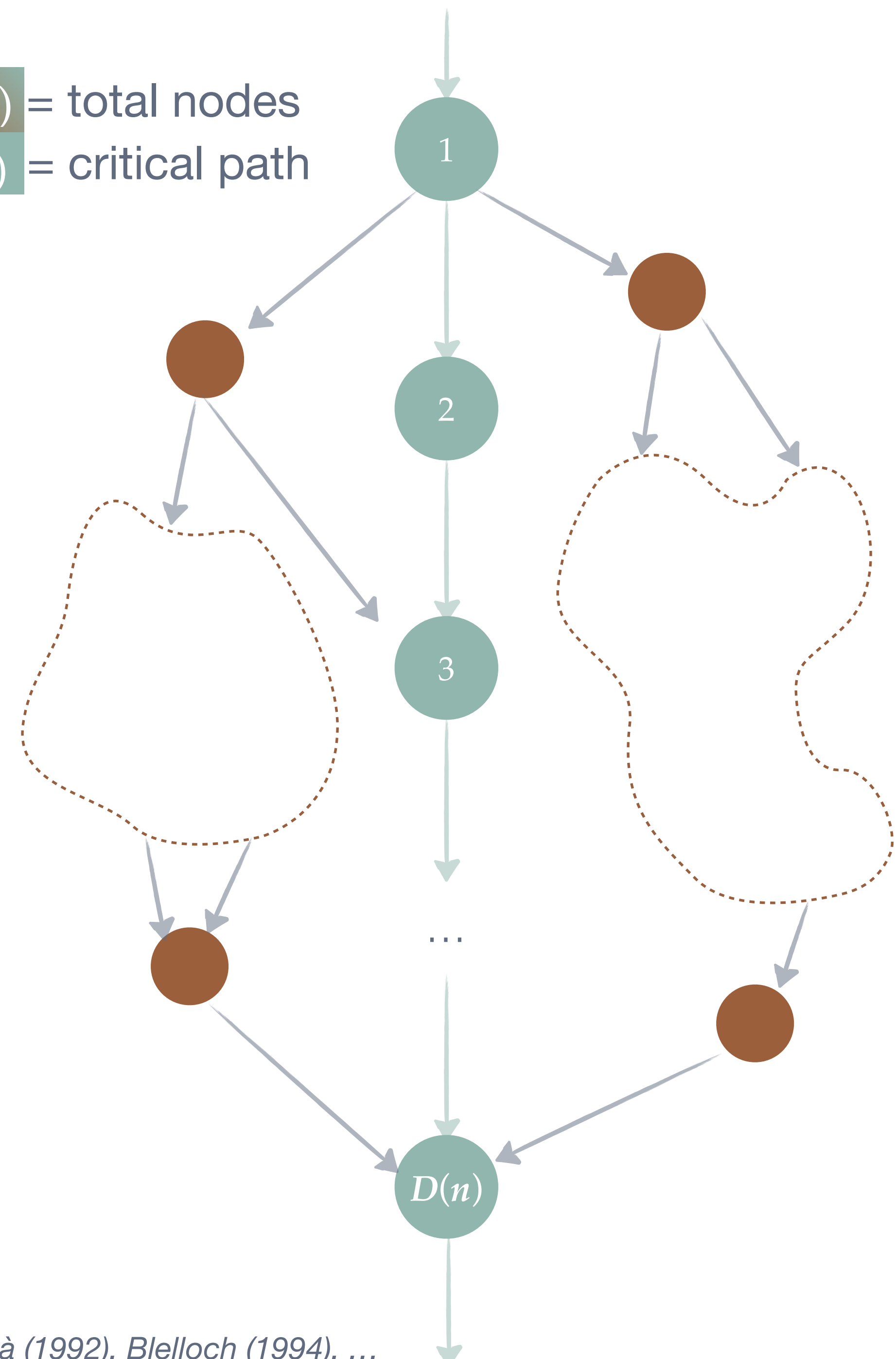
I. Basic abstract models

$W(n)$ = total nodes
 $D(n)$ = critical path



E.g., JàJà (1992), Blelloch (1994), ...

$W(n)$ = total nodes
 $D(n)$ = critical path



E.g., JàJà (1992), Blelloch (1994), ...

+

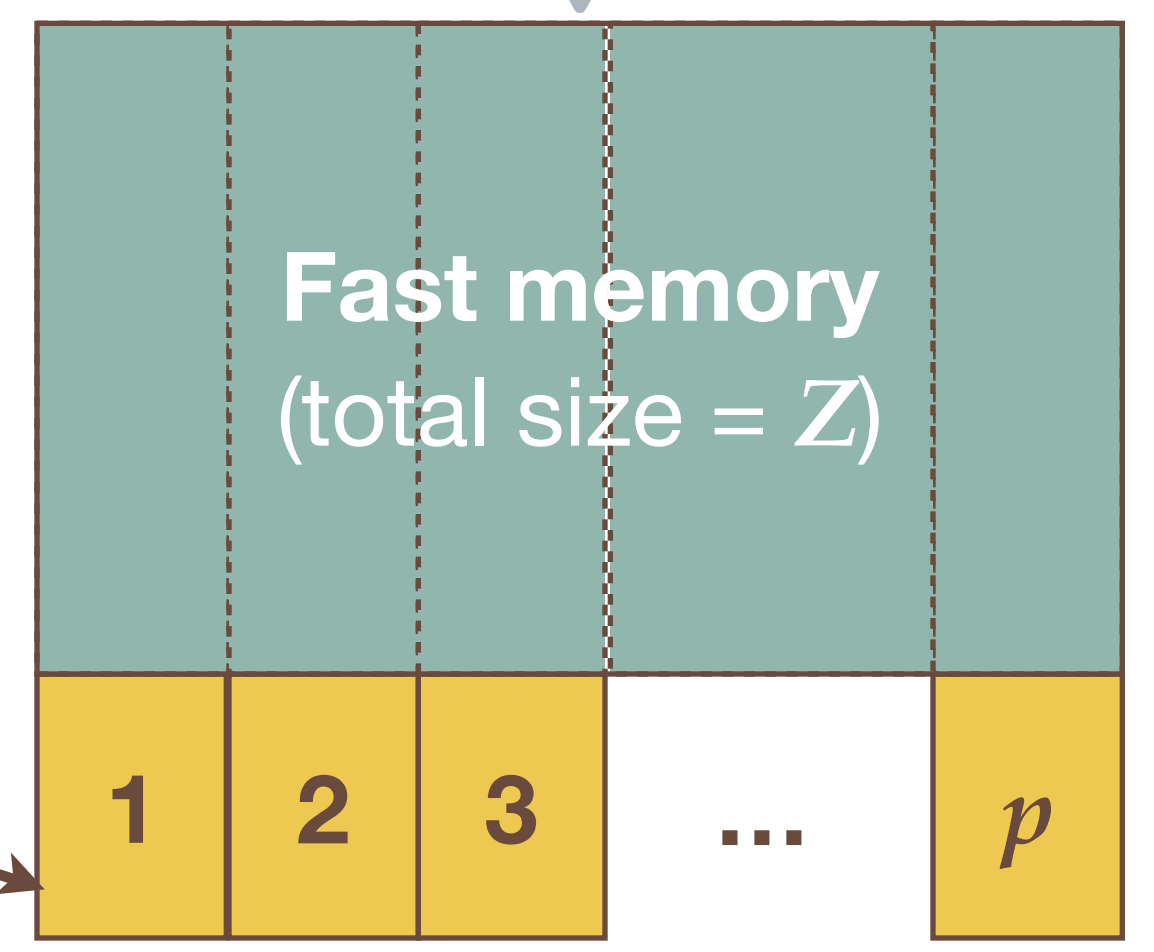
$Q(n; Z, L)$ = words transferred

E.g., Agarwal and Vitter (1988), ...



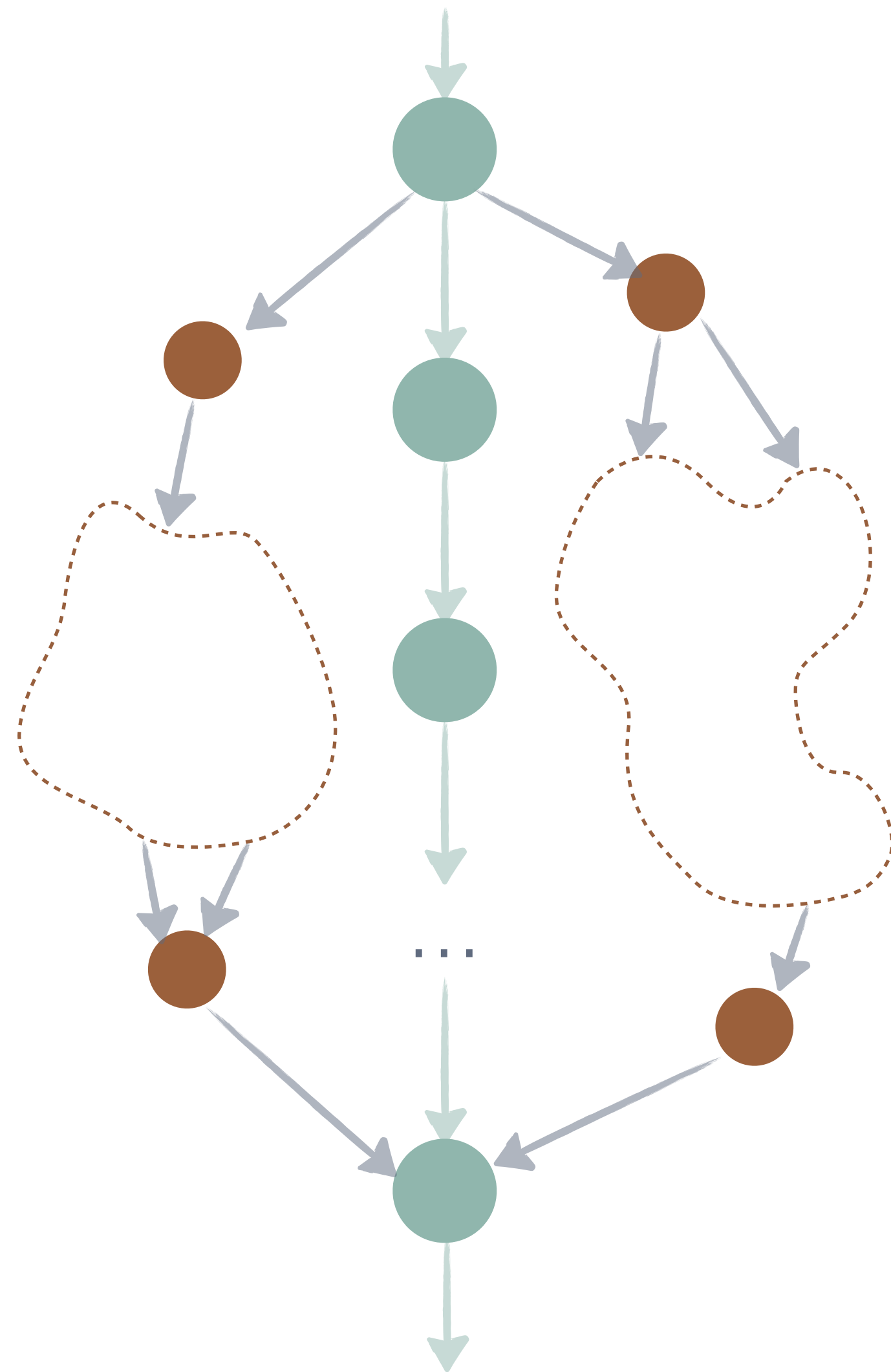
α latency
 β bandwidth

L words per transaction



C_0 op/s per core

Example: Work-span (depth) model



$W(n) = \text{work (total ops)}$

$D(n) = \text{span (critical path)}$

$W(n) / D(n)$

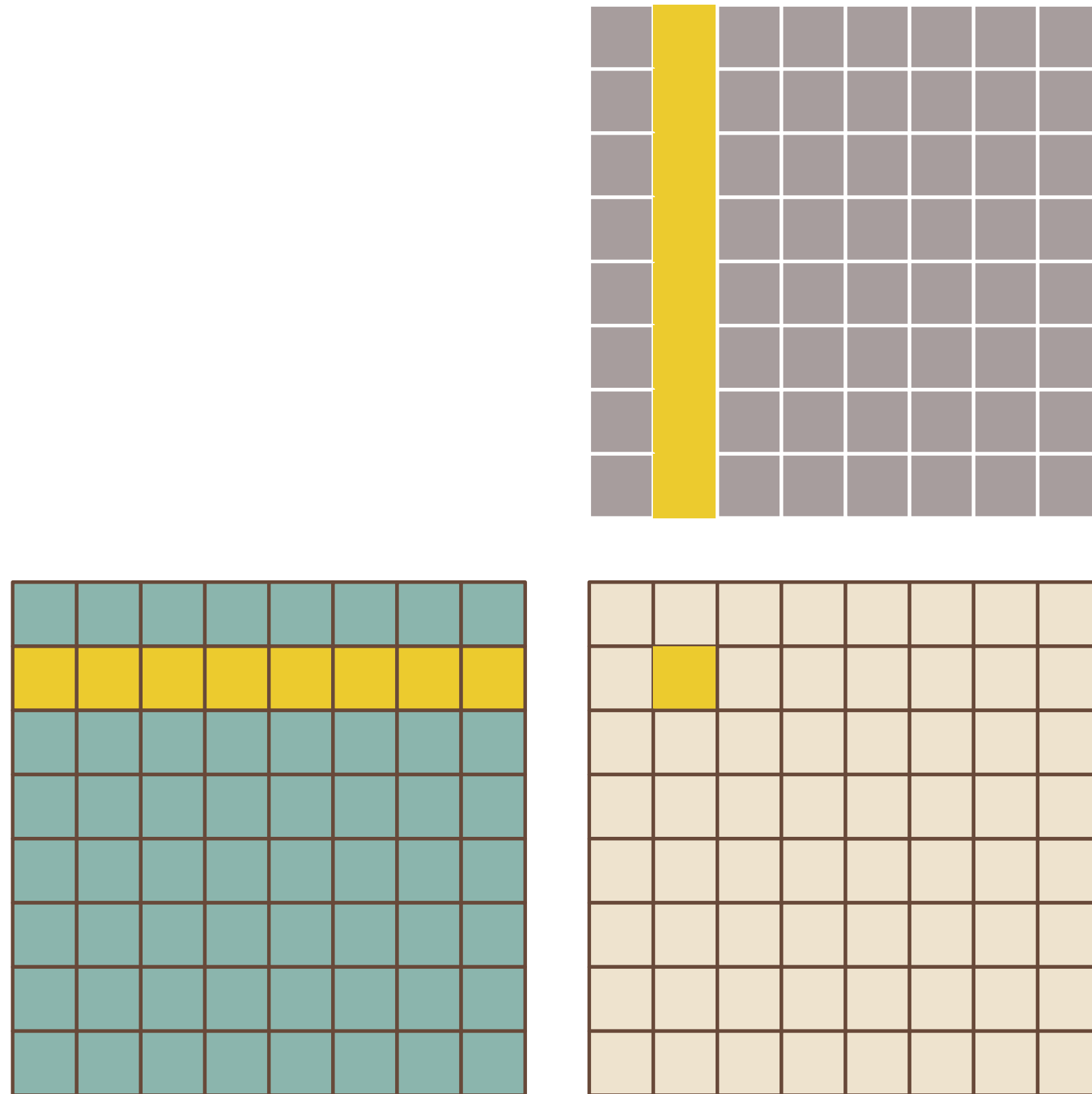
= **inherent parallelism**

Desiderata:

Work optimality

Maximal parallelism

Example: Matrix multiply
(non-Strassen)



$$\frac{W(n)}{D(n)} = \Omega \left(\frac{n^3}{\log^k n} \right)$$

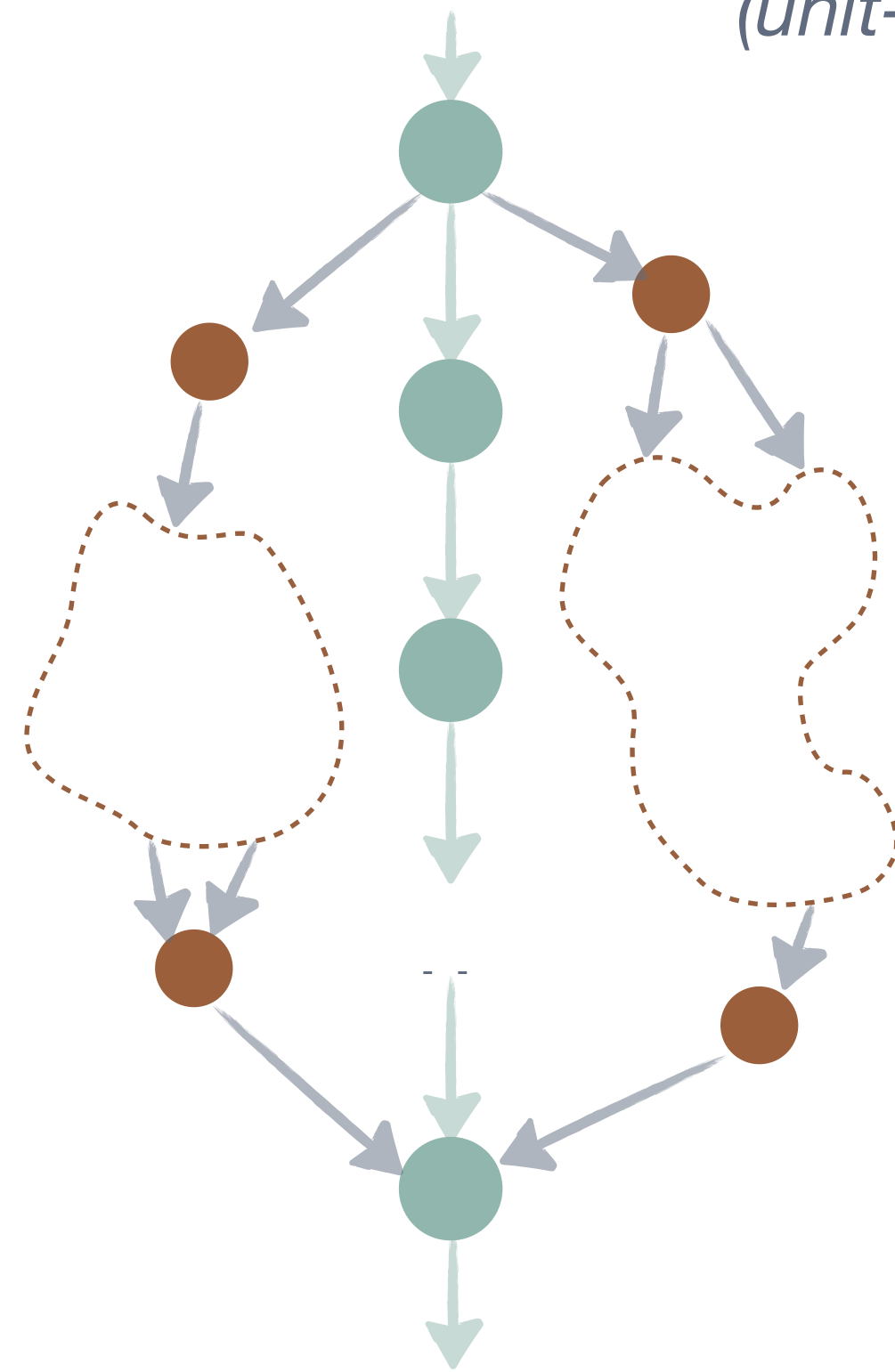
Parallelism

$$\boxed{C} \leftarrow \boxed{C} + \boxed{A} * \boxed{B}$$

II. Moving from abstract to concrete (physical) costs

How much time to execute a DAG on P processors?

(unit-cost operations)



$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

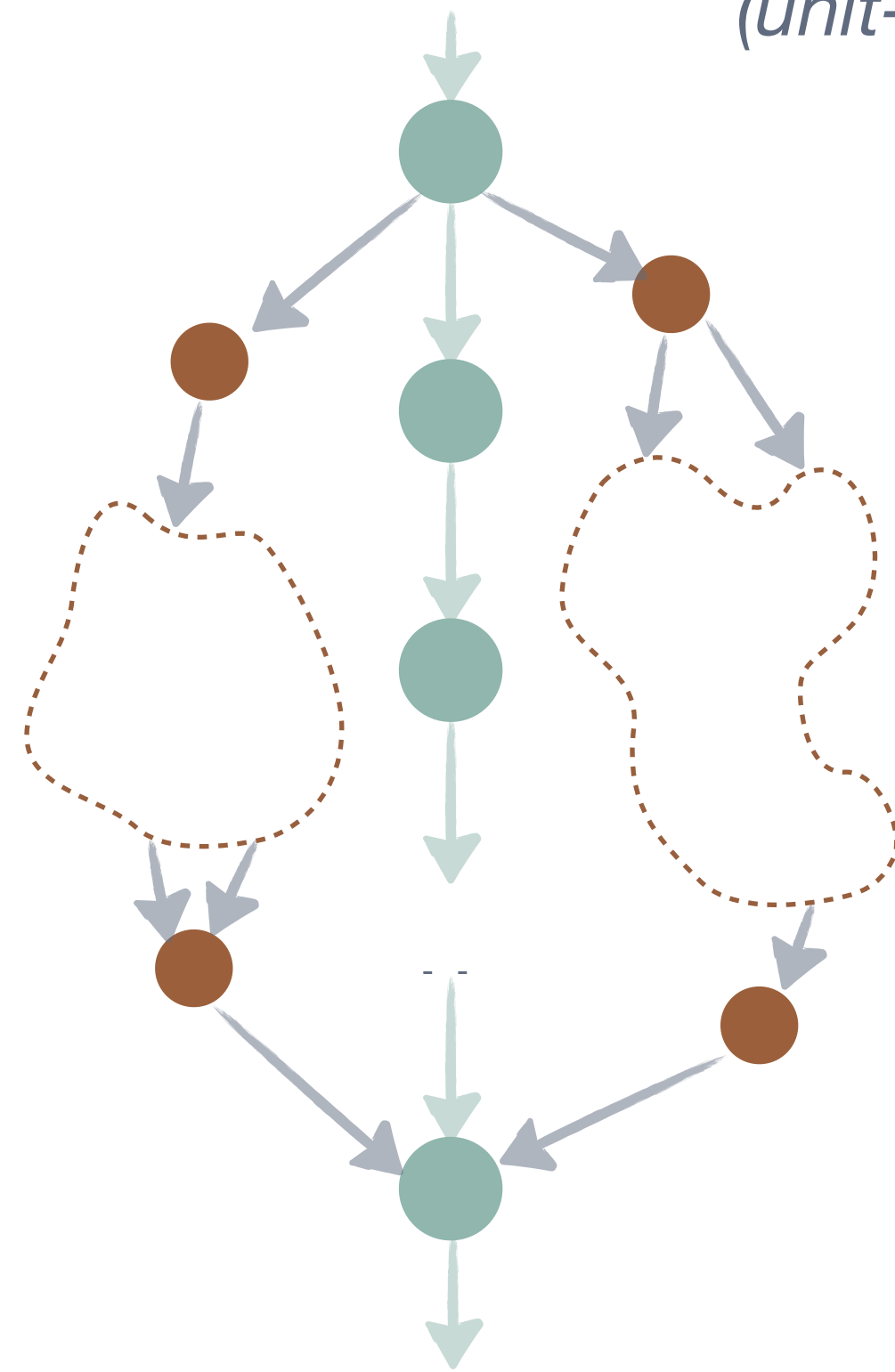
$W(n) / D(n)$

= *inherent parallelism*

$$T(n; P) = ?$$

How much time to execute a DAG on P processors?

(unit-cost operations)



$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

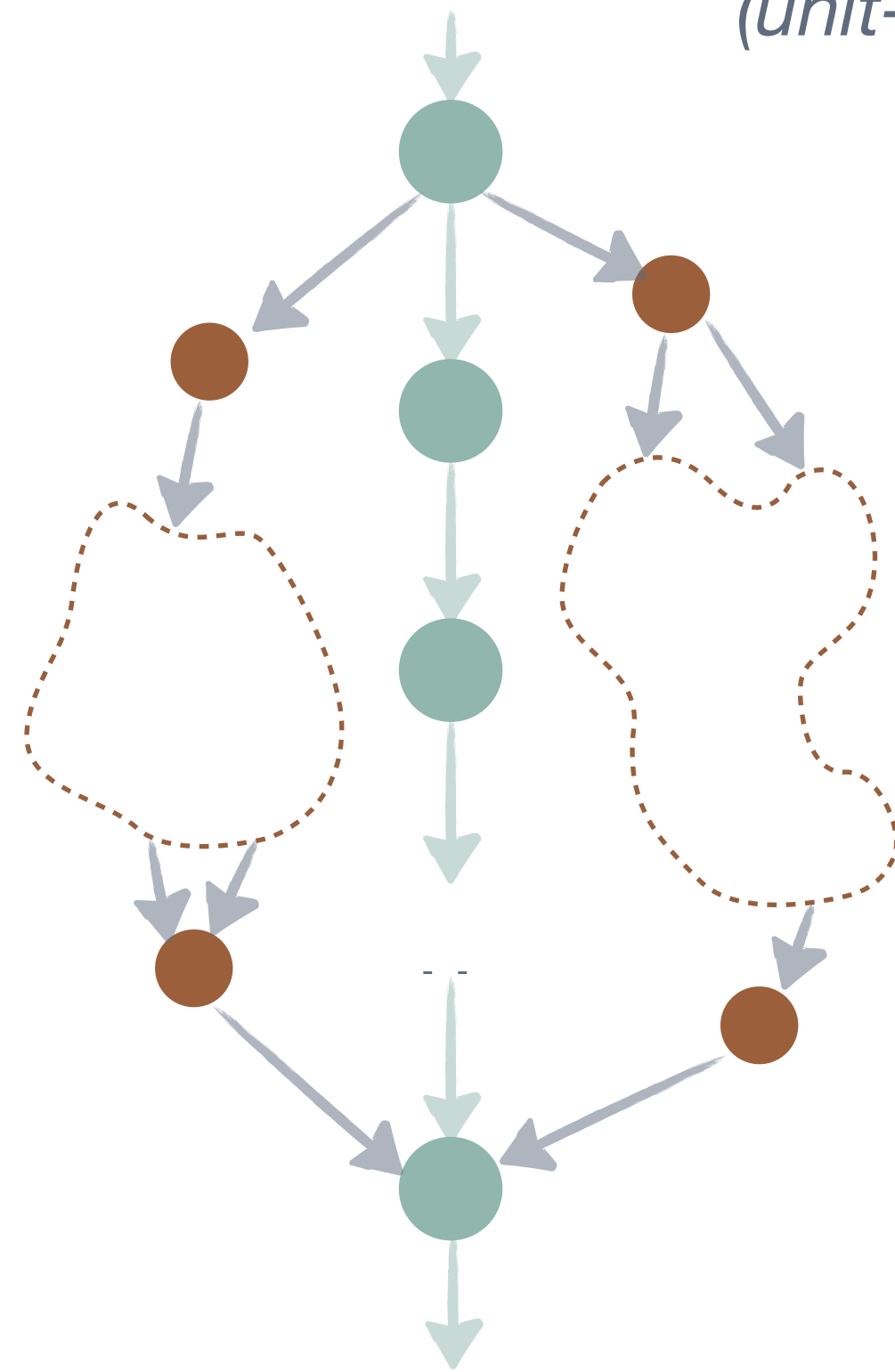
$W(n) / D(n)$

= ***inherent parallelism***

$$T(n; P) \geq \max \left\{ \frac{W(n)}{P}, D(n) \right\}$$

How much time to execute a DAG on P processors?

