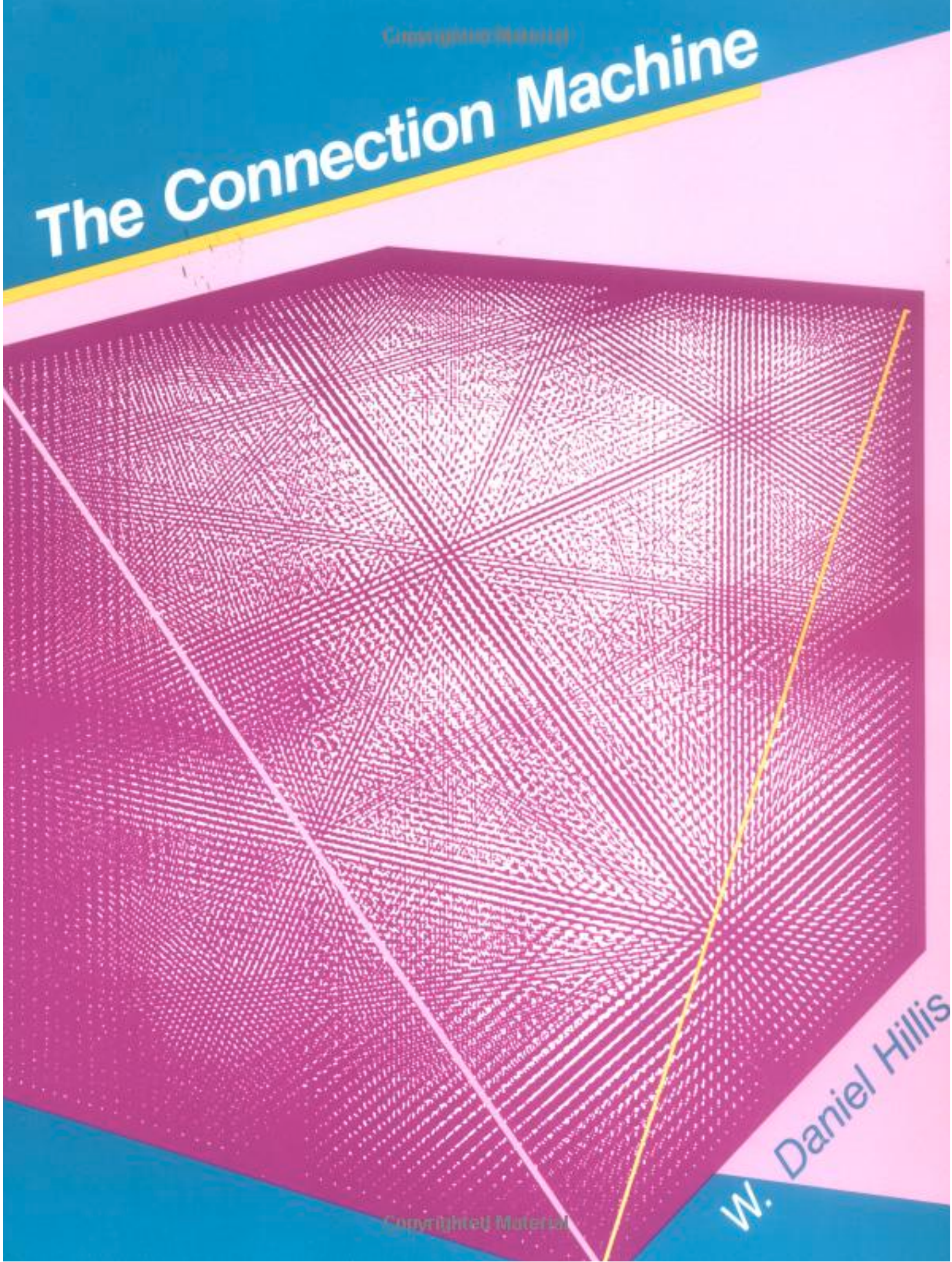# A power-tunable algorithm to compute single-source shortest paths
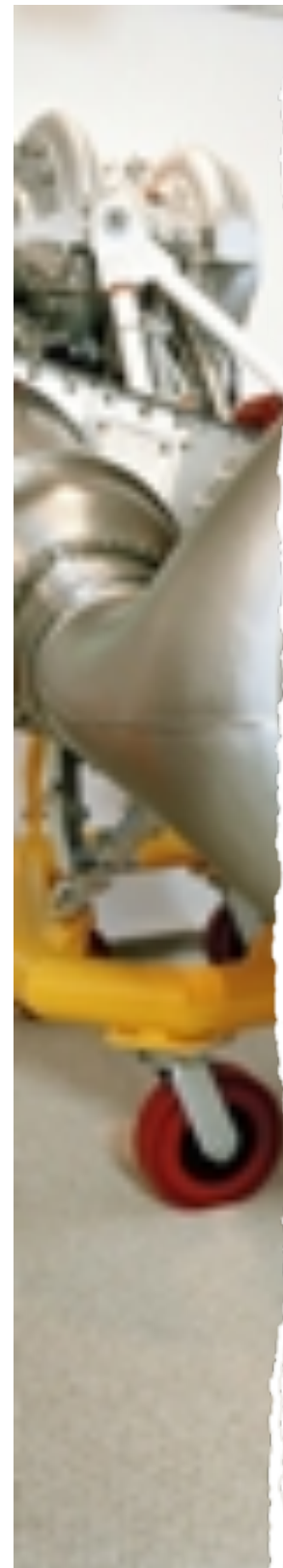
**Sara Karamati** · Jeffrey Young · **Richard (Rich) Vuduc**