(unit-cost operations)



$W(n)$ = work (total ops)

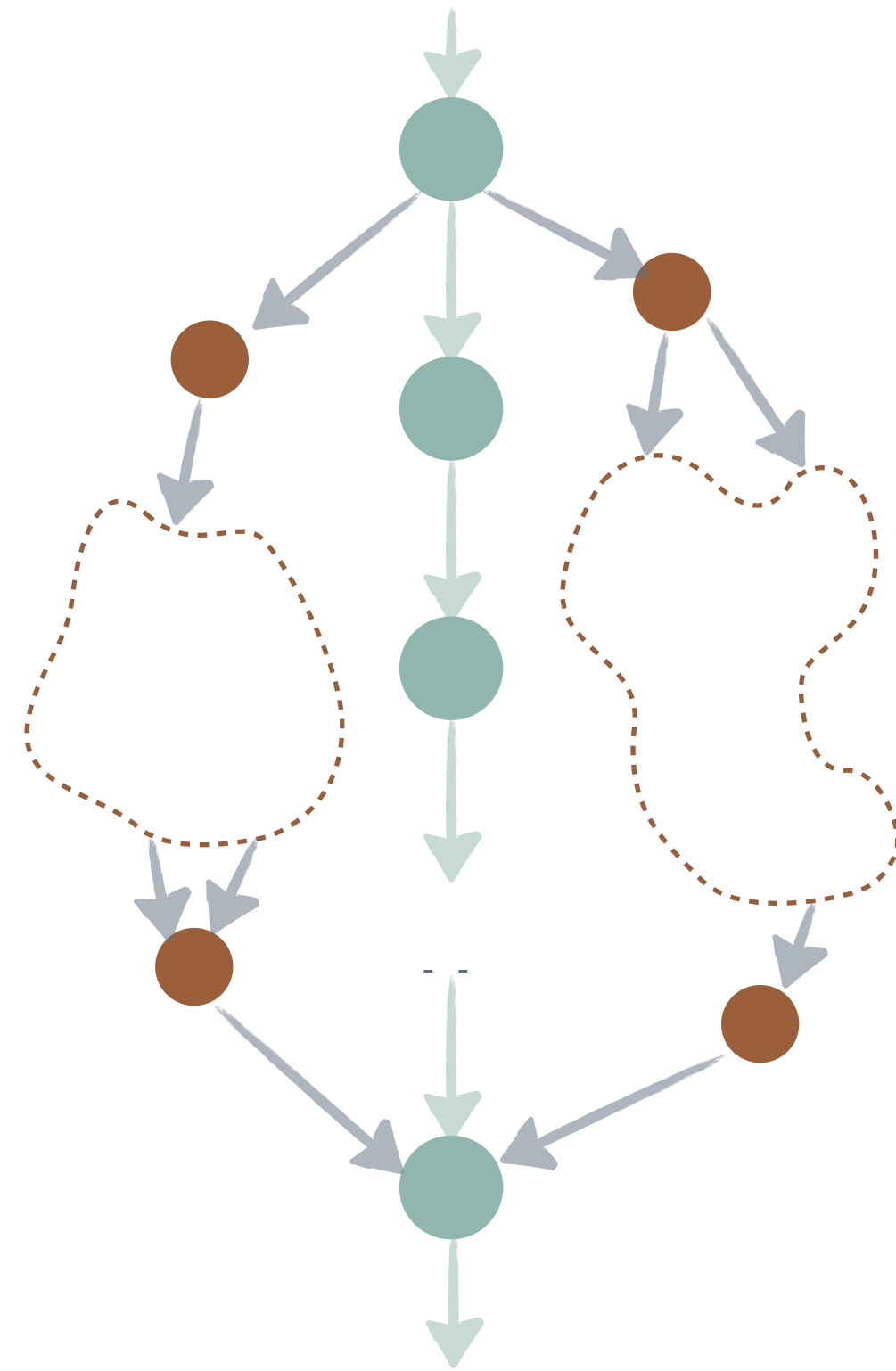
$D(n)$ = span (critical path)

$W(n) / D(n)$

= **inherent parallelism**

$$T(n; P) \leq D(n) + \frac{W(n) - D(n)}{P}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

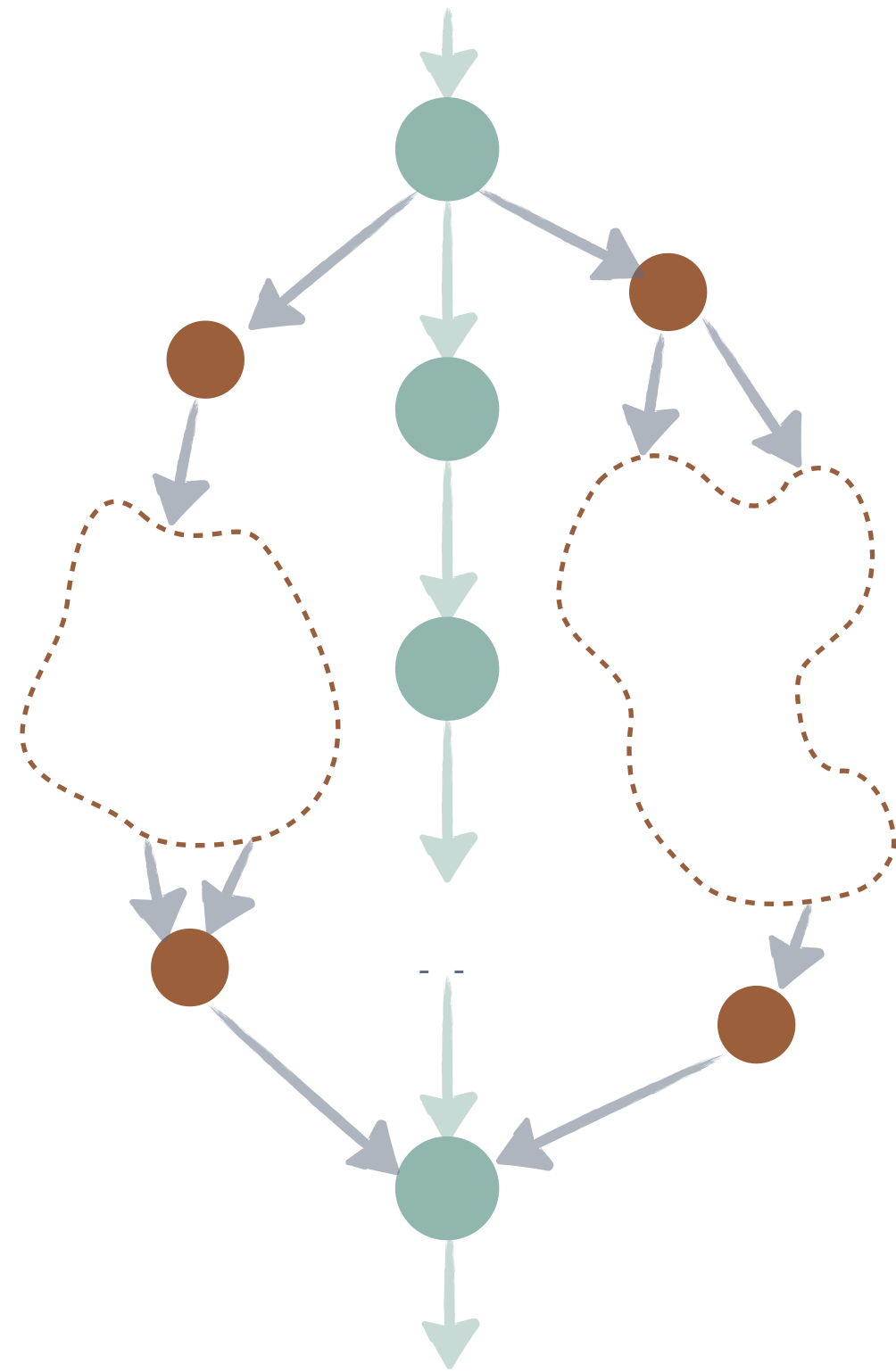
$D(n)$ = span (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$S_*(n; P) \equiv \frac{T_*(n)}{T(n; P)}$$

What is the speedup over the *best* sequential algorithm?



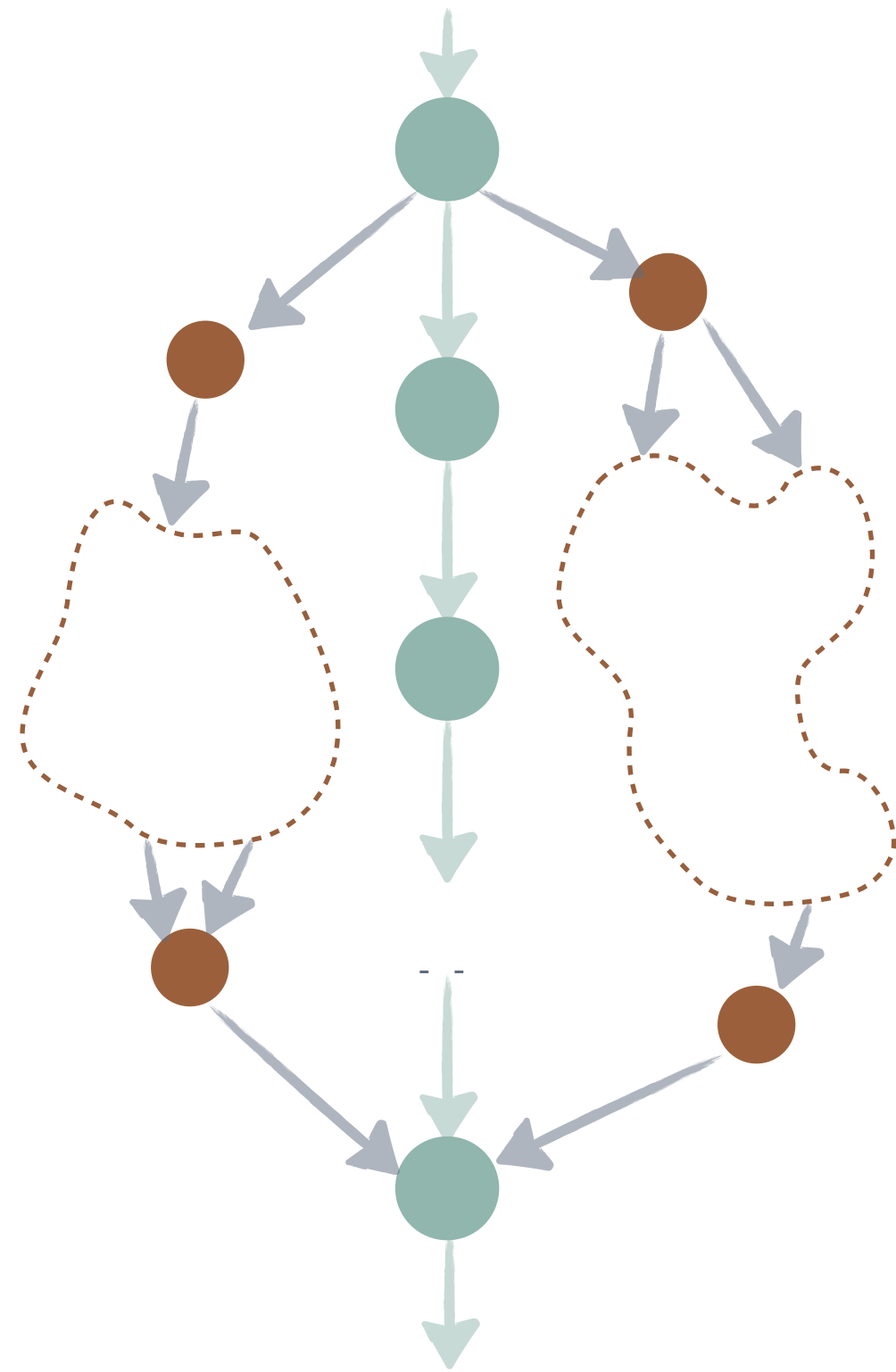
$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

$W(n) / D(n)$
= *inherent parallelism*

$$S_*(n; P) \leq \frac{W_*(n)}{\max \left\{ \frac{W(n)}{P}, D(n) \right\}}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

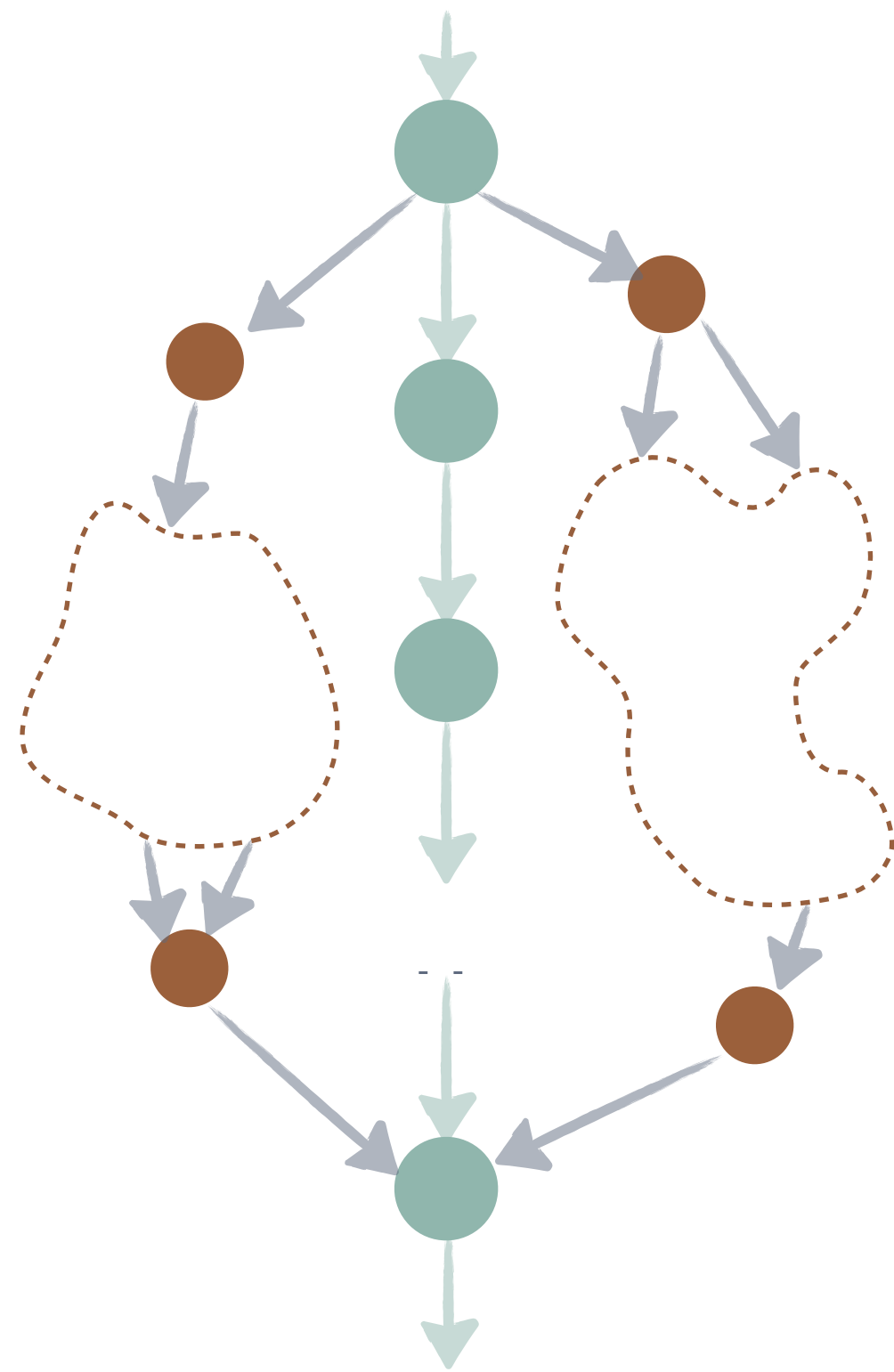
$D(n)$ = span (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$S_*(n; P) \leq P \cdot \min \left\{ \frac{W_*(n)}{W(n)}, \frac{W_*(n)/D(n)}{P} \right\}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

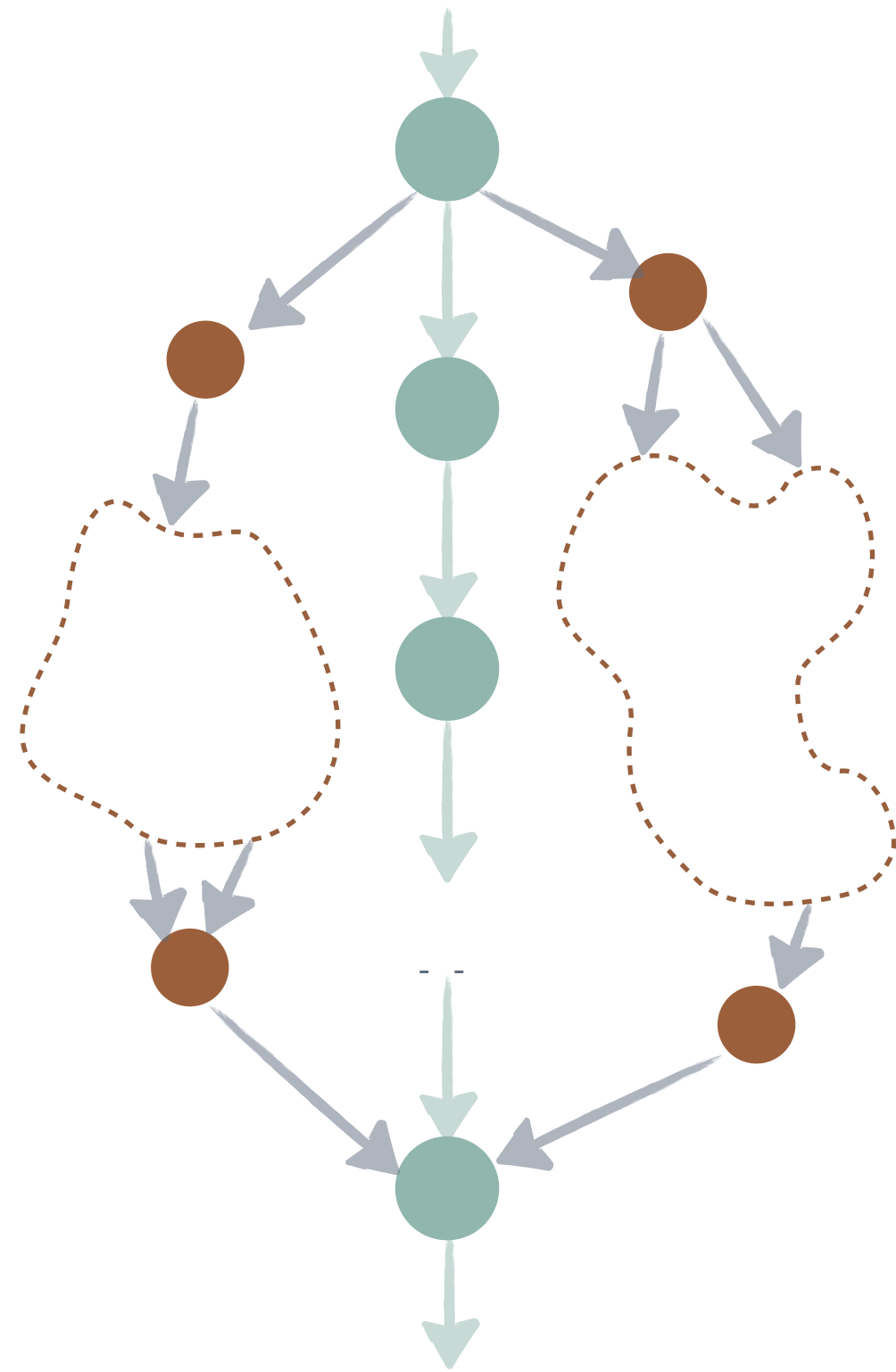
$W(n) / D(n)$

= *inherent parallelism*

Ideal (linear)

$$S_*(n; P) \leq \min \left\{ \frac{W_*(n)}{W(n)}, \frac{W_*(n)/D(n)}{P} \right\}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

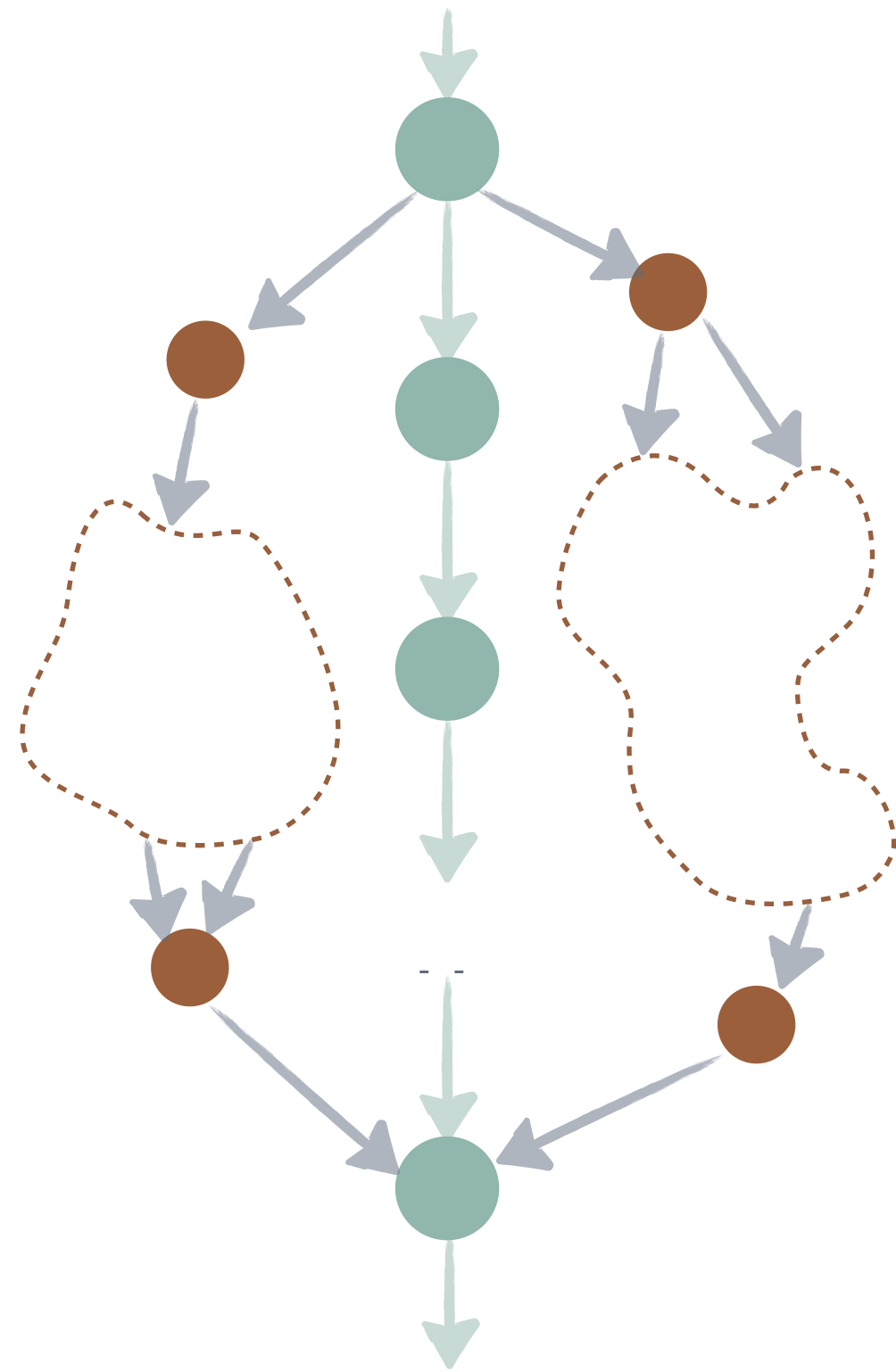
$D(n)$ = span (critical path)

$W(n) / D(n)$
= *inherent parallelism*

Work-optimality

$$S_*(n; P) \leq P \cdot \min \left\{ \frac{W_*(n)}{W(n)}, \frac{W_*(n)/D(n)}{P} \right\}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

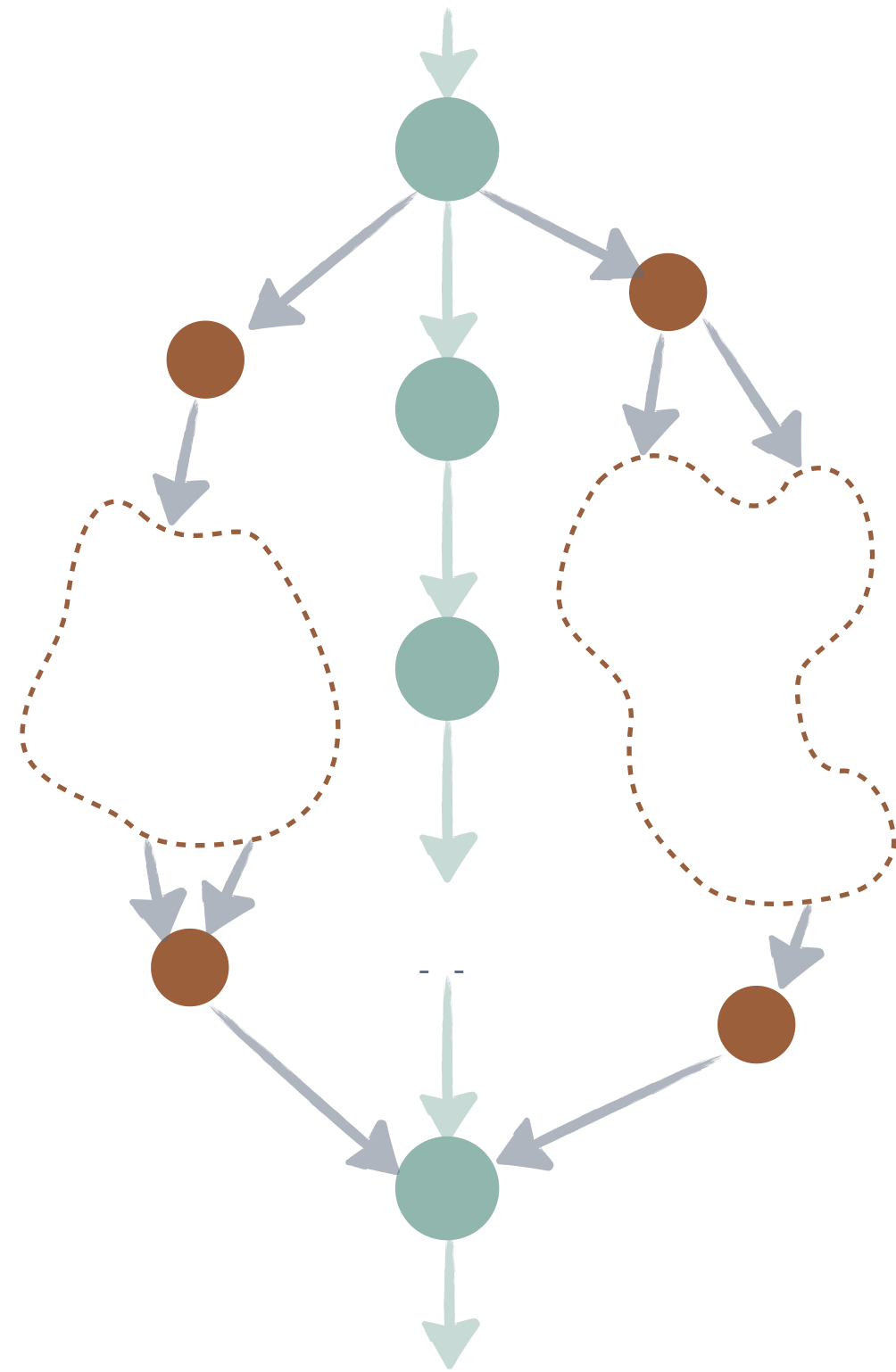
$W(n) / D(n)$

= *inherent parallelism*

Weak scalability

$$S_*(n; P) \leq P \cdot \min \left\{ \frac{W_*(n)}{W(n)}, \frac{W_*(n)/D(n)}{P} \right\}$$

What is the speedup over the *best* sequential algorithm?



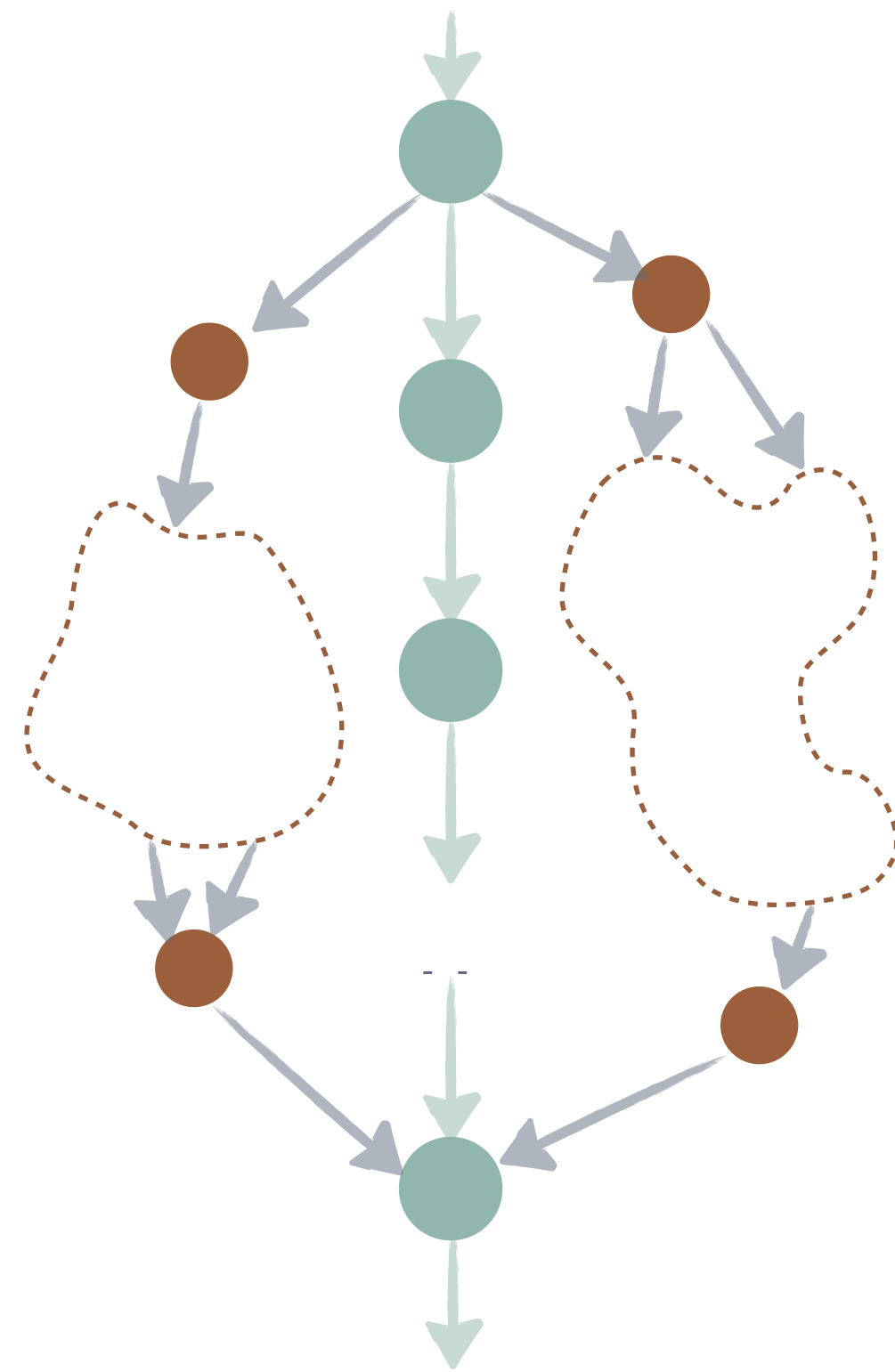
$W(n)$ = work (total ops)

$D(n)$ = span (critical path)

$W(n) / D(n)$
= *inherent parallelism*

$$S_*(n; P) \geq \frac{W_*(n)}{D(n) + \frac{W(n) - D(n)}{P}}$$

What is the speedup over the *best* sequential algorithm?



$W(n)$ = work (total ops)

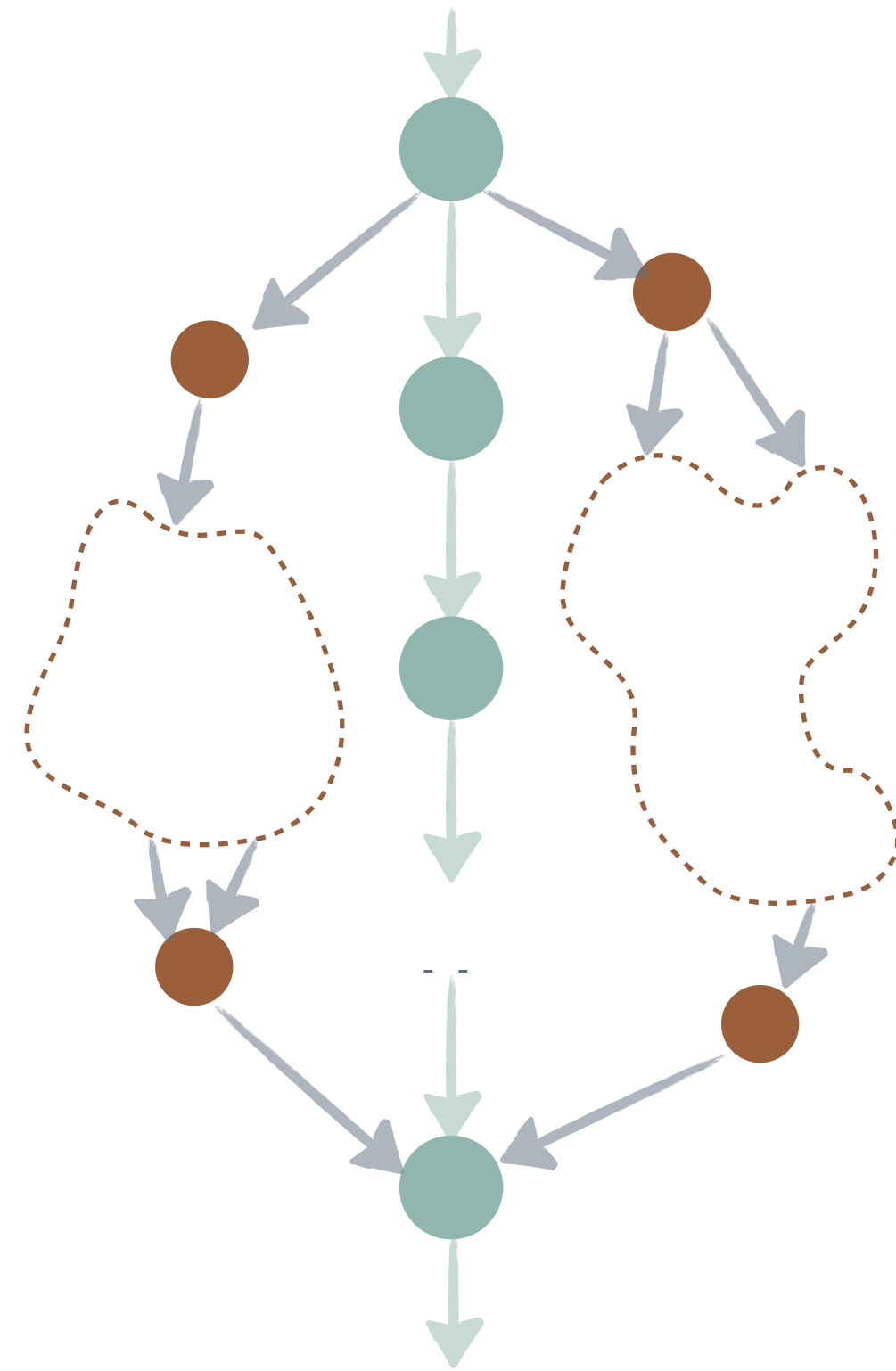
$D(n)$ = span (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$S_*(n; P) \geq \frac{P}{\frac{P-1}{W_*(n)/D(n)} + \frac{W(n)}{W_*(n)}}$$

How much **energy** is required?



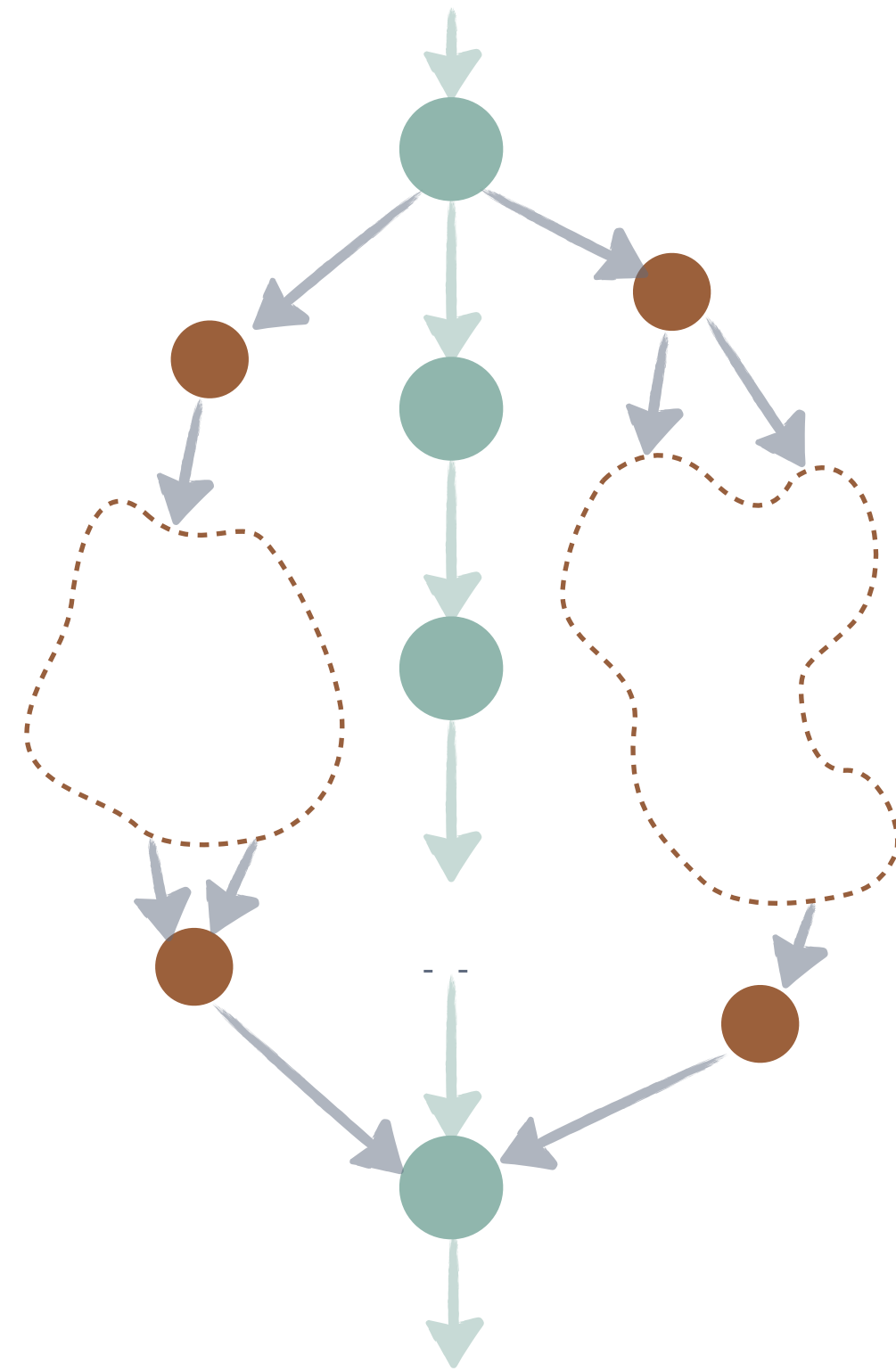
$W(n)$ = *work* (total ops)

$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

How much **energy** is required?



$W(n)$ = *work* (total ops)

$D(n)$ = *span* (critical path)

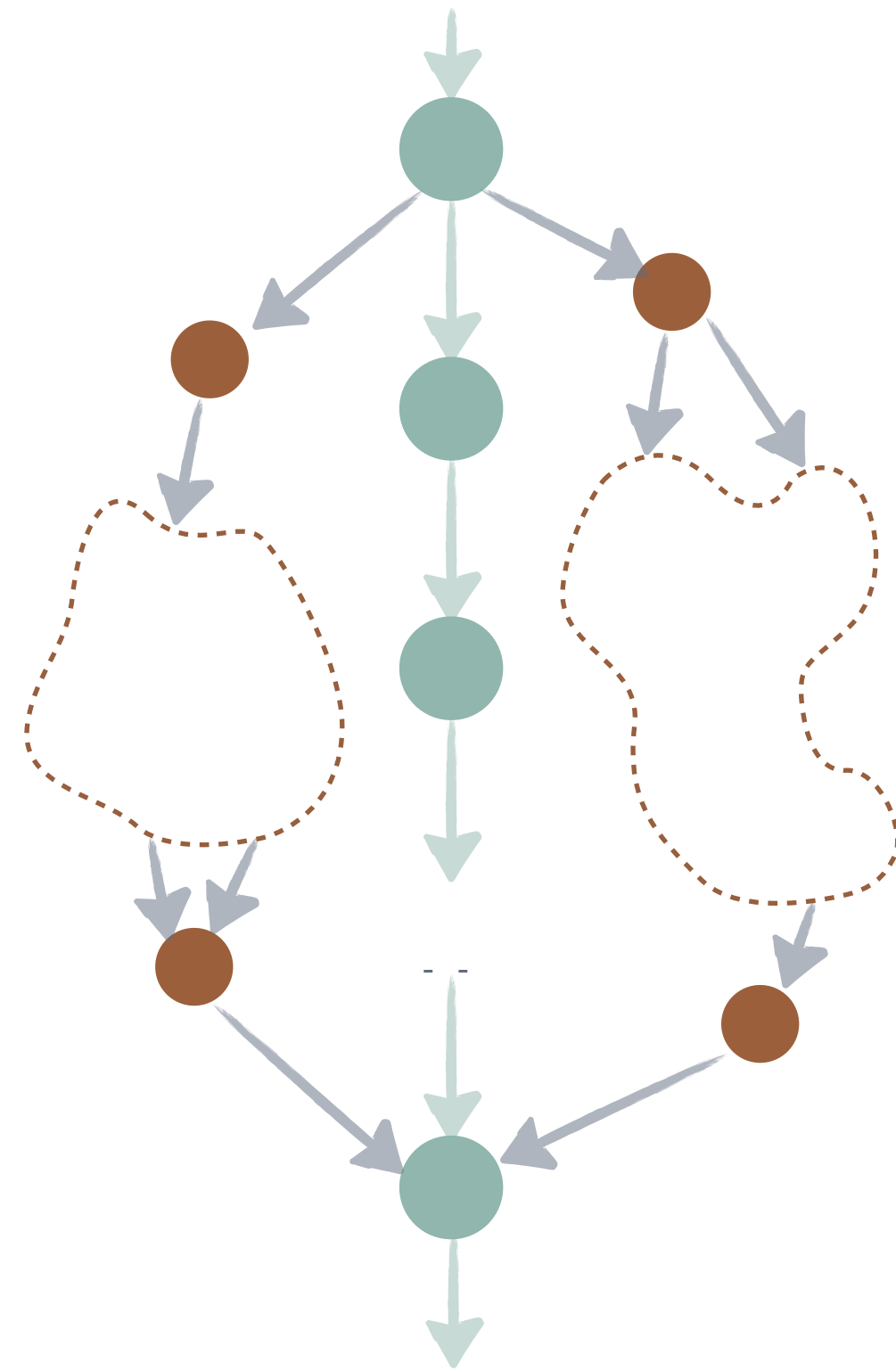
$W(n) / D(n)$

= *inherent parallelism*

Time is a cost you *may hide* by overlap, e.g., parallelism.

Energy is a cost you *must pay* for every operation.

How much **energy** is required?



$W(n)$ = *work* (total ops)

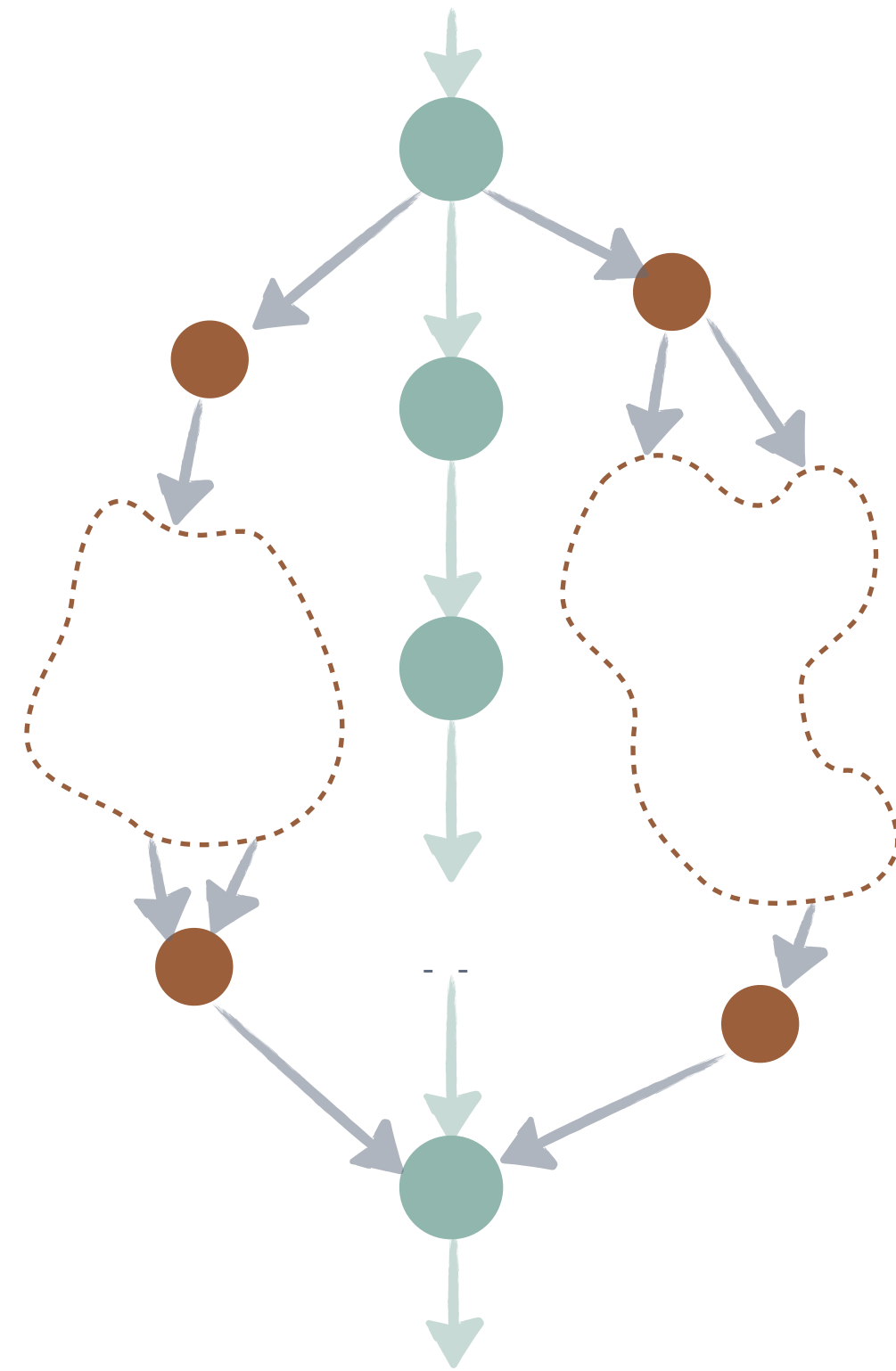
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

Energy \propto **Work**

How much **energy** is required?



$W(n)$ = *work* (total ops)

$D(n)$ = *span* (critical path)

$W(n) / D(n)$

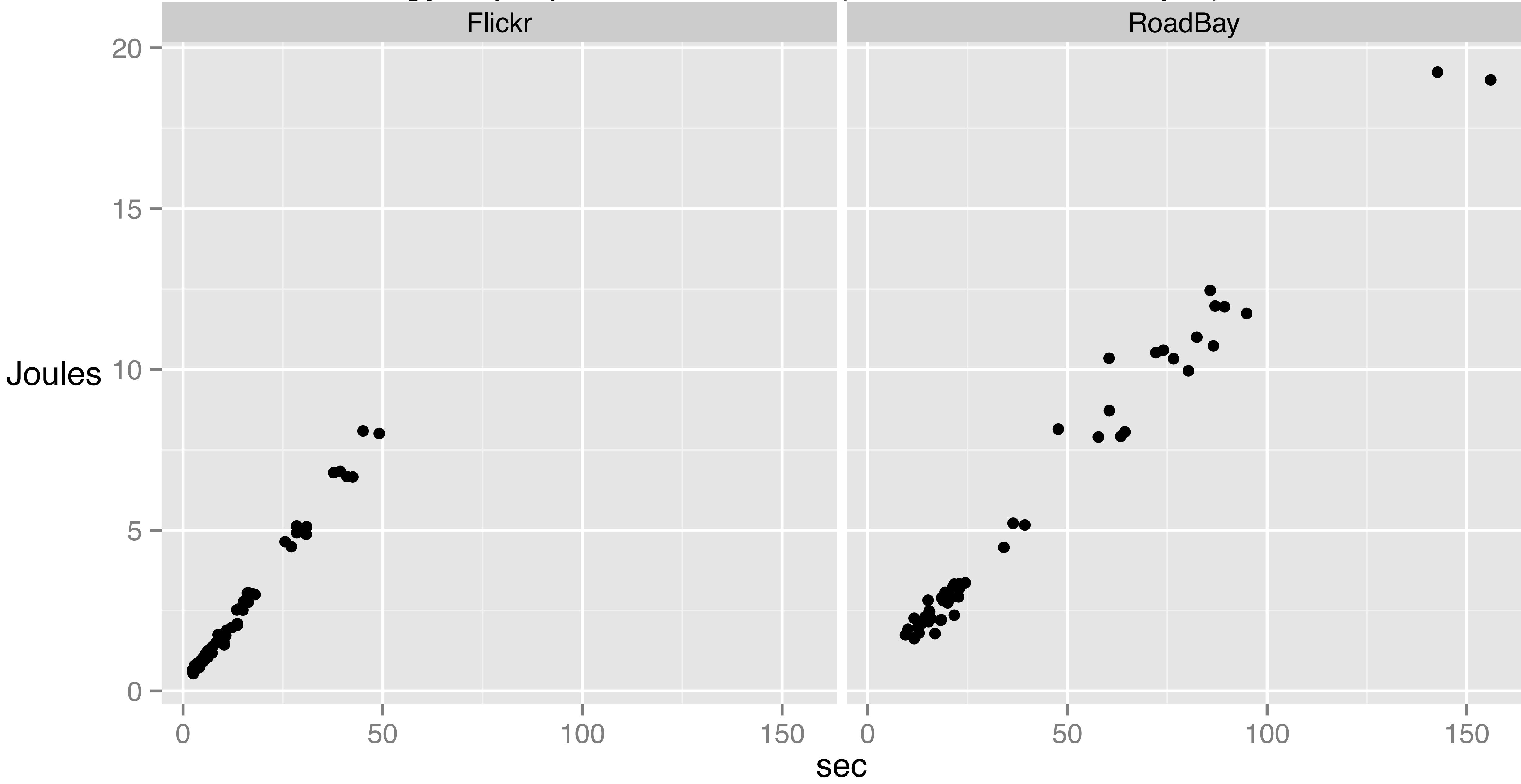
= *inherent parallelism*

Energy \propto **Work**

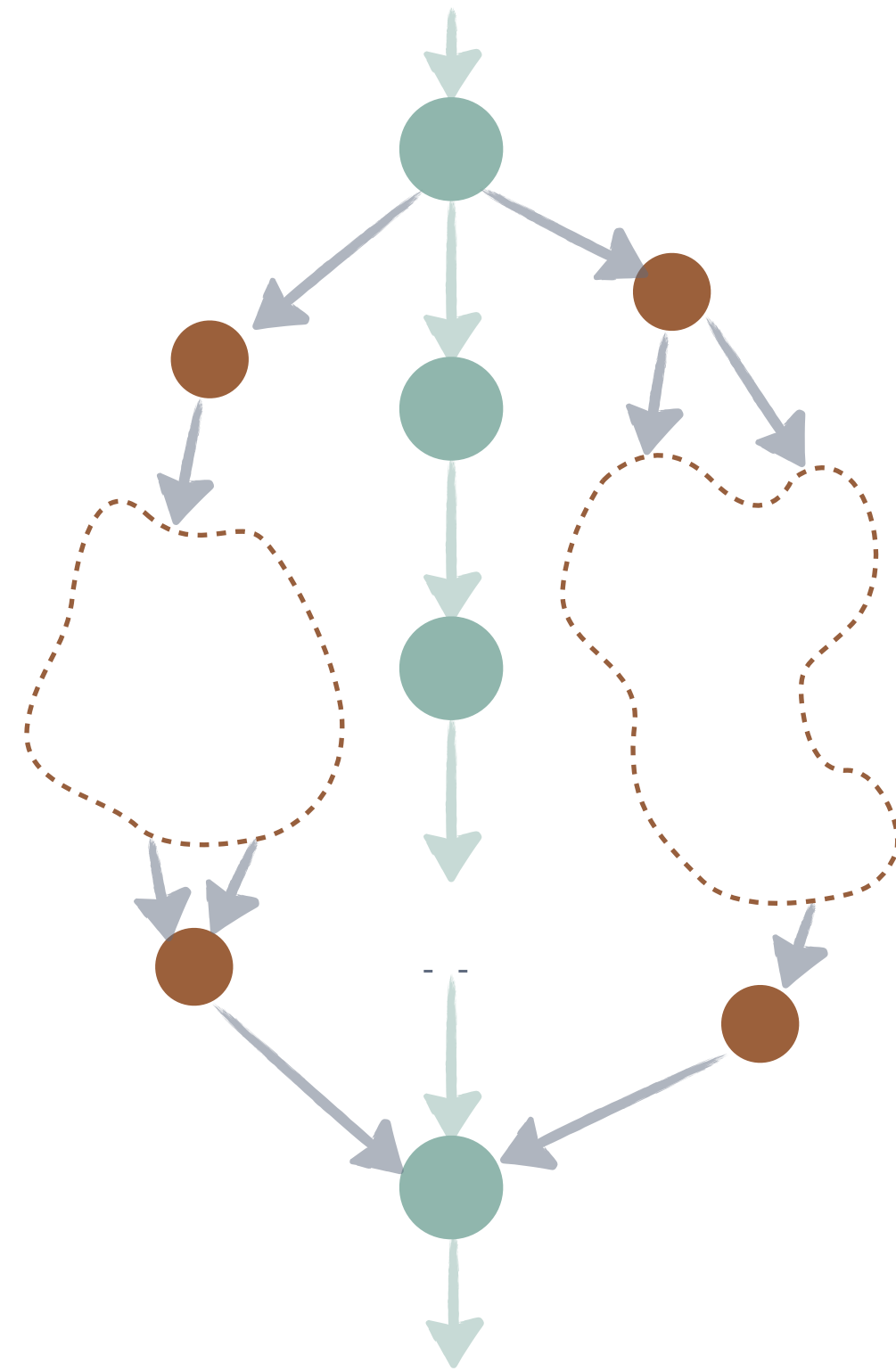
$$T(n; P) \sim \frac{W(n)}{P}$$

Time \propto **Energy**

Execution energy is proportional to time (SSSP+GPU example)



How much **power** is required?