November 17, 2017 — University of Colorado at Colorado Springs

Georgia Tech | College of Computing
Computational Science and Engineering

hpcgarage

# ACM Doctoral Dissertation Award Winner (1985)

The Connection Machine
W. Daniel Hillis

# What **physical constraint** limits speed today?

What **physical constraint** limits speed today?

**Energy** **Power**

COMPETITOR A

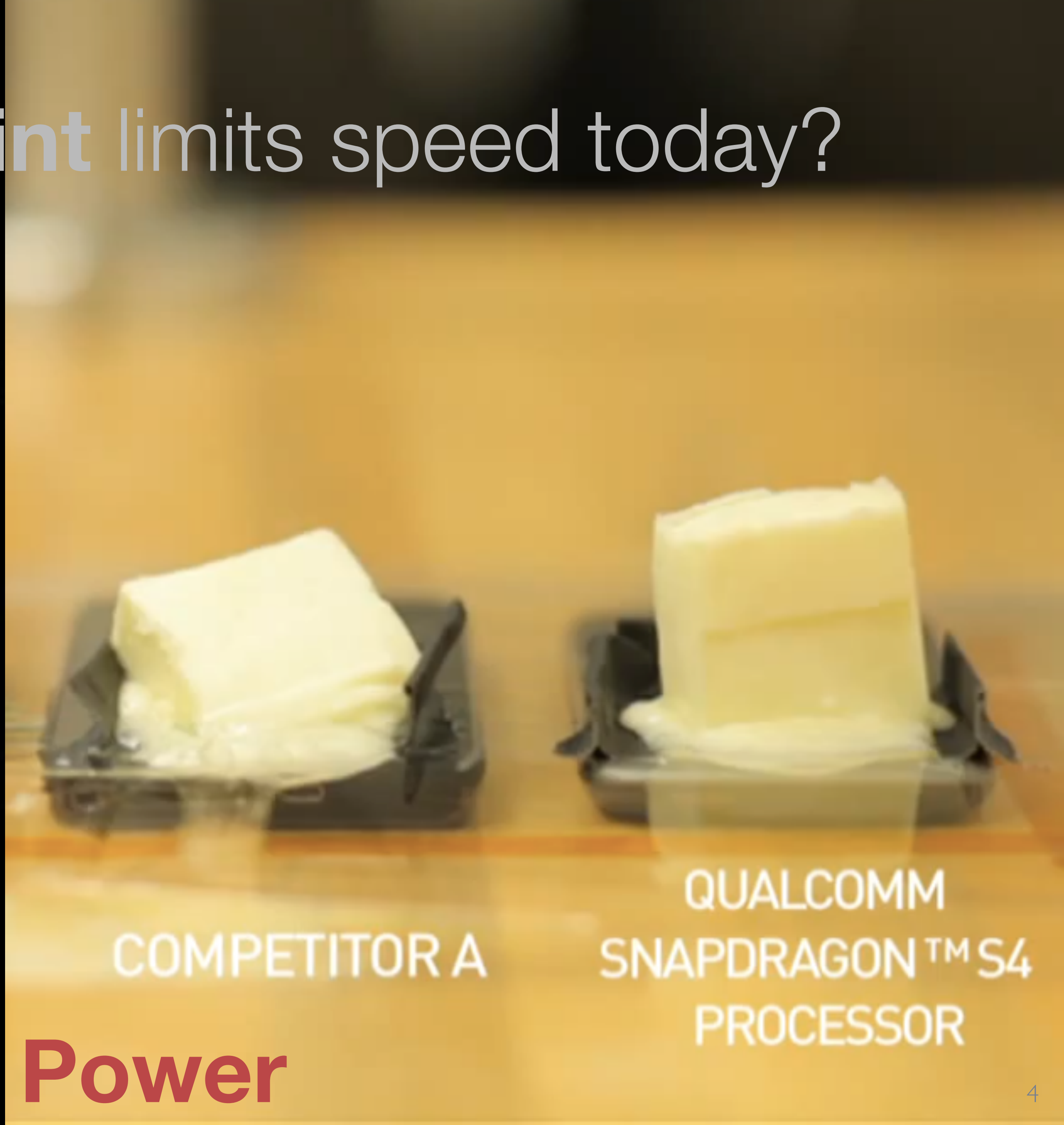QUALCOMM
SNAPDRAGON™ S4
PROCESSOR

4

*(An aside on the relationship between computational performance and power)*