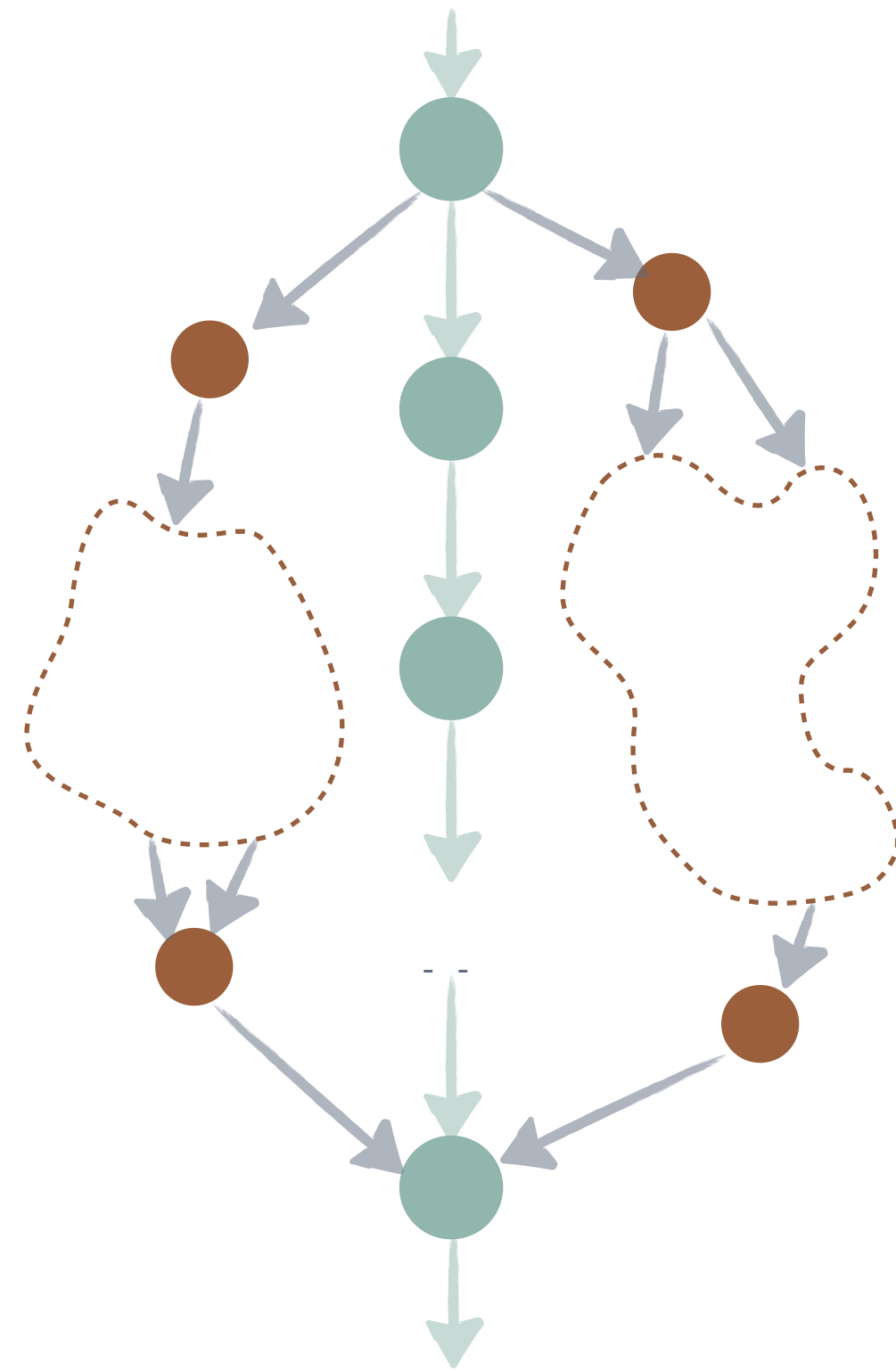
$W(n)$ = *work* (total ops)

$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

How much **power** is required?



$W(n)$ = *work* (total ops)

$D(n)$ = *span* (critical path)

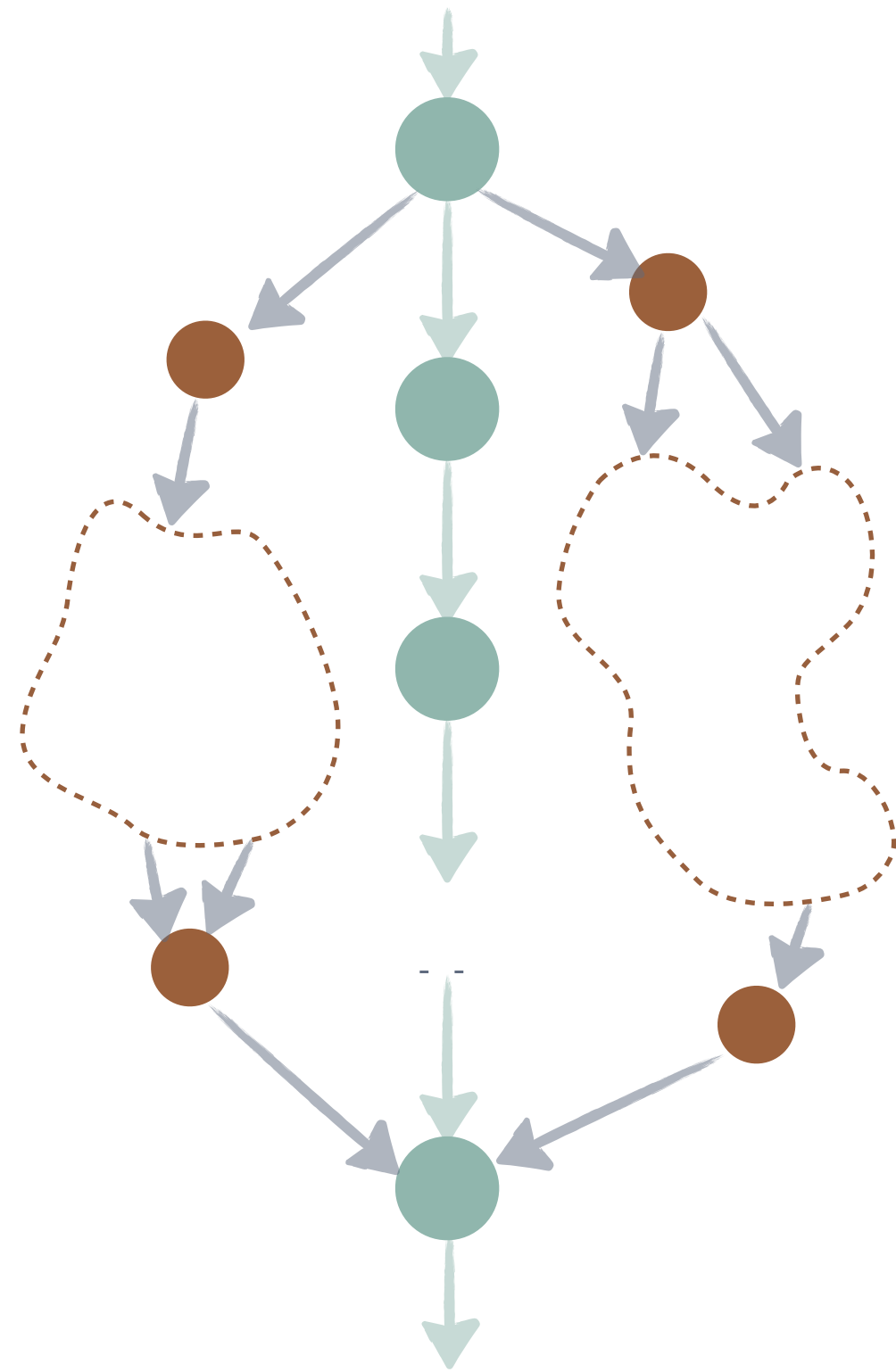
$W(n) / D(n)$

= *inherent parallelism*

Power \equiv **Energy** / **Time**

(average instantaneous)

How much **power** is required?



$W(n)$ = *work* (total ops)

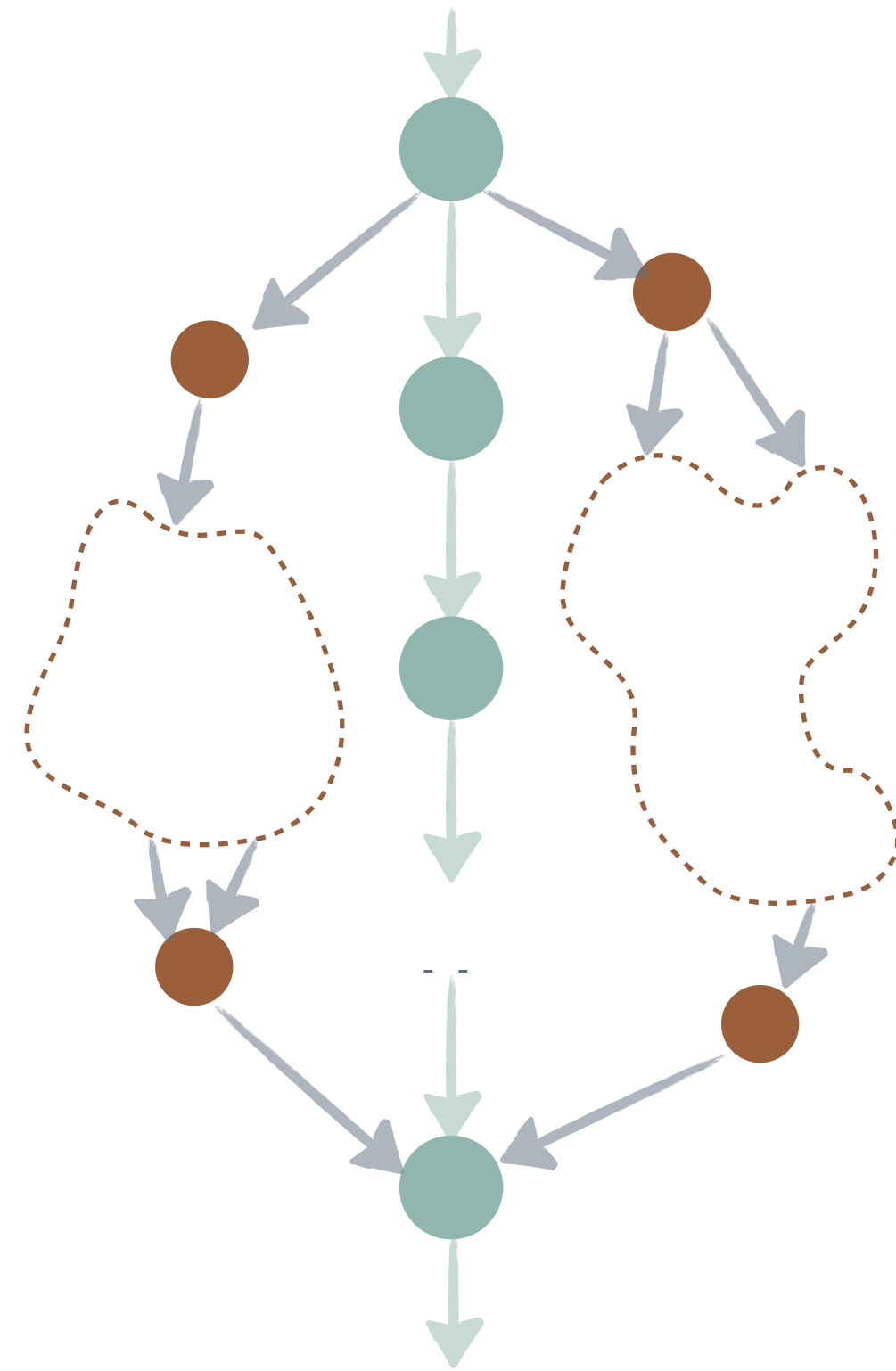
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$\Phi(n; P) \equiv \frac{\Theta(W(n))}{T(n; P)}$$

How much **power** is required?



$W(n)$ = *work* (total ops)

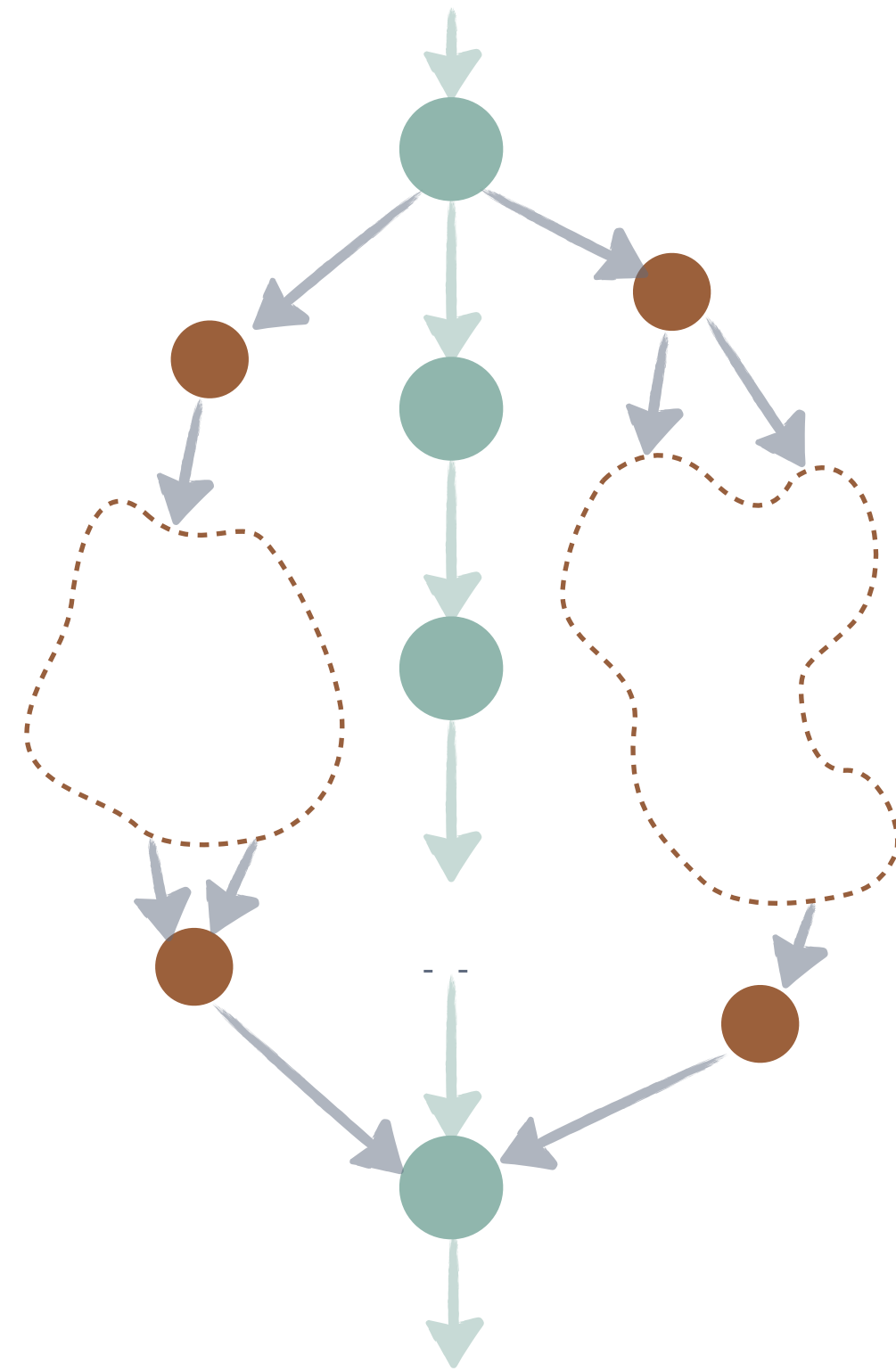
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$\Phi(n; P) = \frac{\Theta(T(n; P = 1))}{T(n; P)} = \Theta(S(n; P))$$

How much **power** is required?



$W(n)$ = *work* (total ops)

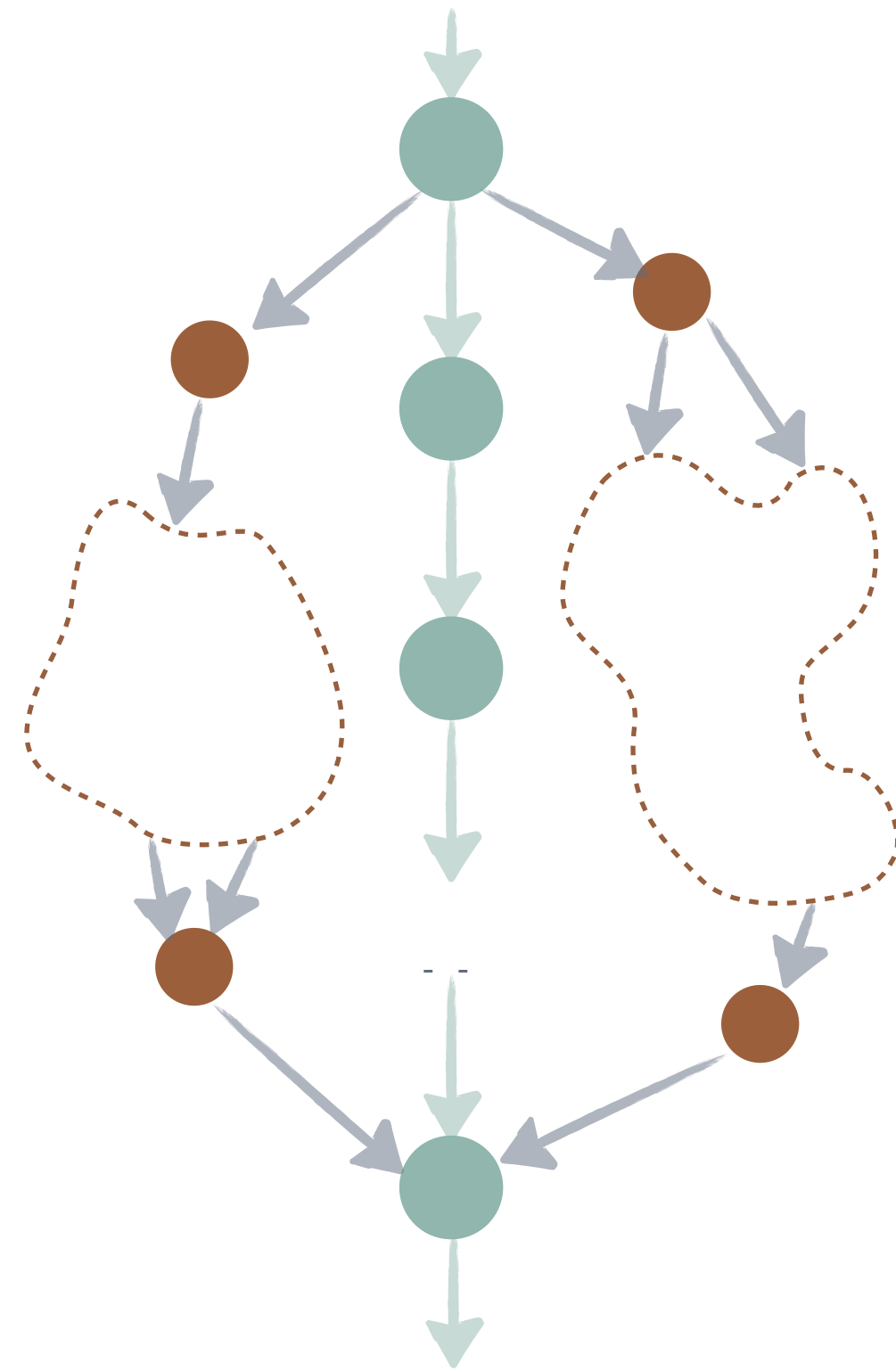
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$\Phi(n; P) = \frac{\Theta(T(n; P = 1))}{T(n; P)} = \Theta(S(n; P))$$

How much **power** is required?



$W(n)$ = *work* (total ops)

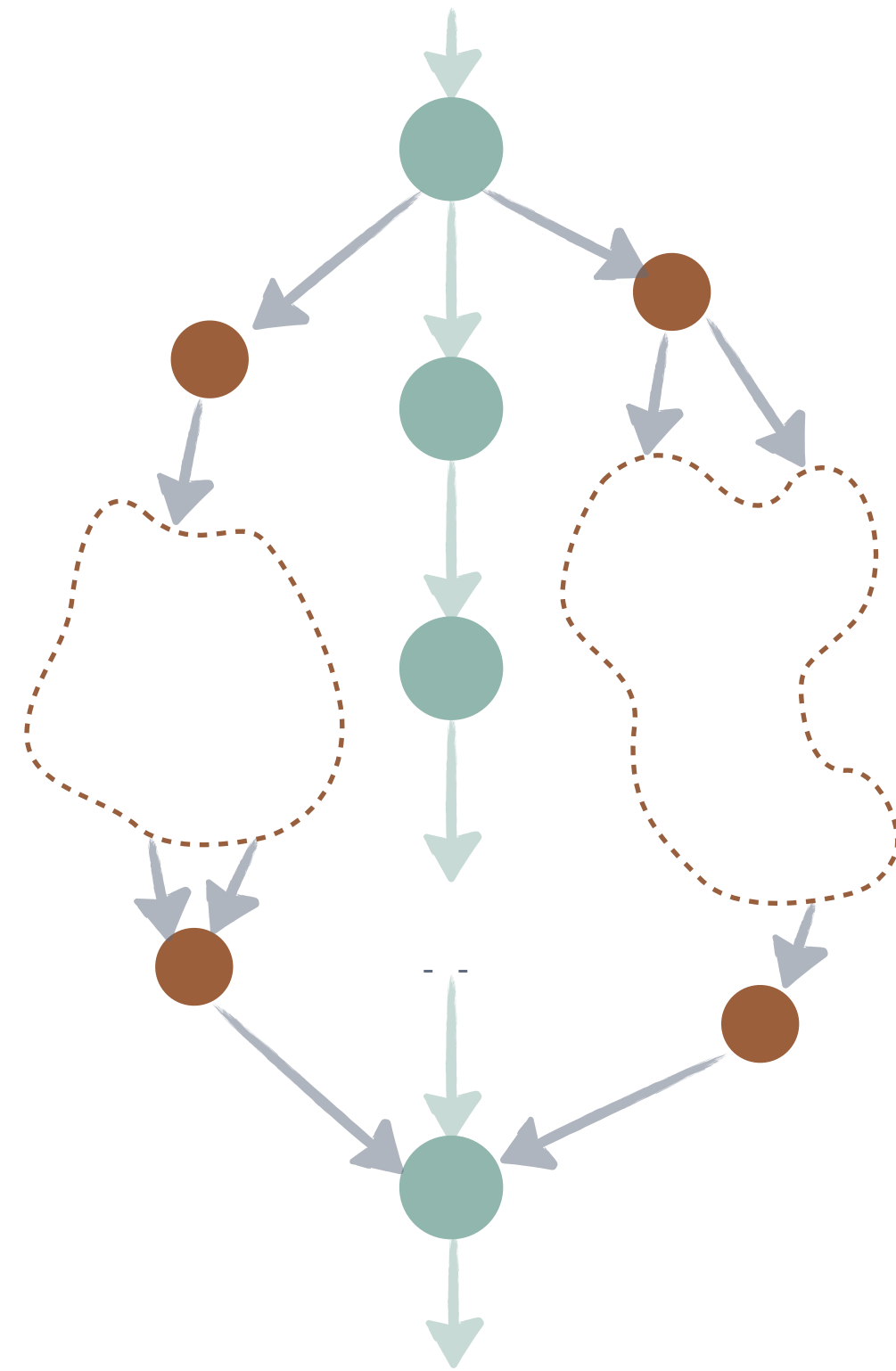
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

$$\Phi(n; P) = \frac{\Theta(T(n; P = 1))}{T(n; P)} = \Theta(S(n; P))$$

How much **power** is required?



$W(n)$ = *work* (total ops)

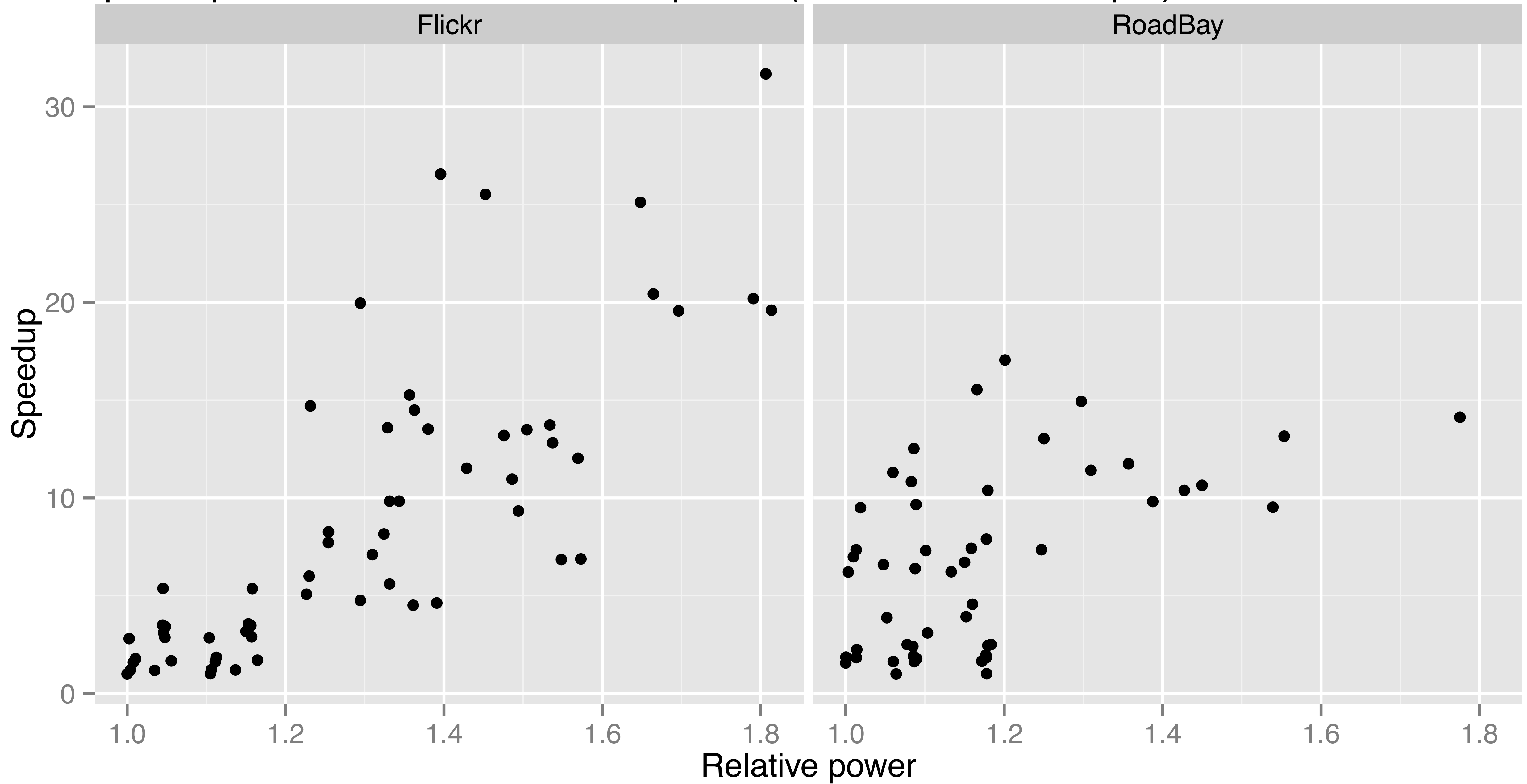
$D(n)$ = *span* (critical path)

$W(n) / D(n)$

= *inherent parallelism*

⇒ One expects **algorithmic trade-offs** between **time** and **power**.

Speedup increases with additional power (SSSP+GPU example)



Summary so far:

Energy optimality ~ work optimality

Time ~ energy

Time and power trade-off

Assumes uniform time and energy costs.

What if the costs are **non-uniform**?

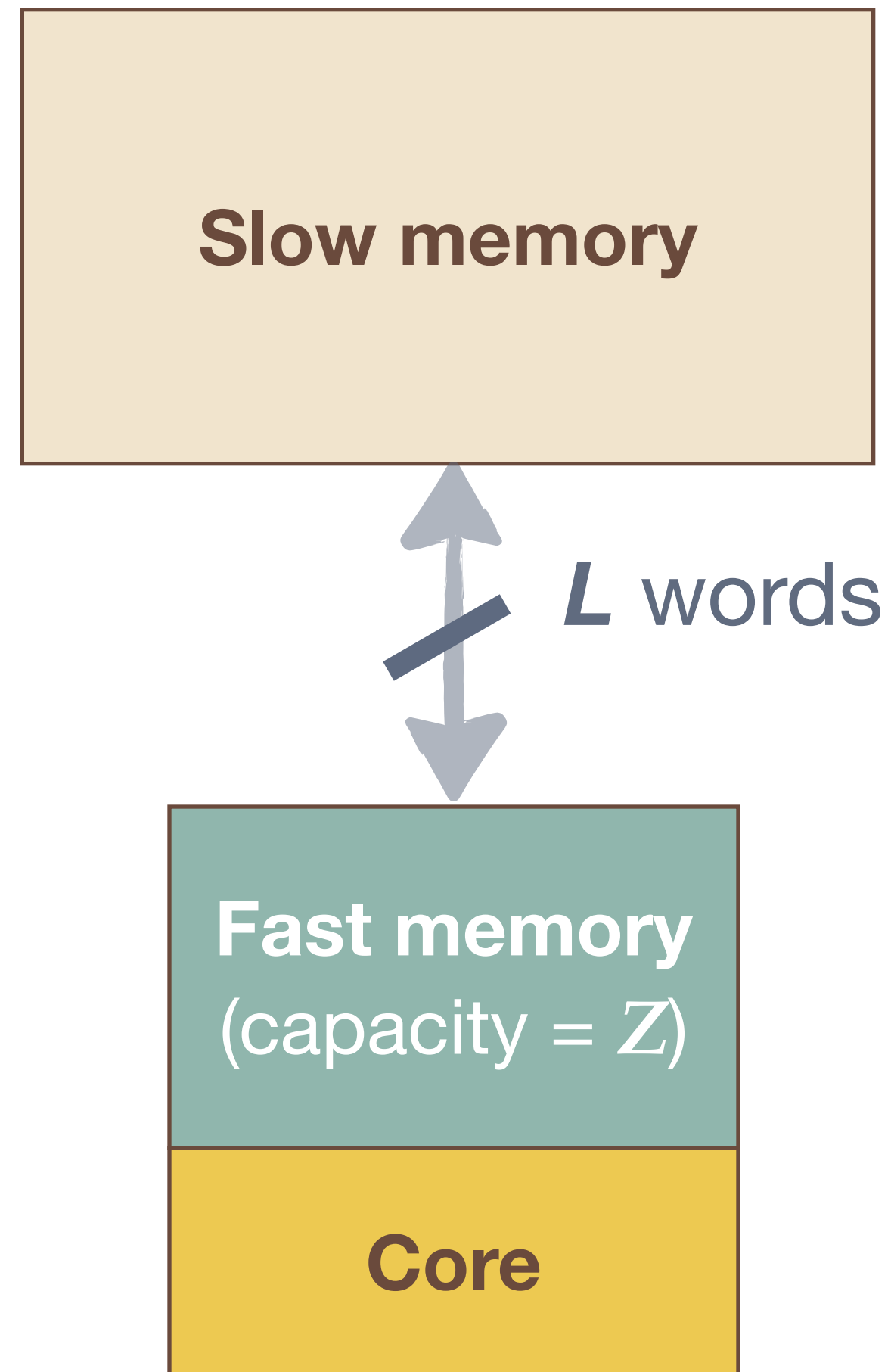
(At least) two interesting cases:

- 1. Operation (“op”) costs may change, e.g., under DVFS.**
- 2. Ops differ in cost, e.g., computation vs. communication.**

(At least) two interesting cases:

1. **Operation (“op”) costs may change**, e.g., under DVFS.
2. **Ops differ in cost**, e.g., computation vs. communication.

Example: External memory model



$W(n)$ = work (total ops)

$Q(n; Z, L)$ = no. transfers

$W(n) / (Q(n; Z, L) \cdot L)$
= **computational intensity**
(ops / words)

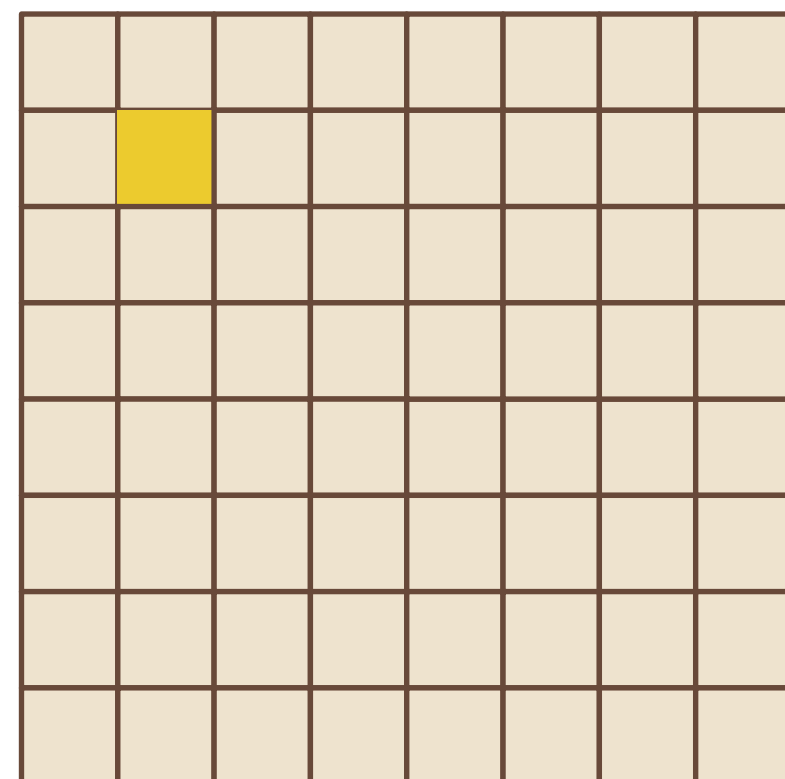
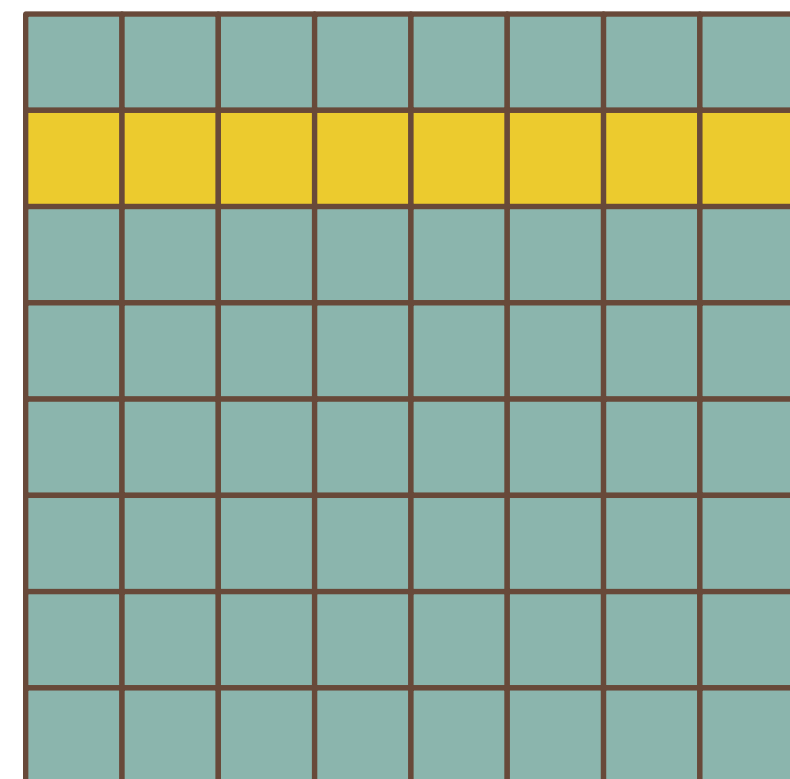
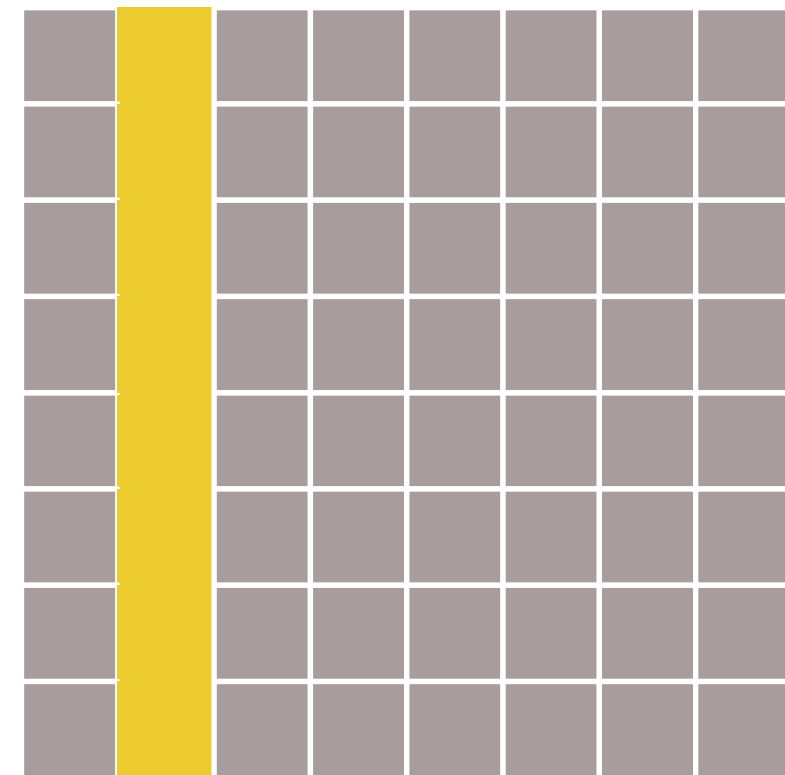
Desiderata:

Work optimality

Maximal intensity

Example: Matrix multiply

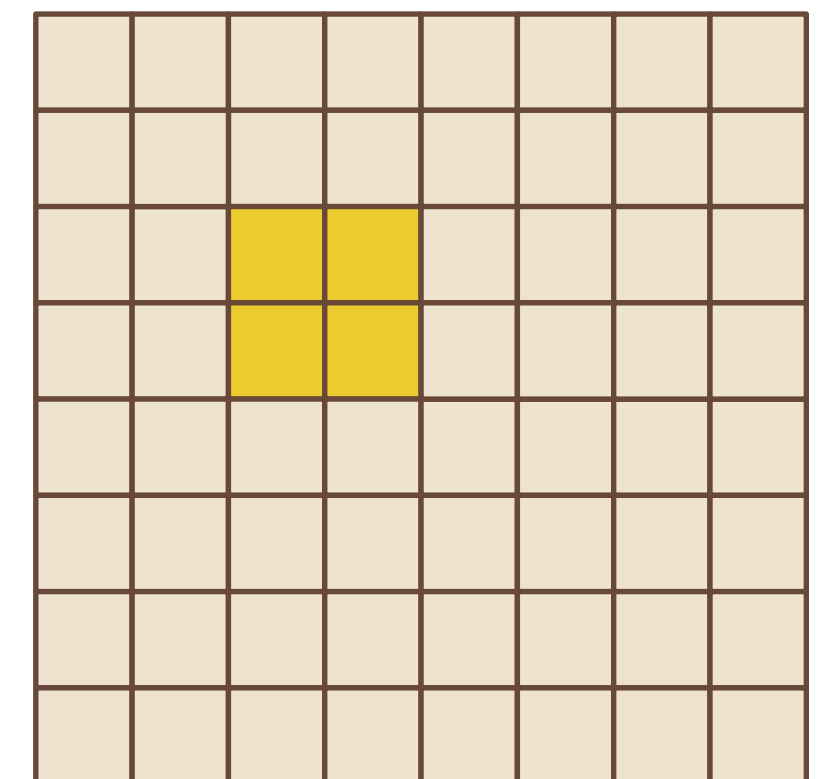
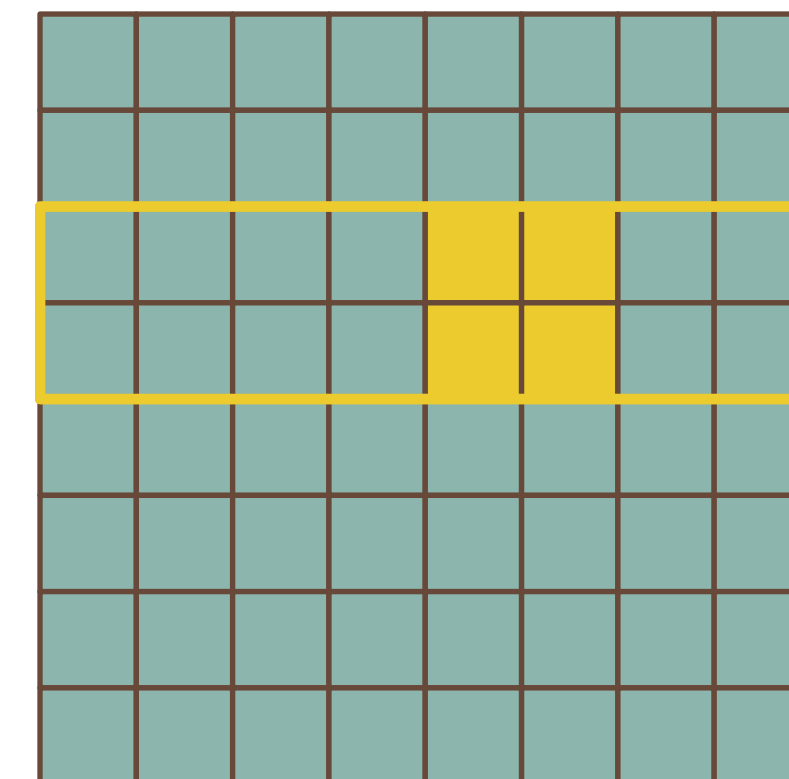
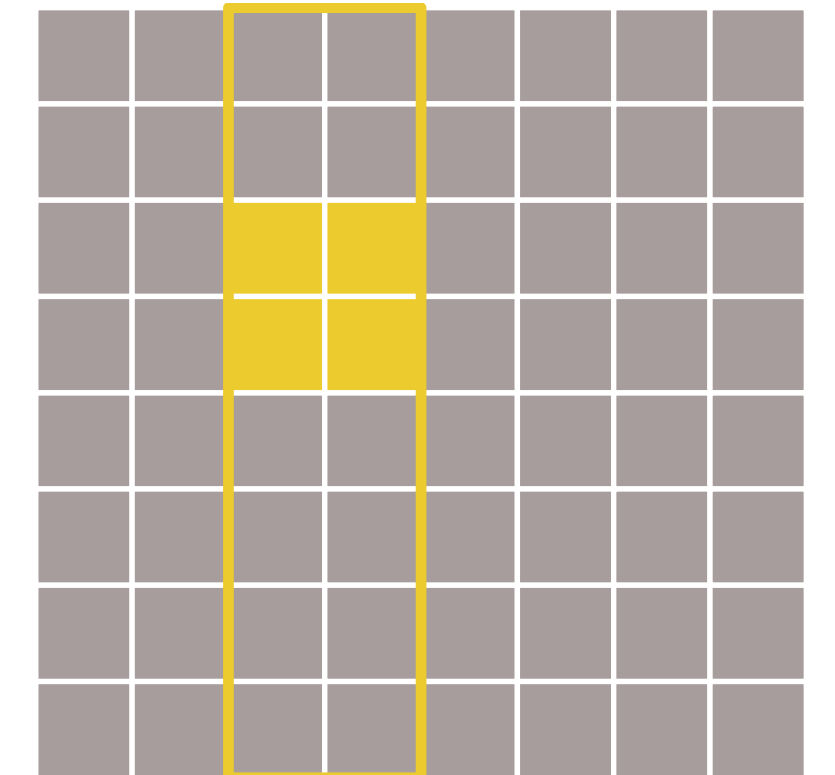
(non-Strassen)



$$C \leftarrow C + A * B$$

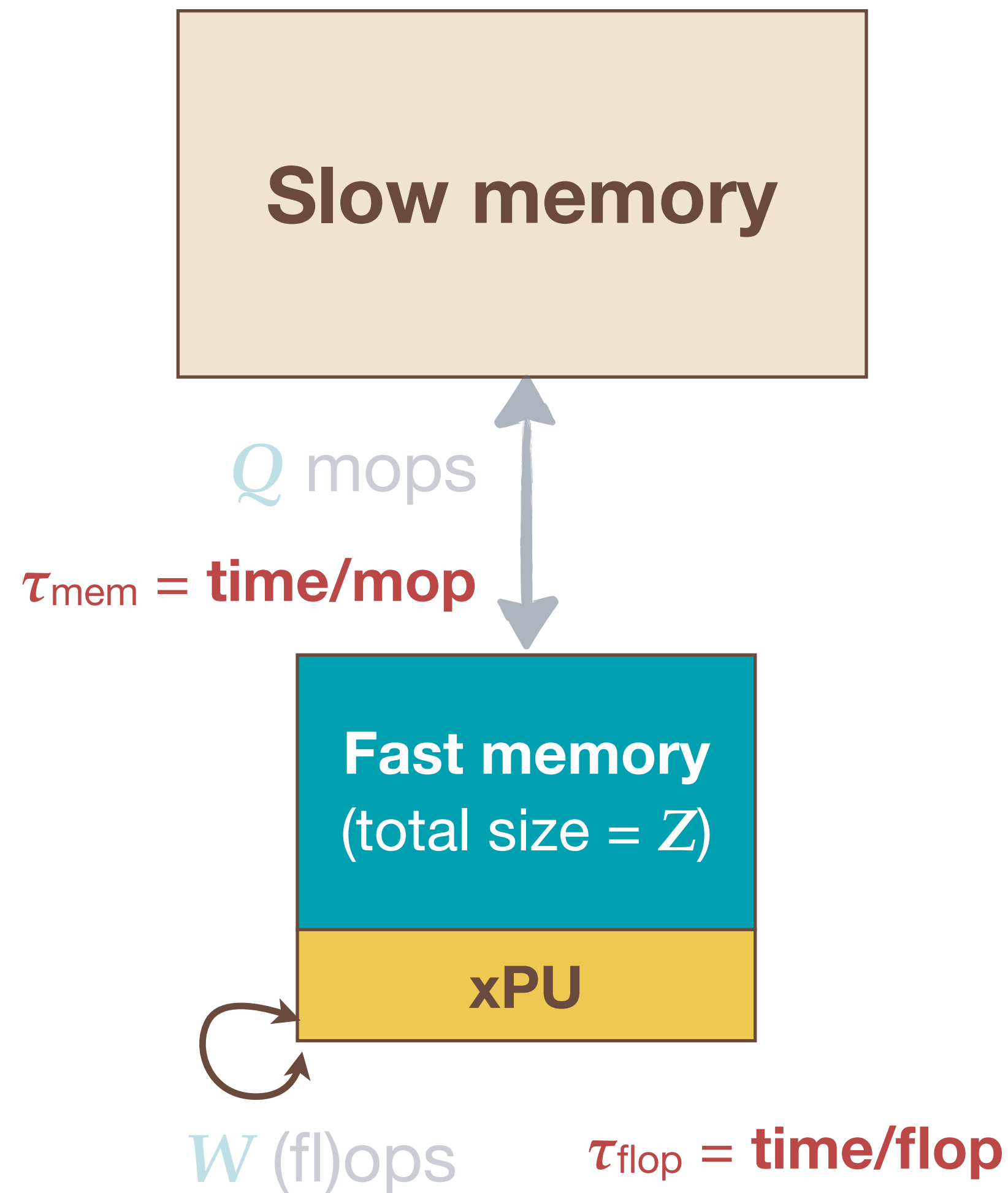
$$\frac{W(n)}{Q(n; Z, L) \cdot L} = \Theta(1)$$

Intensity



$$C \leftarrow C + A * B$$

$$\frac{W}{Q \cdot L} = \mathcal{O}(\sqrt{Z})$$



$W \equiv \# \text{ (fl)ops}$

$Q \equiv \# \text{ mem. ops (mops)} = Q(Z)$

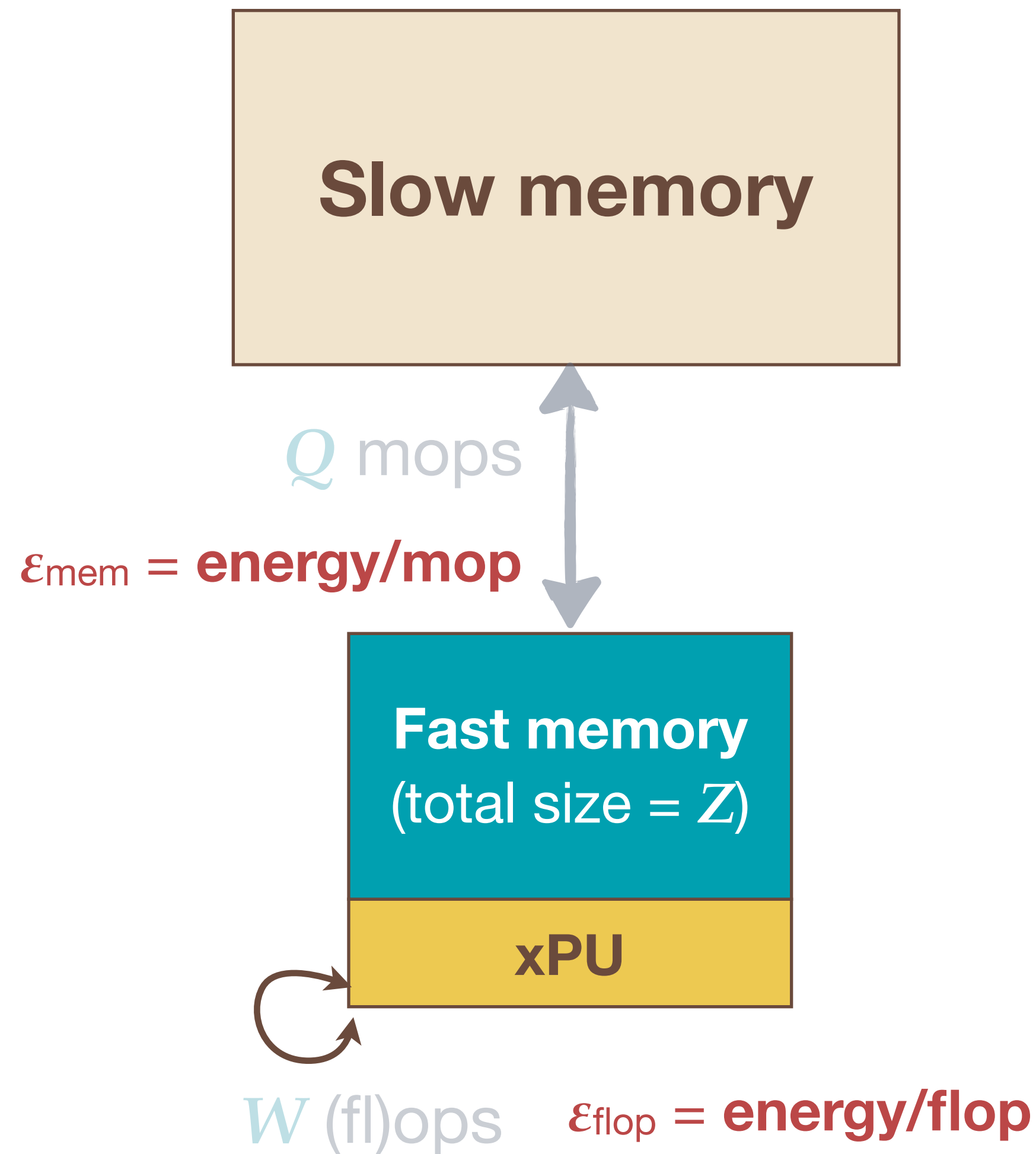
$I \equiv \frac{W}{Q} = \text{Intensity (flop:mop)}$