# ImageNet Dataset



Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). Imagenet large scale visual recognition challenge. arXiv preprint arXiv:1409.0575. [web]

(4 GPUs)
x (250 Watts / GPU)
x (1 week)
**~ 0.6 billion Joules**

(4 GPUs)
x (250 Watts / GPU)
x (1 week)
**~ 0.6 billion Joules**

(1 brain)
x (20 Watts / brain)
x (1 year)
**~ 0.6 billion Joules**

(4 GPUs)
x (250 Watts / GPU)
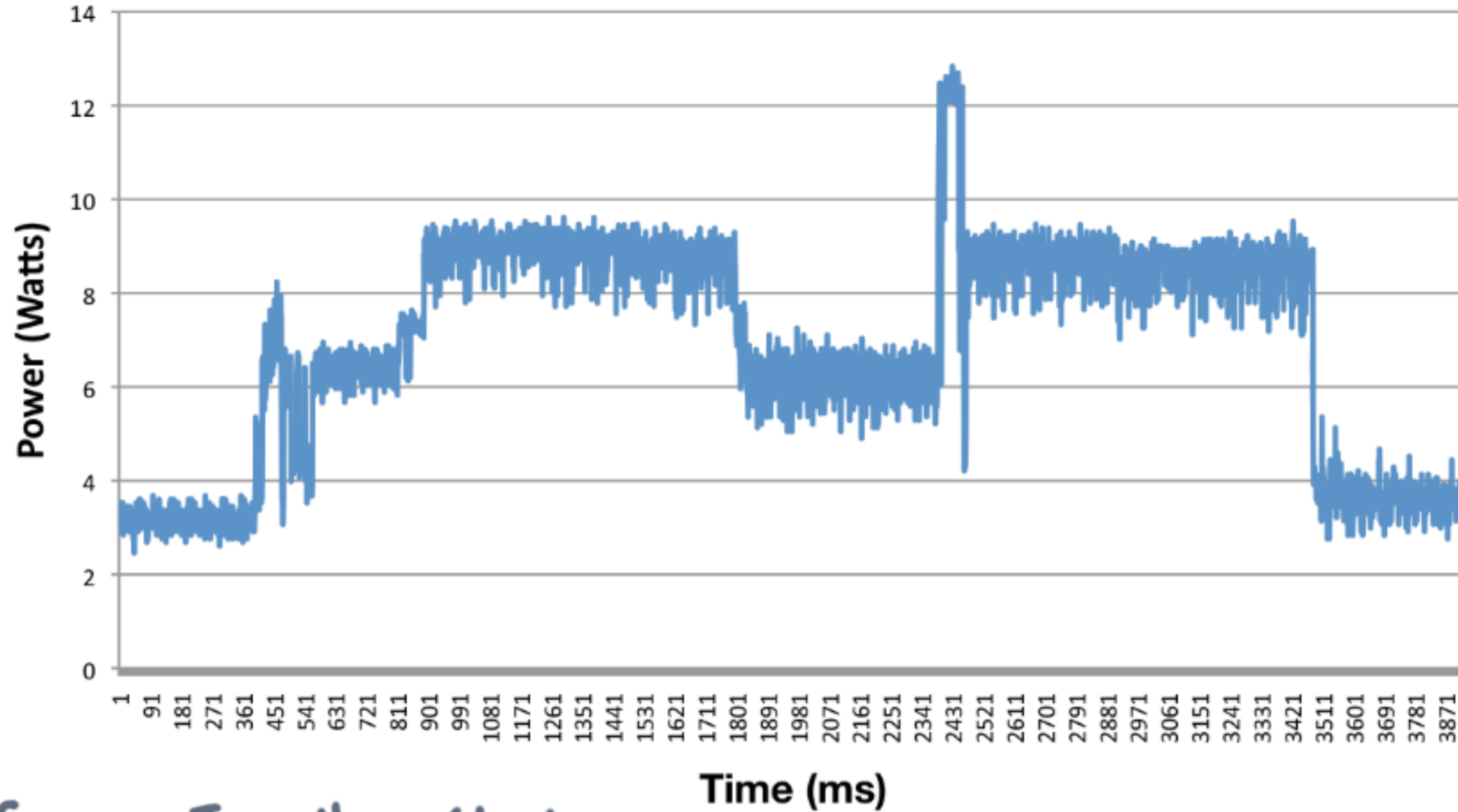x (1 week)
**~ 0.6 billion Joules**