$\tau_{\text{flop}} \equiv \text{time per (fl)op}$

$\tau_{\text{mem}} \equiv \text{time per mop}$

$B_{\tau} \equiv \frac{\tau_{\text{mem}}}{\tau_{\text{flop}}} = \text{Balance (flop:mop)}$

Running time



$W \equiv \#$ (fl)ops

$Q \equiv \#$ mem. ops (mops) = $Q(Z)$

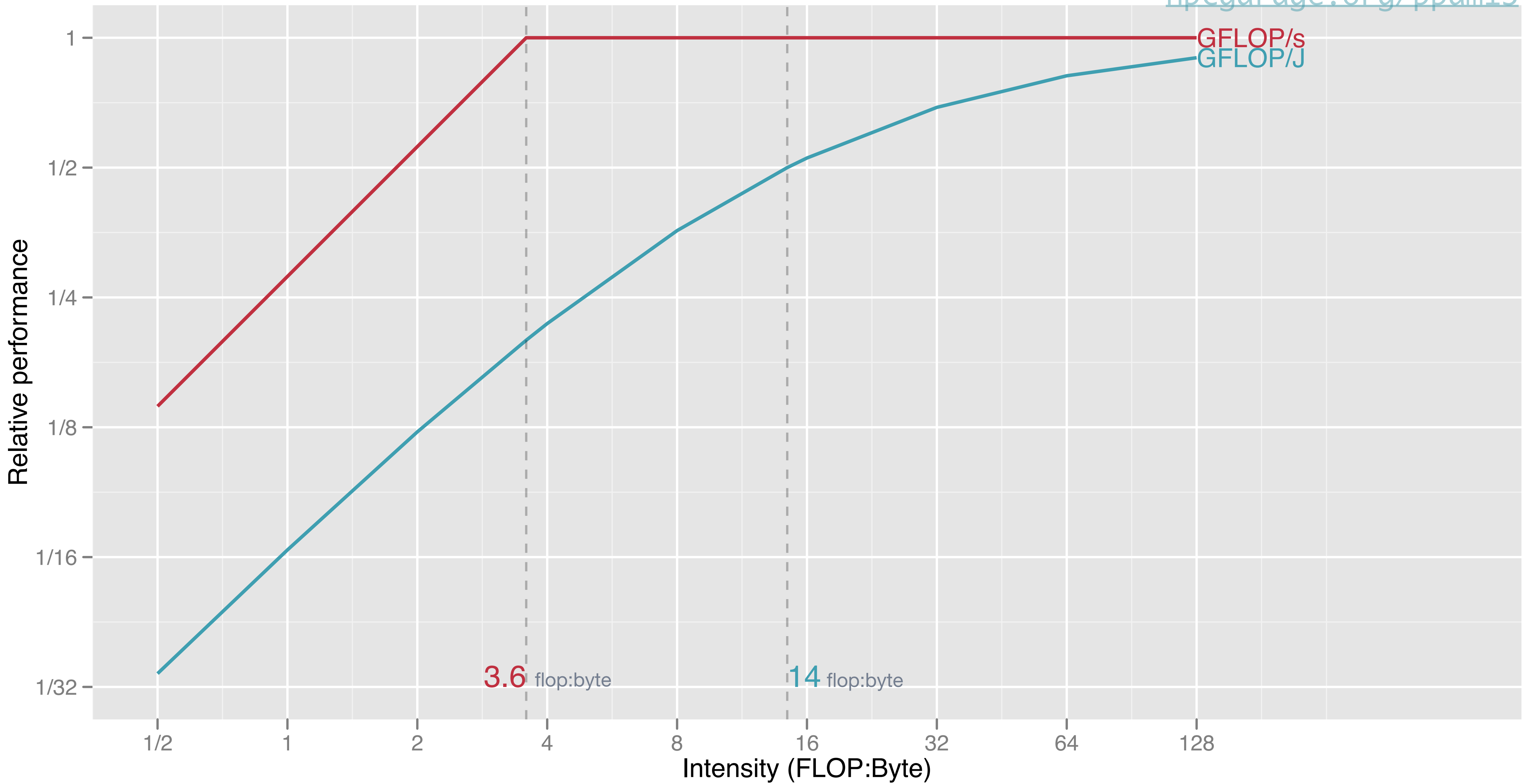
$I \equiv \frac{W}{Q} = \text{Intensity}$ (flop:mop)

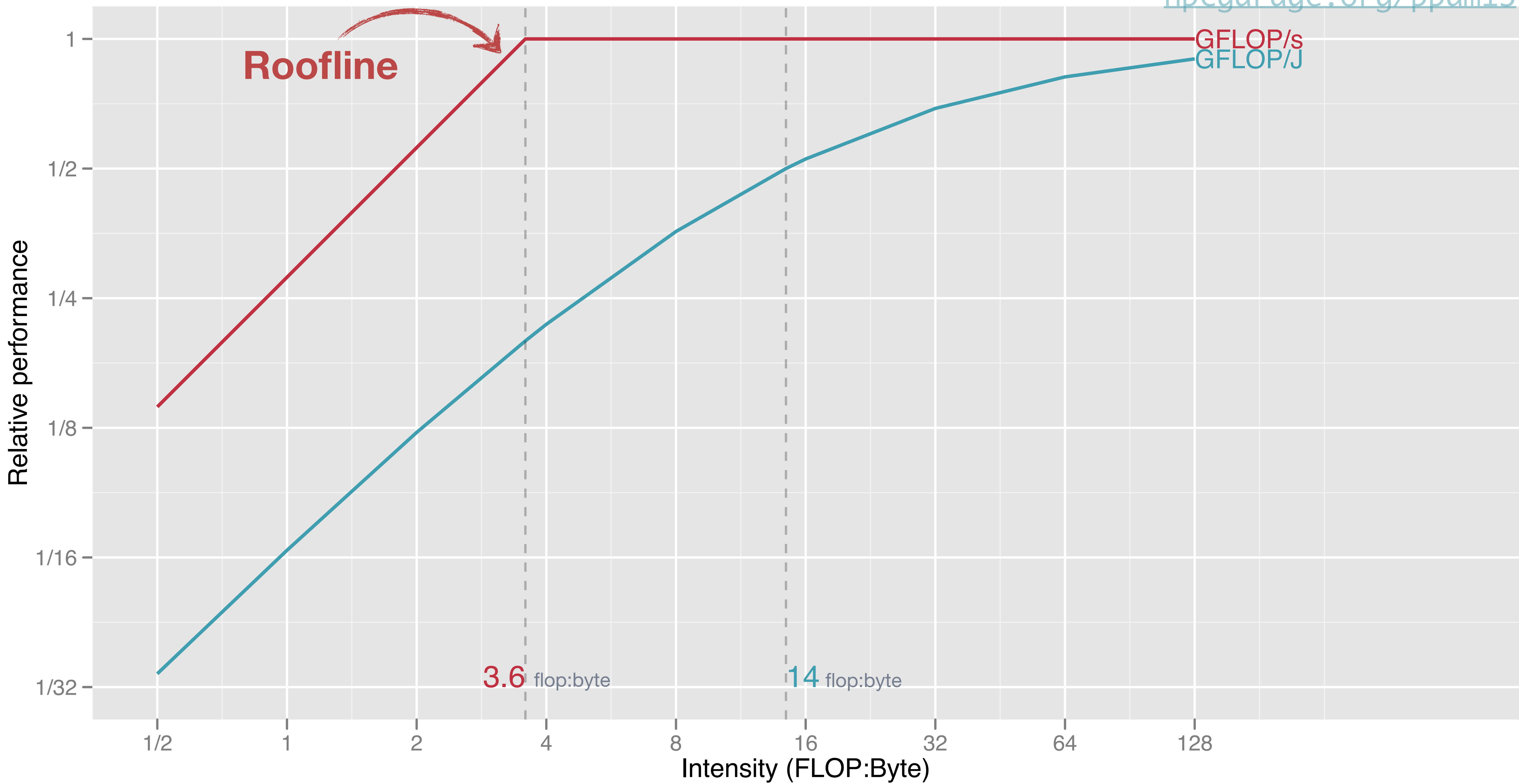
$\epsilon_{\text{flop}} \equiv$ energy per (fl)op

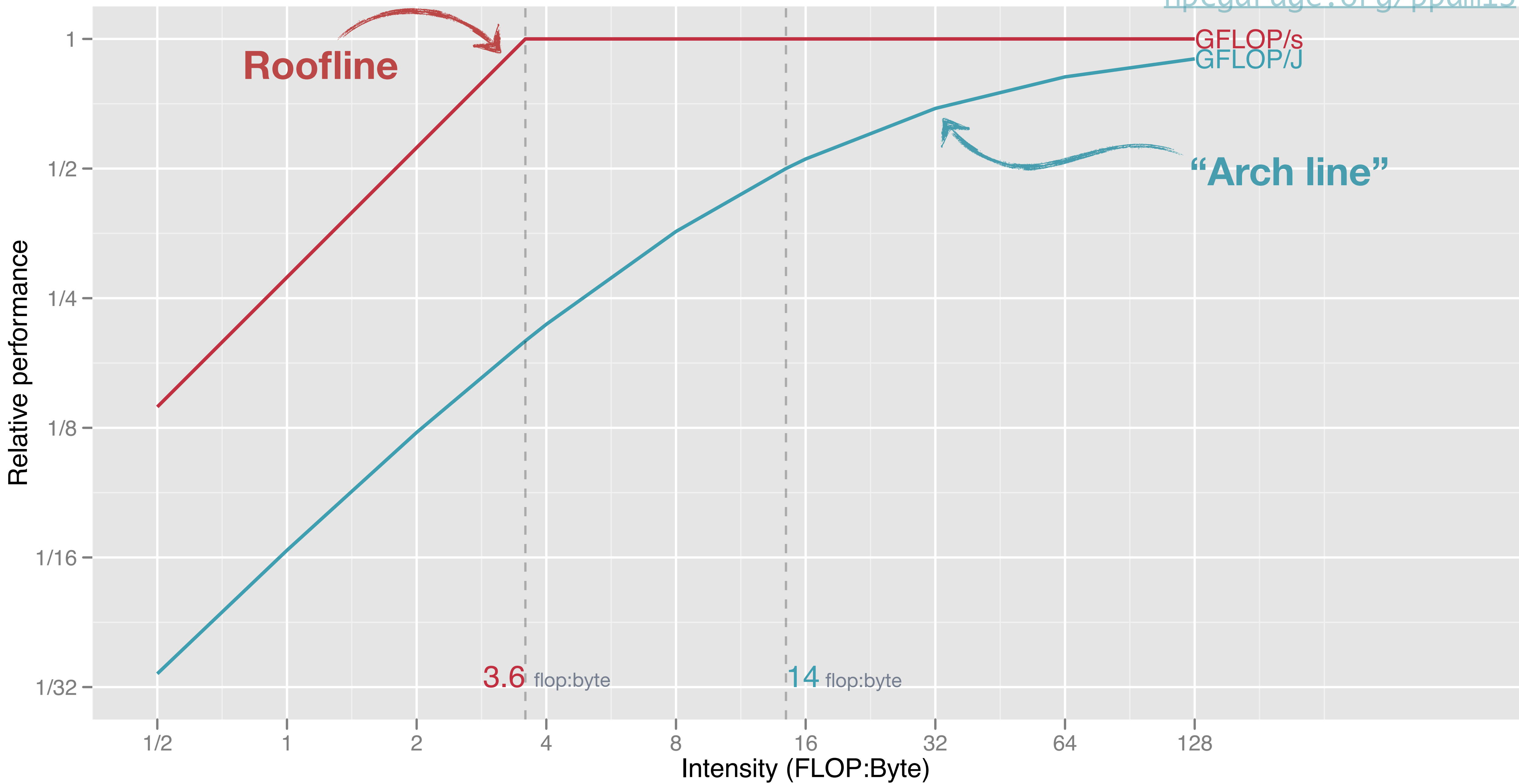
$\epsilon_{\text{mem}} \equiv$ energy per mop

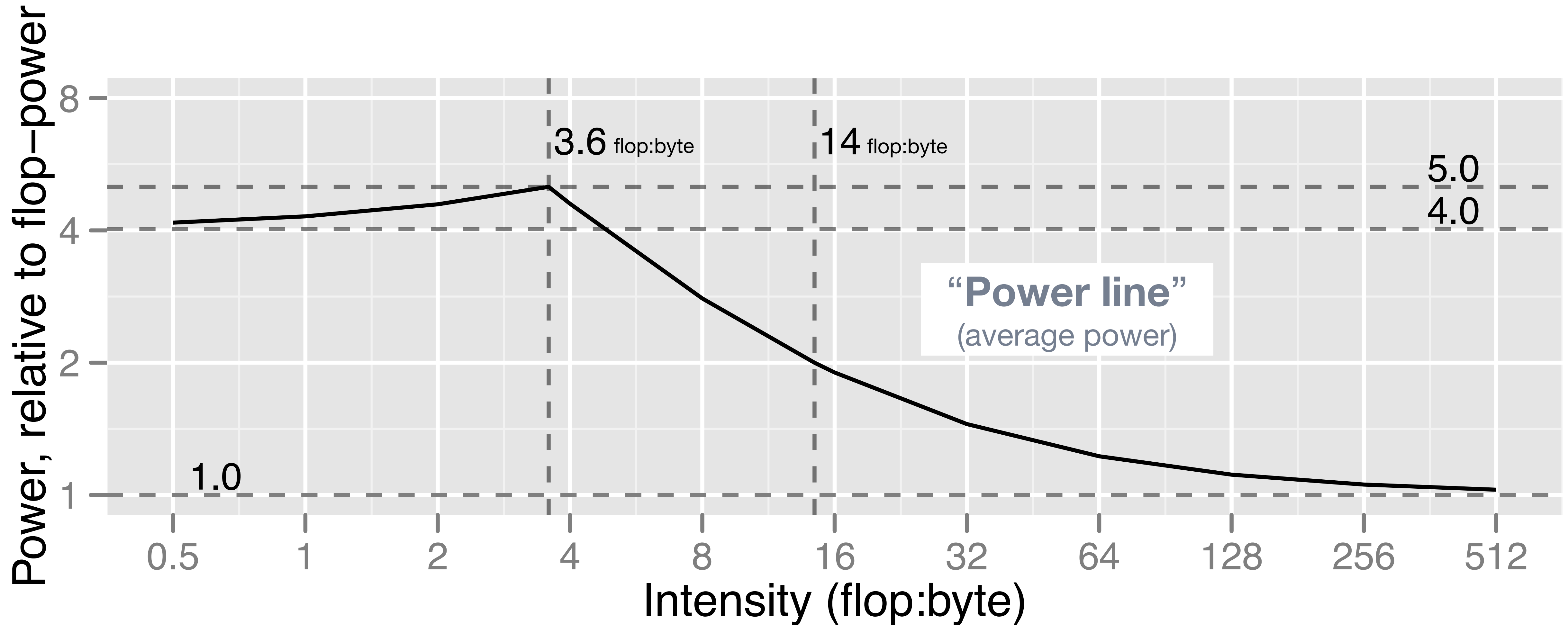
$B_{\epsilon} \equiv \frac{\epsilon_{\text{mem}}}{\epsilon_{\text{flop}}} = \text{Energy balance}$ (flop:mop)

Running energy

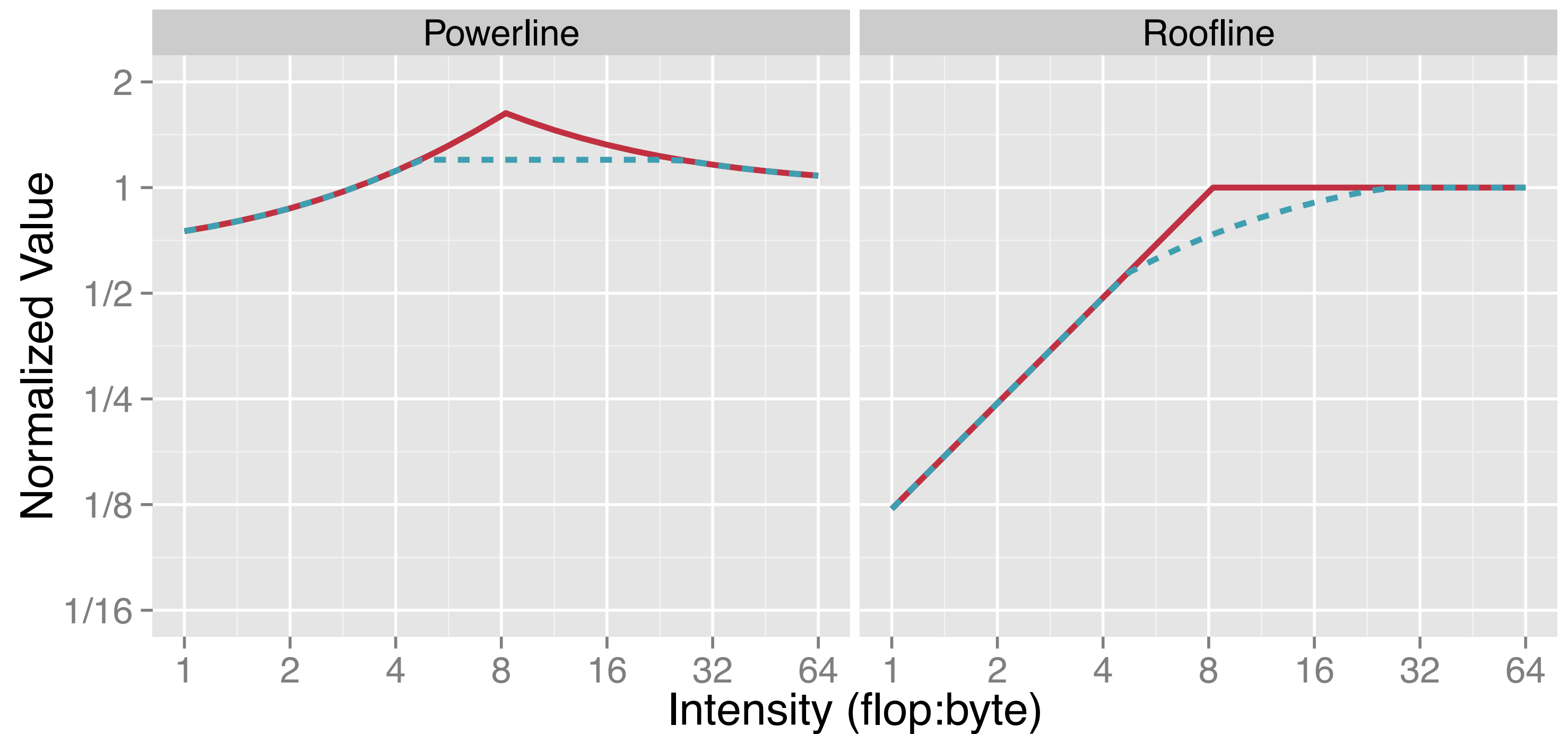




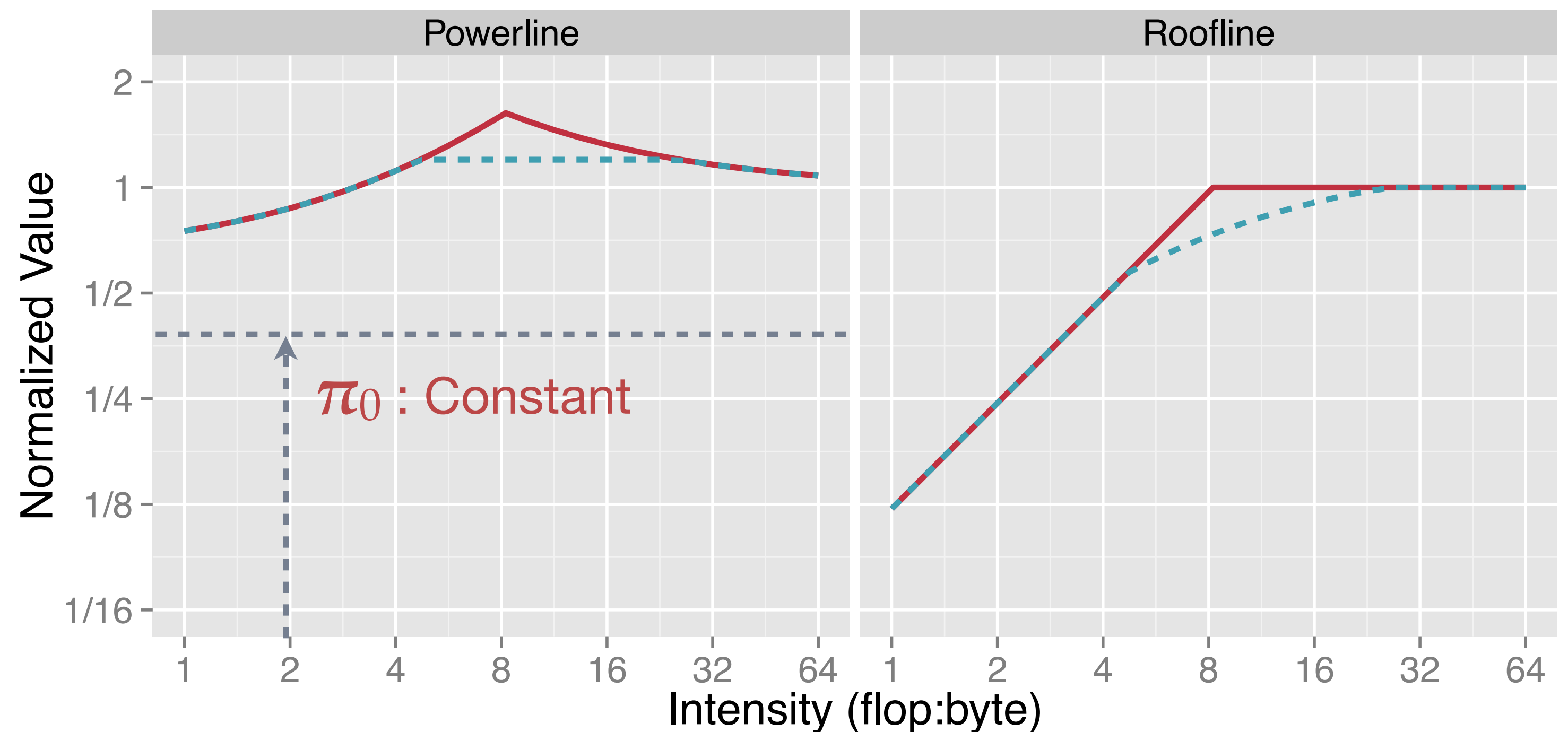




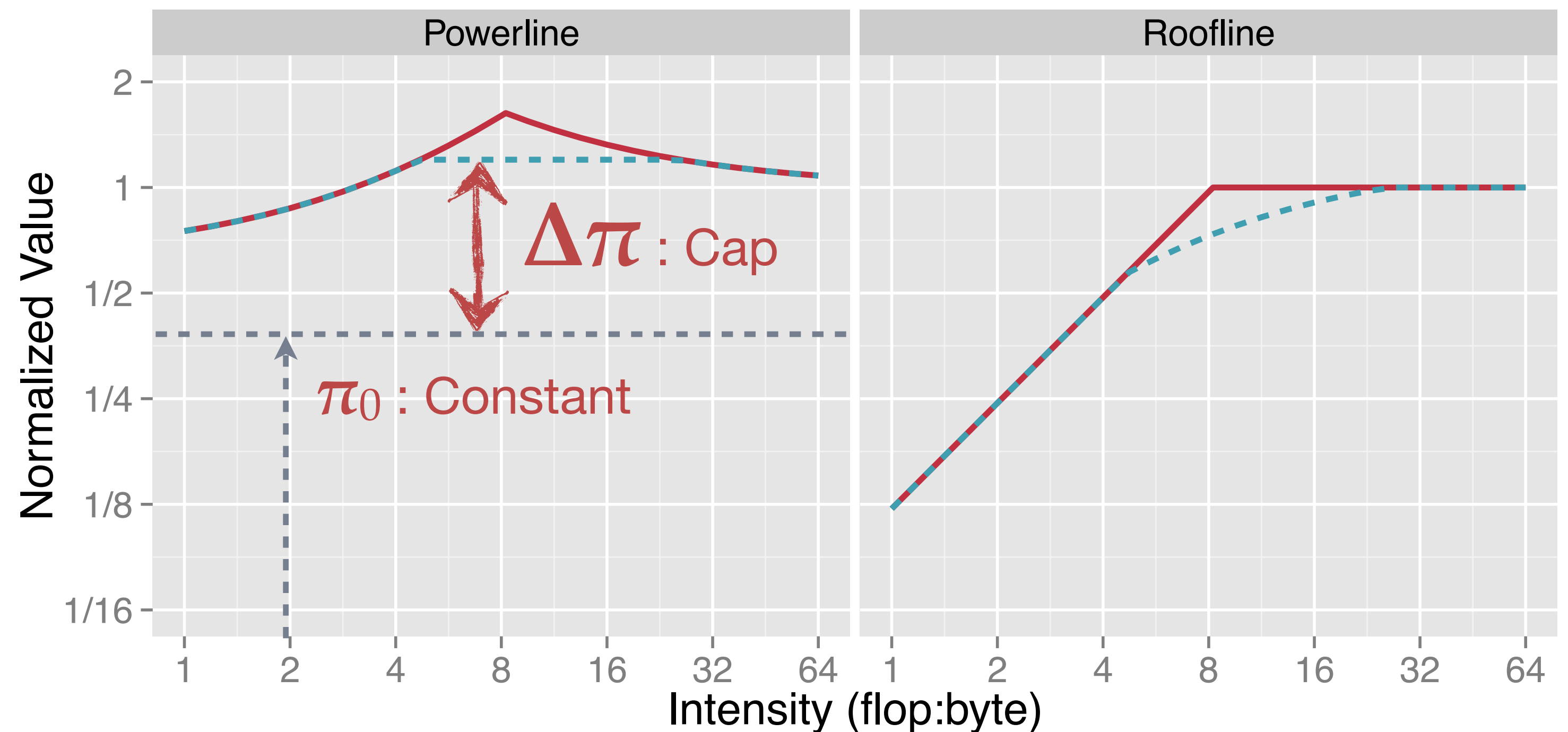
For real systems, the model should account for “**constant power**” and “**power capping.**”

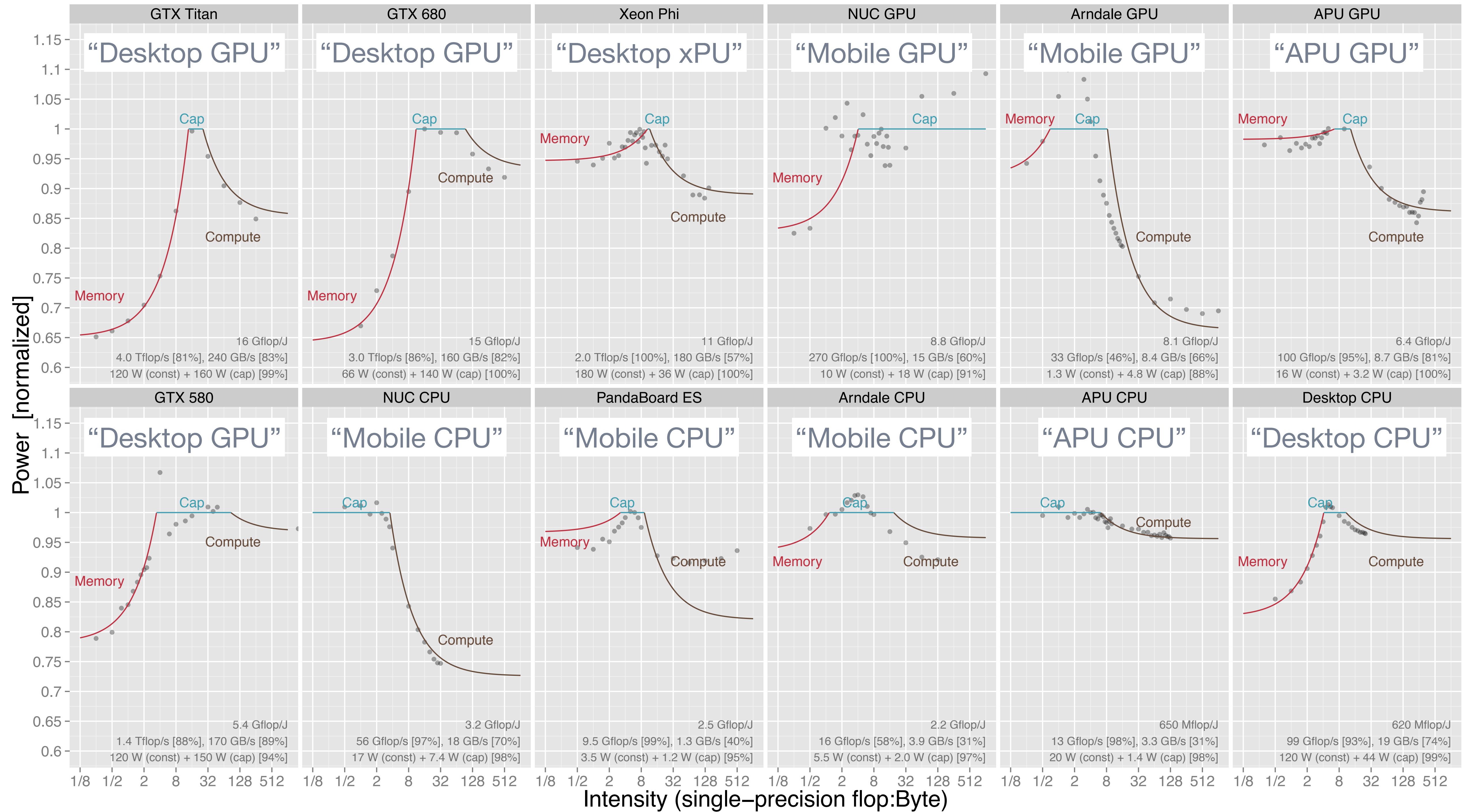


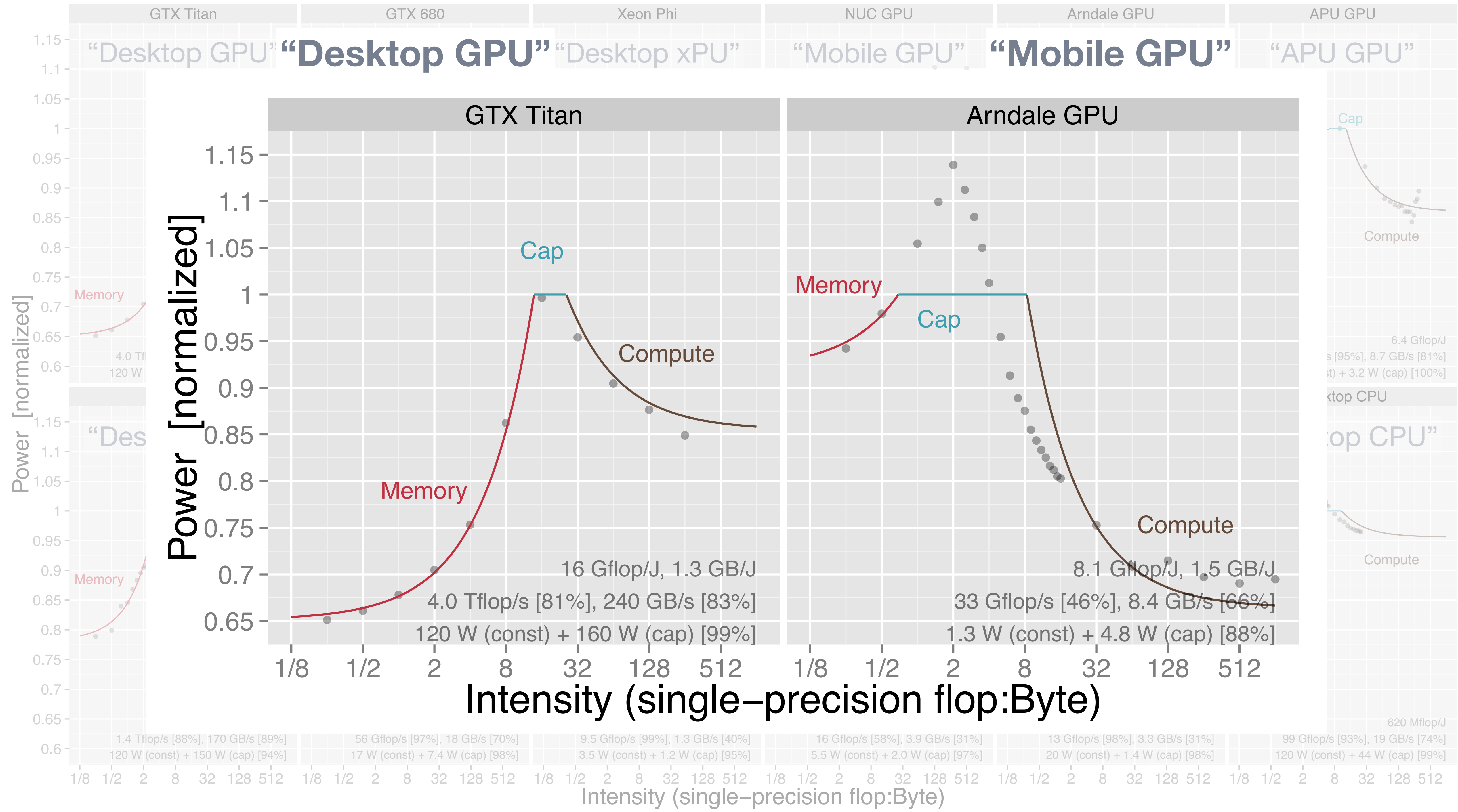
For real systems, the model should account for “**constant power**” and “**power capping.**”

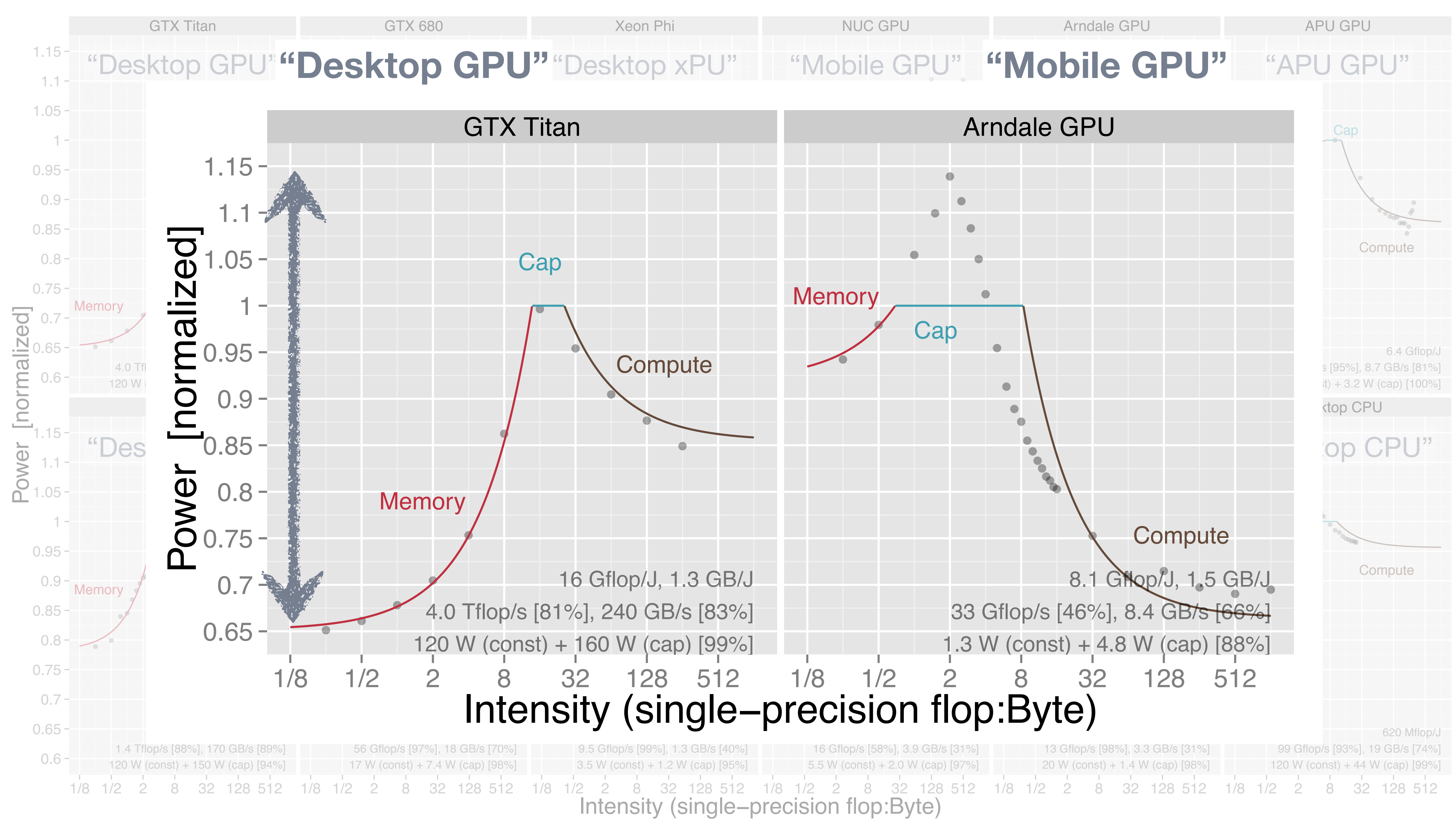


For real systems, the model should account for “**constant power**” and “**power capping.**”









Work-communication trade-offs

Abstract work-communication trade-offs

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

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

Abstract work-communication trade-offs

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

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

More work

Less
communication

Abstract work-communication trade-offs

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

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

More work

Less
communication

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

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

Abstract work-communication trade-offs

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

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

More work

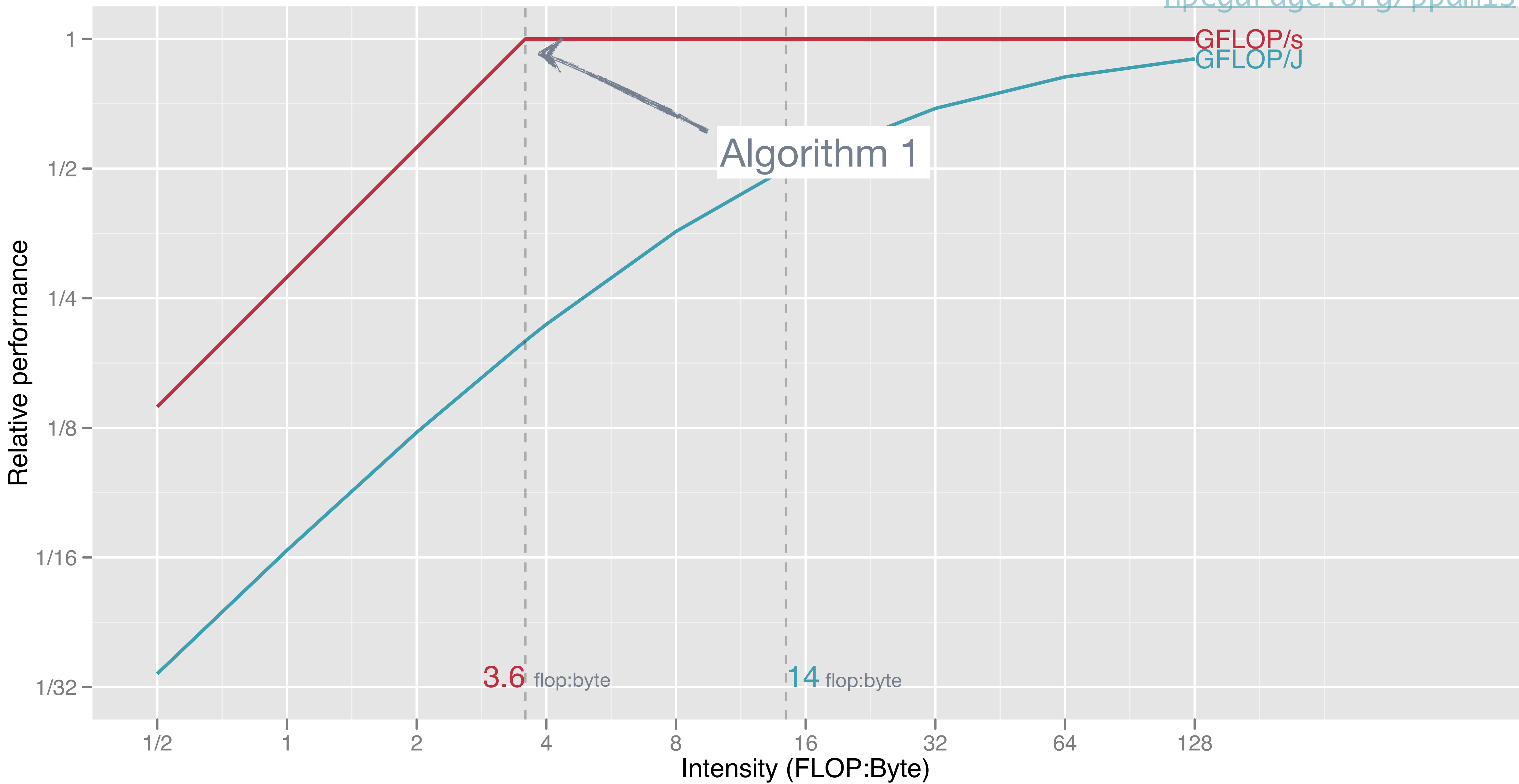
Less
communication

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

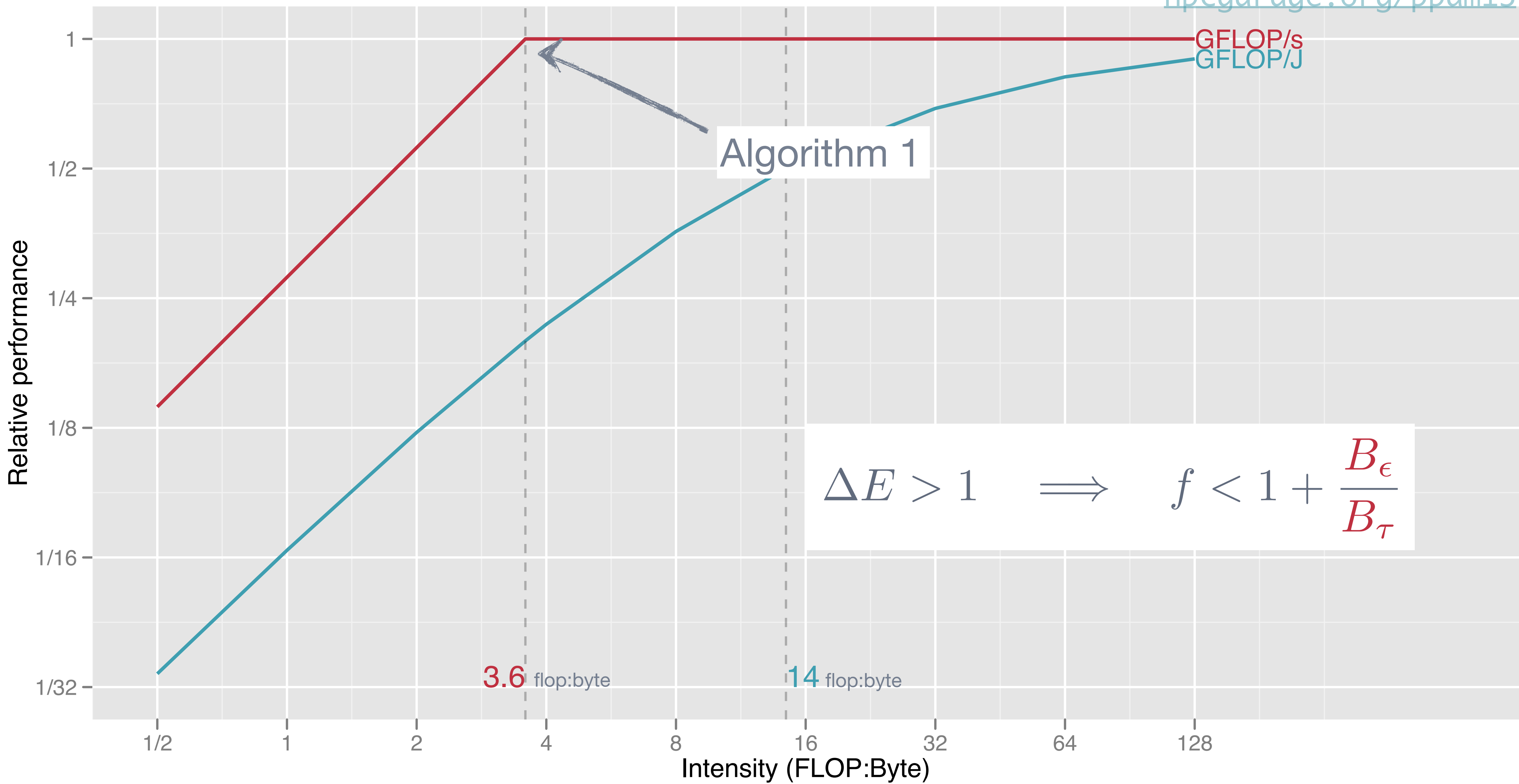
$$\text{"Greenup" } \Delta E = \frac{E_{1,1}}{E_{f,m}}$$

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

A general "greenup" condition



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



III. Case study: **Single-source shortest path**



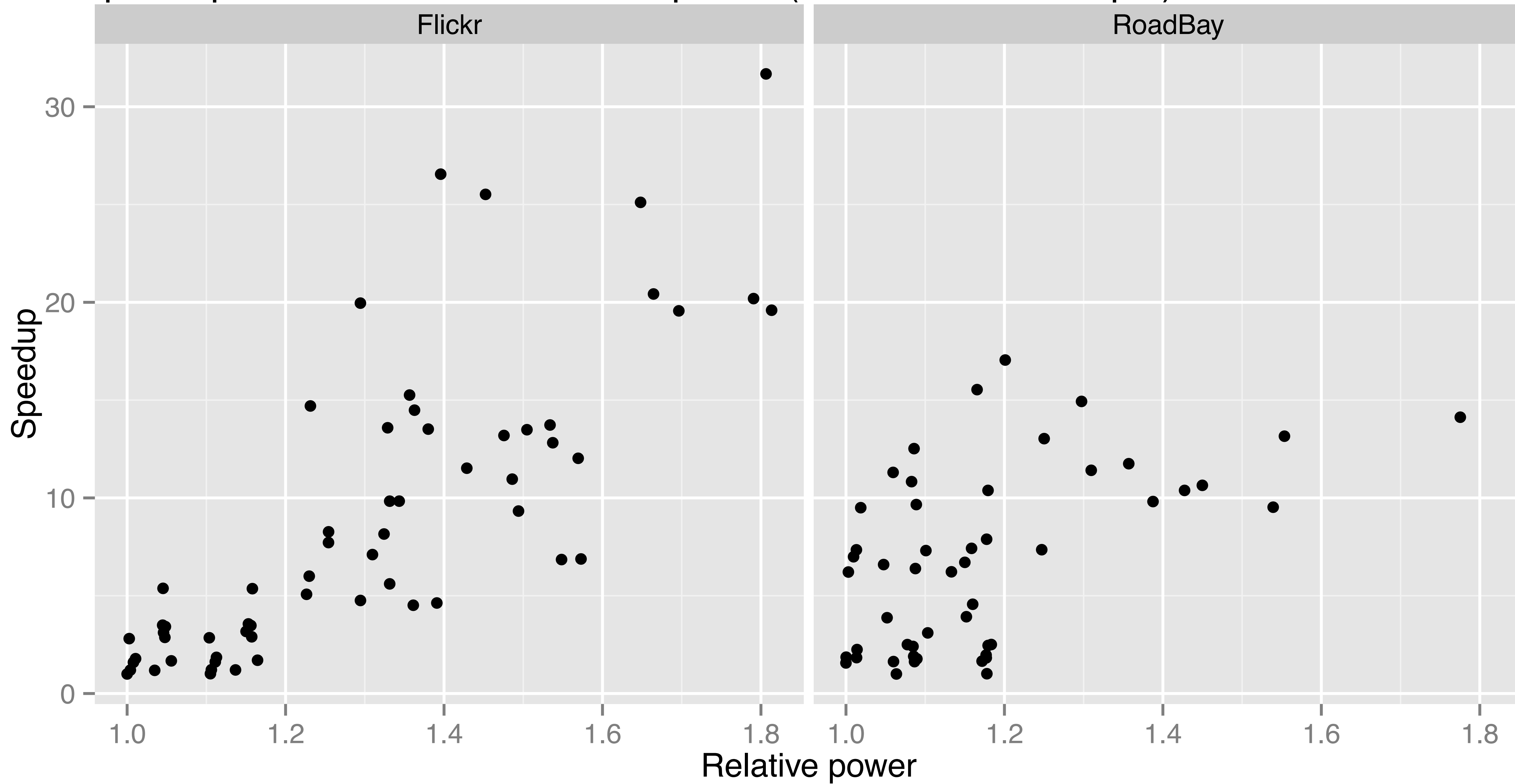
Sara Karamati, Jeff Young (PhD), R. Vuduc

- ⦿ Consider two implementation variants*
 - ⦿ “Bellman-Ford-like” — Highly parallel but not work-optimal
 - ⦿ “Delta-stepping-like” — Tunable work-parallelism tradeoff
 - ⦿ No preprocessing shortcuts, a la PHAST**
- ⦿ Both are tuned* for a GPU and we run them on an NVIDIA Jetson TK1, which has tunable core frequencies (**10x**) and memory frequencies (**3x**)

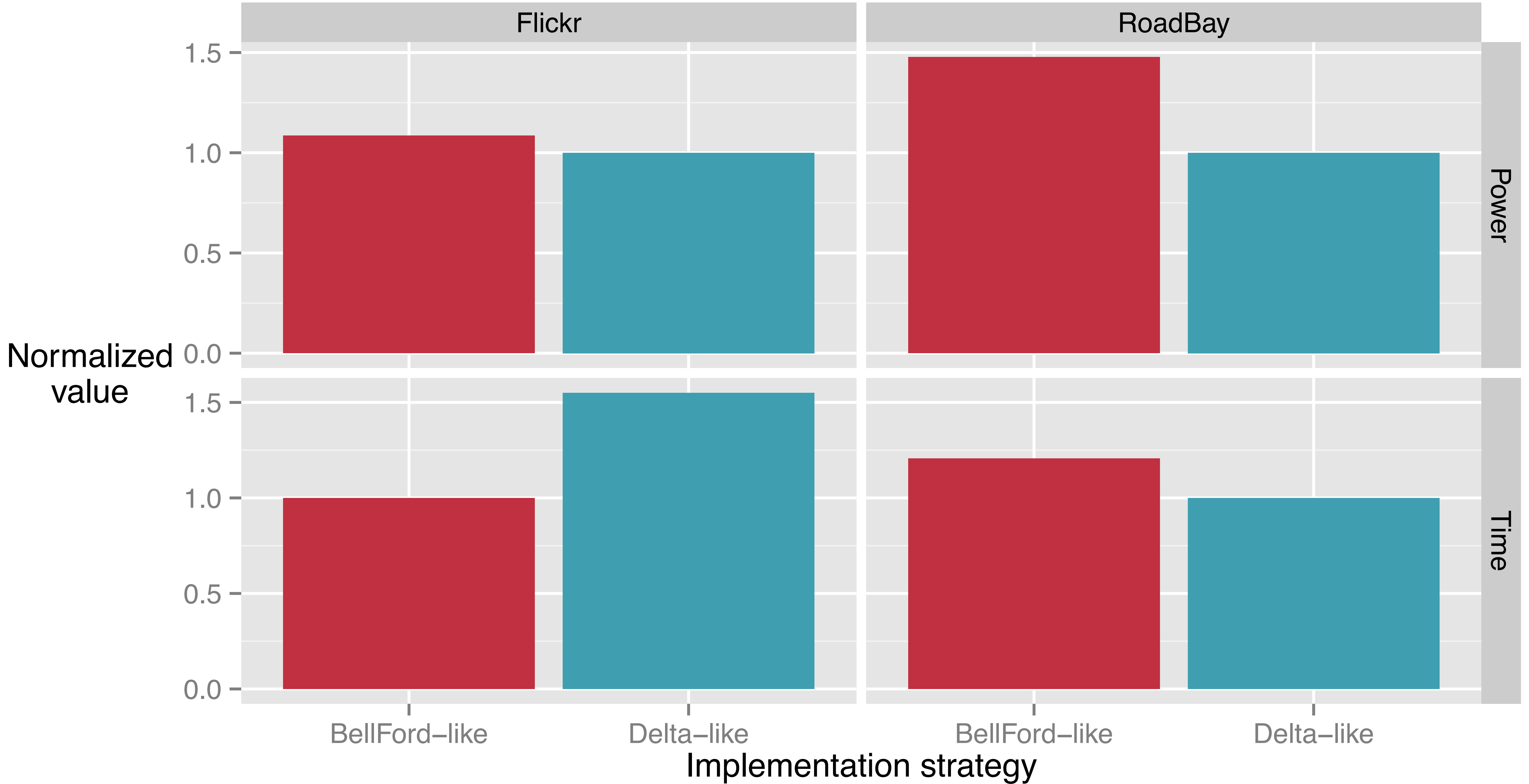
* These are GunRock implementations of Davidson, Baxter, Garland, and Owens (IPDPS'14)

** Delling et al. “PHAST: Hardware-accelerated shortest path trees” (JPDC'10)

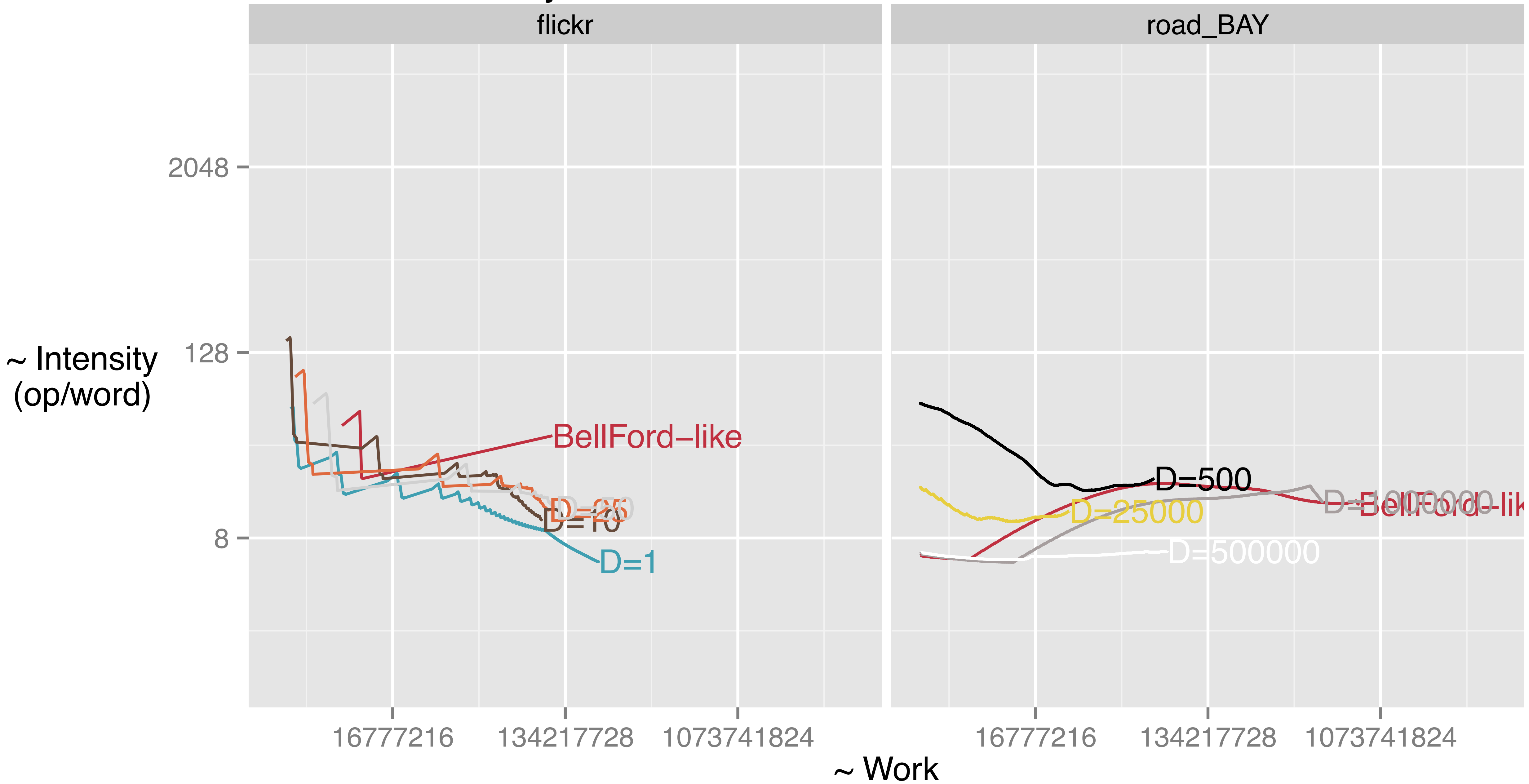
Speedup increases with additional power (SSSP+GPU example)



SSSP: Time and Power



SSSP: How intensity evolves



IV. Case study: **Branch-avoidance**

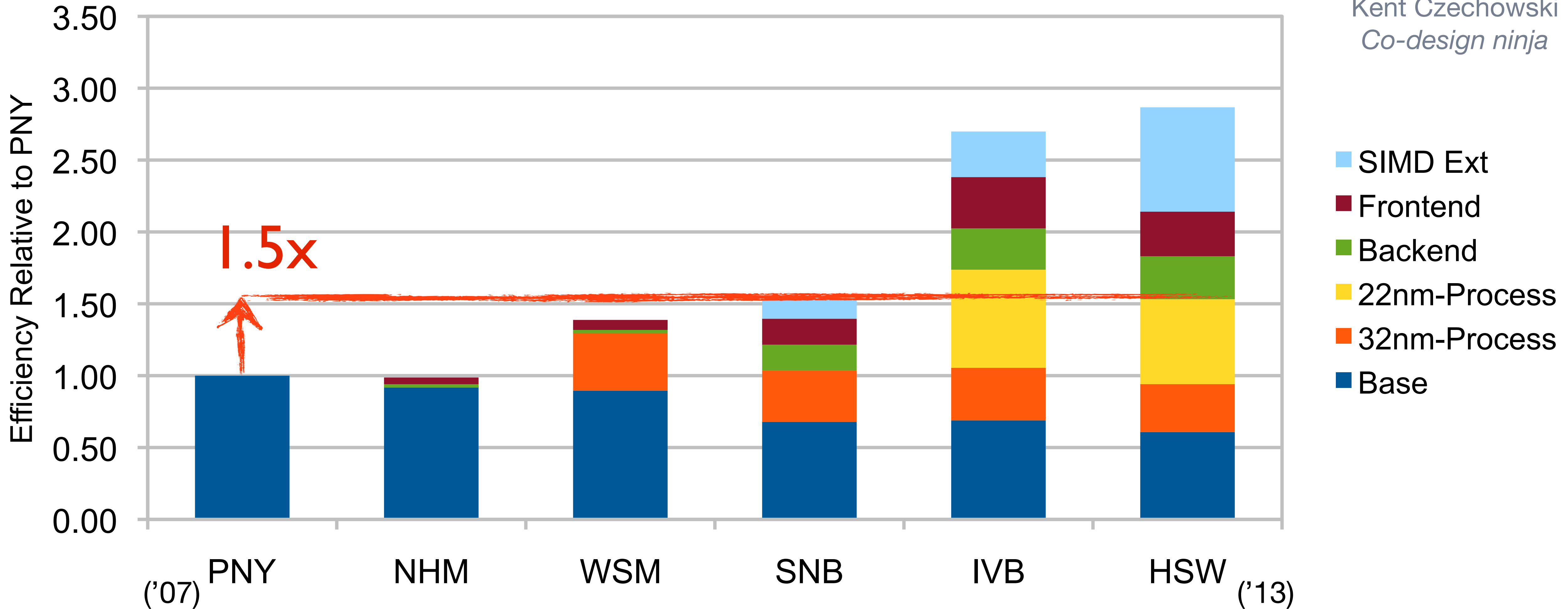


*Oded Green, Marat Dukhan, RV. "Branch-avoiding graph algorithms." In SPAA'15.
+ Post-paper analysis help from Anita Zakrzewska*

Improvement in Energy Efficiency Livermore Loops



Kent Czechowski
Co-design ninja



*Shiloach-Vishkin algorithm to compute
connected components (as labels)*

forall $v \in V$ **do**

label[v] \leftarrow int(v)

while ... **do**

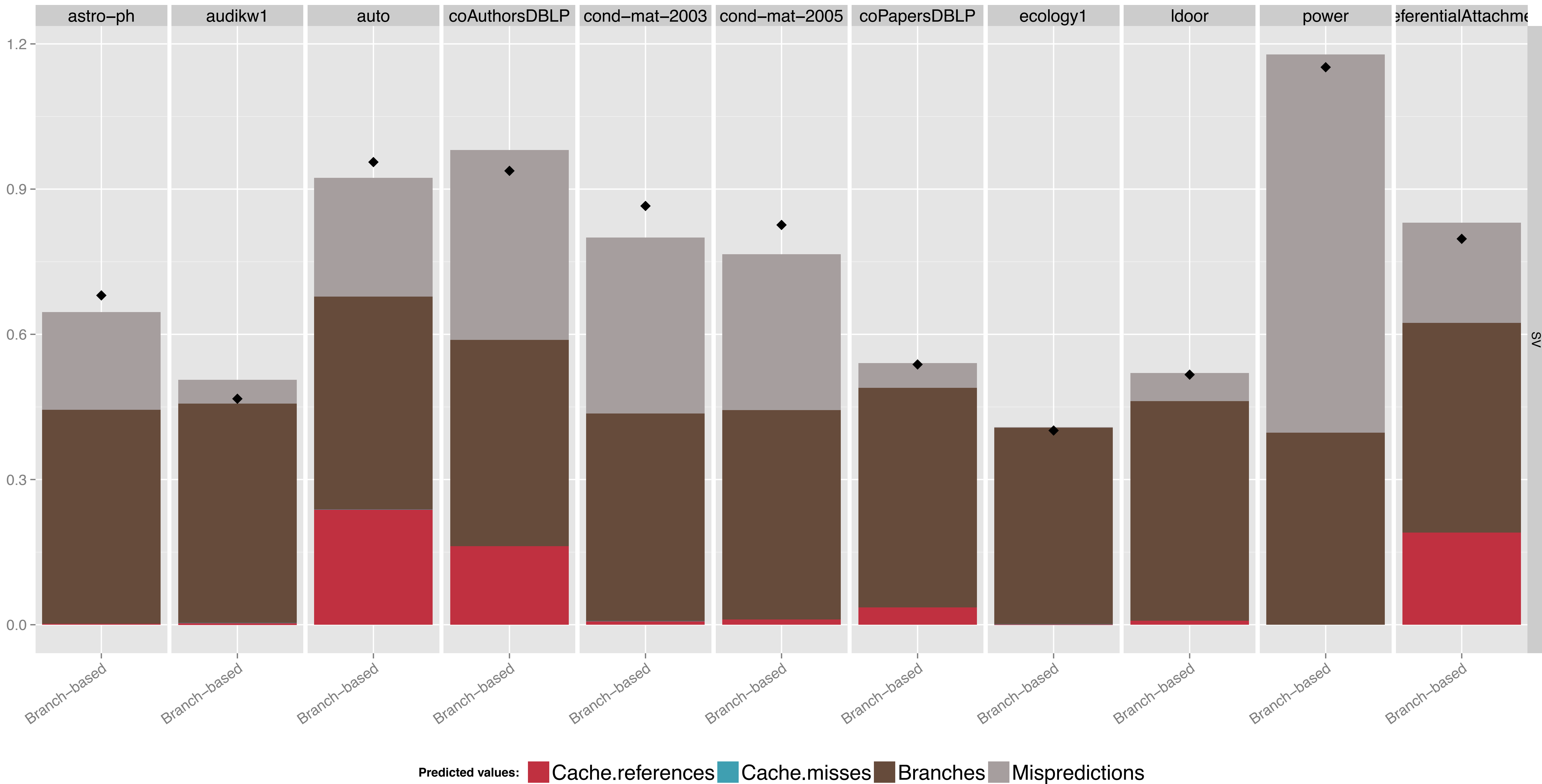
forall $v \in V$ **do**

forall $(v, u) \in E$ **do**

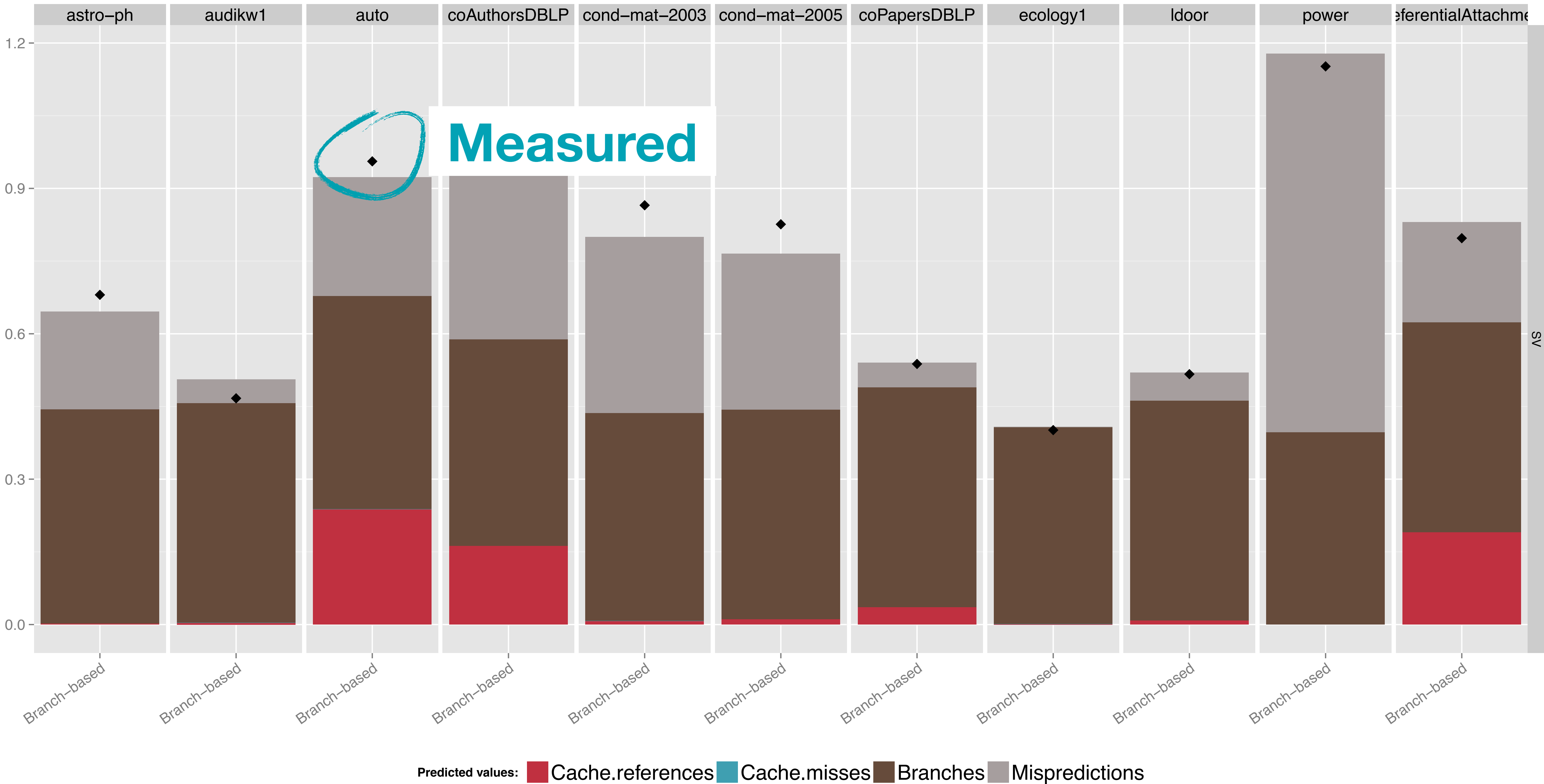
if label[v] < label[u] **then**

label[v] \leftarrow label[u]

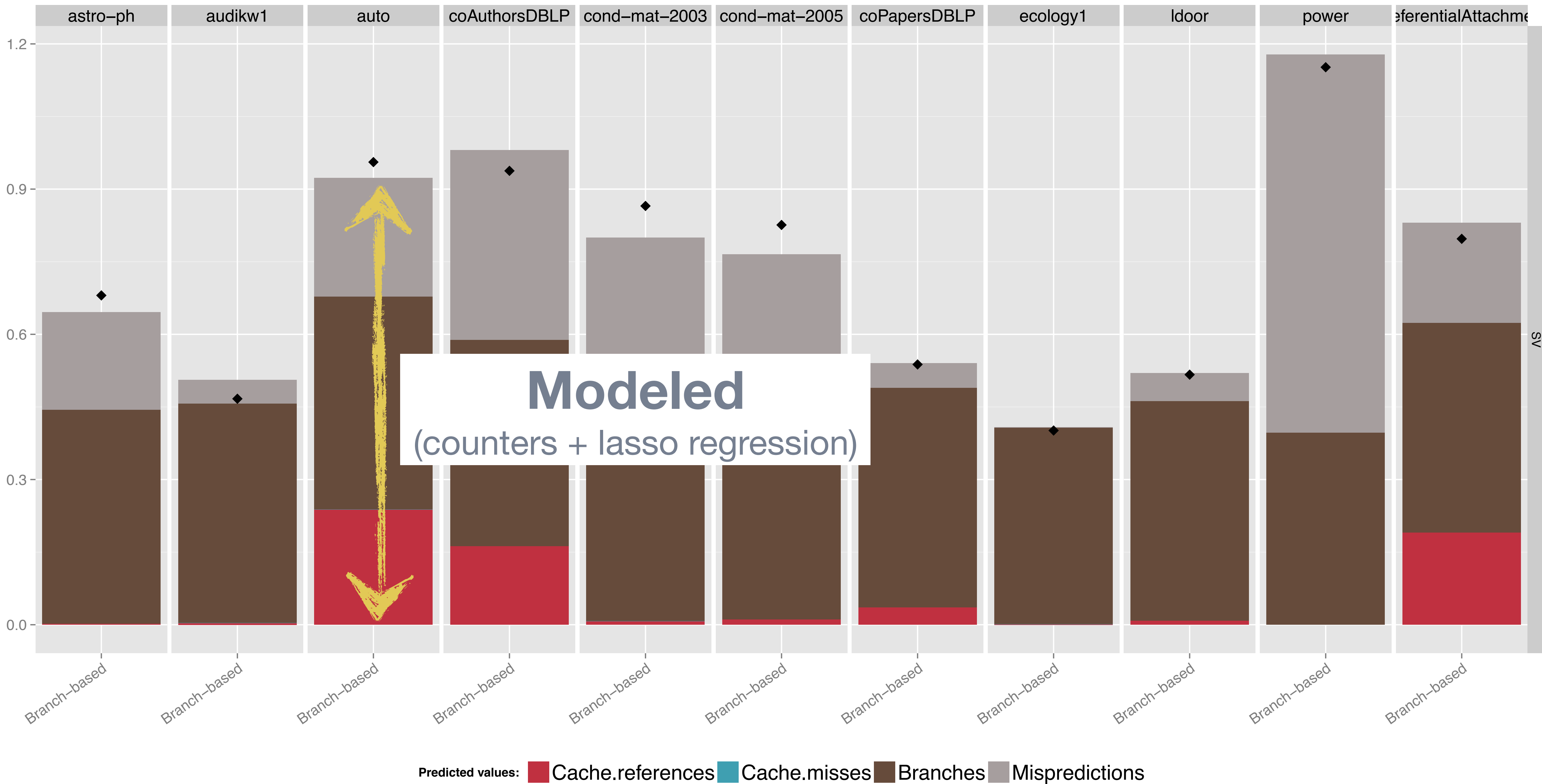
Predicted Cycles per instruction [Ivy Bridge]



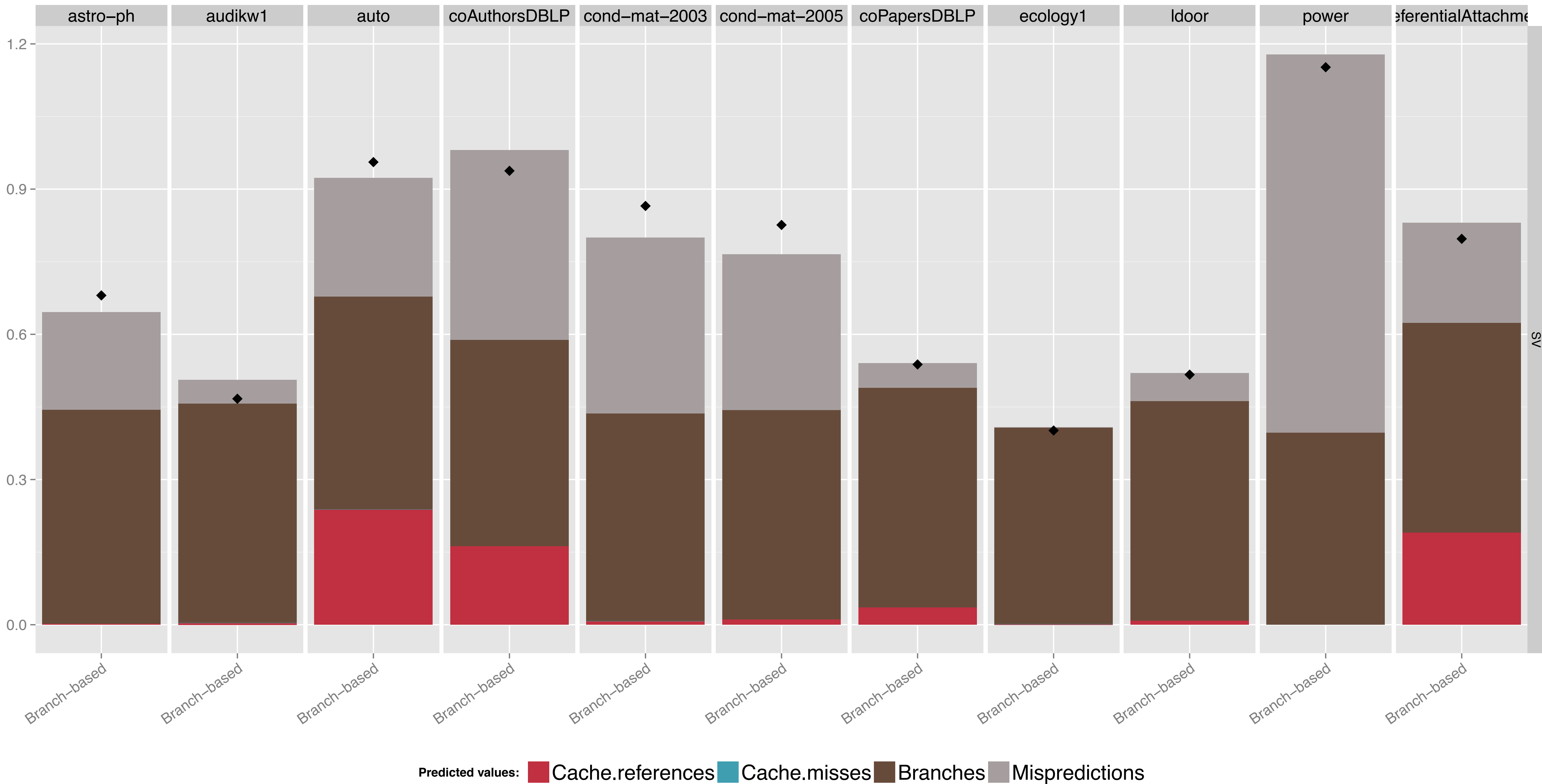
Predicted Cycles per instruction [Ivy Bridge]



Predicted Cycles per instruction [Ivy Bridge]



Predicted Cycles per instruction [Ivy Bridge]



Branch-based (original):

```
forall  $v \in V$  do  
  label[v]  $\leftarrow$  int(v)
```

```
while ... do  
  forall  $v \in V$  do  
    forall  $(v, u) \in E$  do  
      if label[u] < label[v] then  
        label[v]  $\leftarrow$  label[u]
```

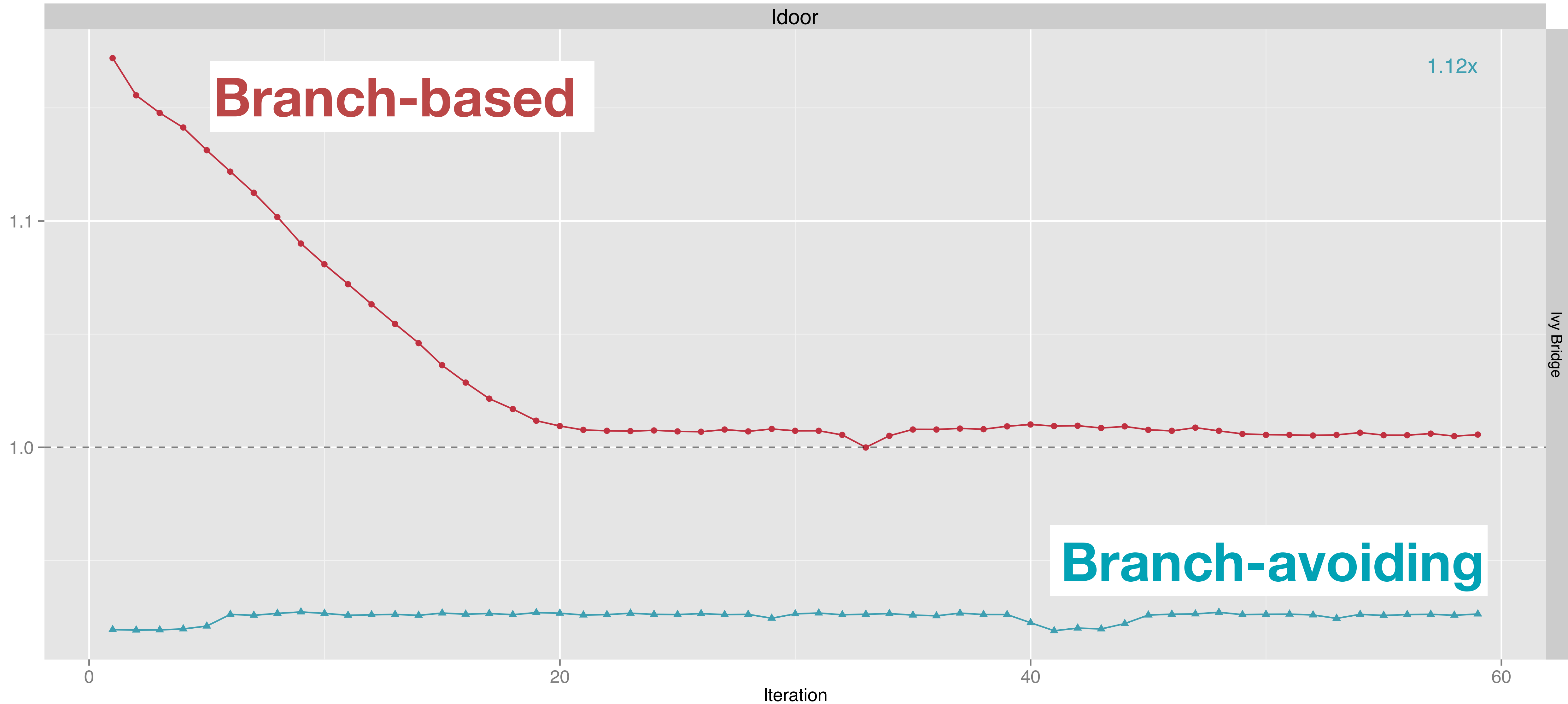
Branch-avoiding:

```
forall  $v \in V$  do  
  label[v]  $\leftarrow$  int(v)
```

```
while ... do  
  forall  $v \in V$  do  
    forall  $(v, u) \in E$  do  
      flag  $\leftarrow$  (label[u] < label[v])  
      cmov (label[v], label[u], flag)
```

Shiloach–Vishkin Connected Components: Cycles

[Normalized to branch-based minimum]



Summary

- The key high-level claim of this talk is that **“classical” principles** of algorithm and software design are not only relevant, but **even more important** when considering metrics beyond time, such as energy and power.
- My main suggested direction for **future research** is to explore **work-X** tradeoffs, where X in this talk included communication, parallelism, and branching behavior, but there may be many others!