**??**

(1 brain)
x (20 Watts / brain)
x (1 year)
**~ 0.6 billion Joules**

# Power Limits

$$Power \equiv \frac{Energy}{Time}$$



Source: Jee Whan Choi

# Power Limits



Source: Jee Whan Choi

| | Jetson TK1 |
|---|---|
| CPU | ARM A15 (*32-bit, 2.3 GHz, 4+1 cores*) |
| GPU | 192 core Kepler, 326 GF/s (*peak*) |
| Memory | 2 GB LPDDR3 |
| Storage | 16 GB eMMC |
| Networking | Ethernet |
| Form Factor | Dev board |
| I/O | USB, HDMI, Serial |
| Release Date | 2014 |

# Power Limits

$$Power \equiv \frac{Energy}{Time}$$



Source: Jee Whan Choi

# Power Limits

$$Power \equiv \frac{Energy}{Time}$$



Source: Jee Whan Choi

# Limit

*or "cap", from a user or the system*

# Power Limits

$$Power \equiv \frac{Energy}{Time}$$



**Limit**

*or "cap", from a user or the system*

Source: Jee Whan Choi

# Power Limits

$$Power \equiv \frac{Energy}{Time}$$



## Limit

*or "cap", from a user or the system*

Source: Jee Whan Choi

Power Limits

$$Power \equiv \frac{Energy}{Time}$$

**Limit**

*or "cap", from a
user or the system*

**Slow down?**

*Main question of this talk:*

# Can you **design an algorithm** in a way that you can **control its power**?

*We are interested in "algorithmic" methods that **complement** techniques available in hardware, like DVFS, and systems software or middleware.*

*A first principle:*

Relationships among **time**, **energy**, and **power**.

*J. Choi, D. Bedard, R. Fowler, R. Vuduc. **"A roofline model of energy."** In IPDPS'13.*
*J. Choi, M. Dukhan, X. Liu, R. Vuduc. "Algorithmic time, energy, and power on candidate HPC building blocks." In IPDPS'14.*

*Time ~ ?*

*Energy ~ ?*

*Power = Energy / Time*

Time ~ (# of operations) / (number of processors)

Energy ~ (# of operations)

Power = Energy / Time ~ (number of processors) = Speedup

## Time ~ (# of operations) / (number of processors)

Energy ~ (# of operations)

Power = Energy / Time ~ (number of processors) = Speedup

*Time ~ (# of operations) / (number of processors)*

***Energy ~ (# of operations)***

*Power = Energy / Time ~ (number of processors) = Speedup*

Time ~ (# of operations) / (number of processors)

Energy ~ (# of operations)

**Power = Energy / Time ~ (number of processors) = Speedup**

*Time ~ (# of operations) / (number of processors)*

*Energy ~ (# of operations)*

*Power = Energy / Time ~ (number of processors) = Speedup*

Conclusion:
To save time & energy: *Must reduce work (# ops or cost/op)*
To save power: *Must slow down (e.g., use fewer cores)*

*Time ~* **(# of operations)** *(/ (number of processors)*

*Energy ~* **(# of operations)**

*Power = Energy / Time ~ (number of processors) = Speedup*

Conclusion:

To save time & energy: *Must* **reduce work (# ops** *or* **cost/op)**
To save power: *Must slow down (e.g., use fewer cores)*

23

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

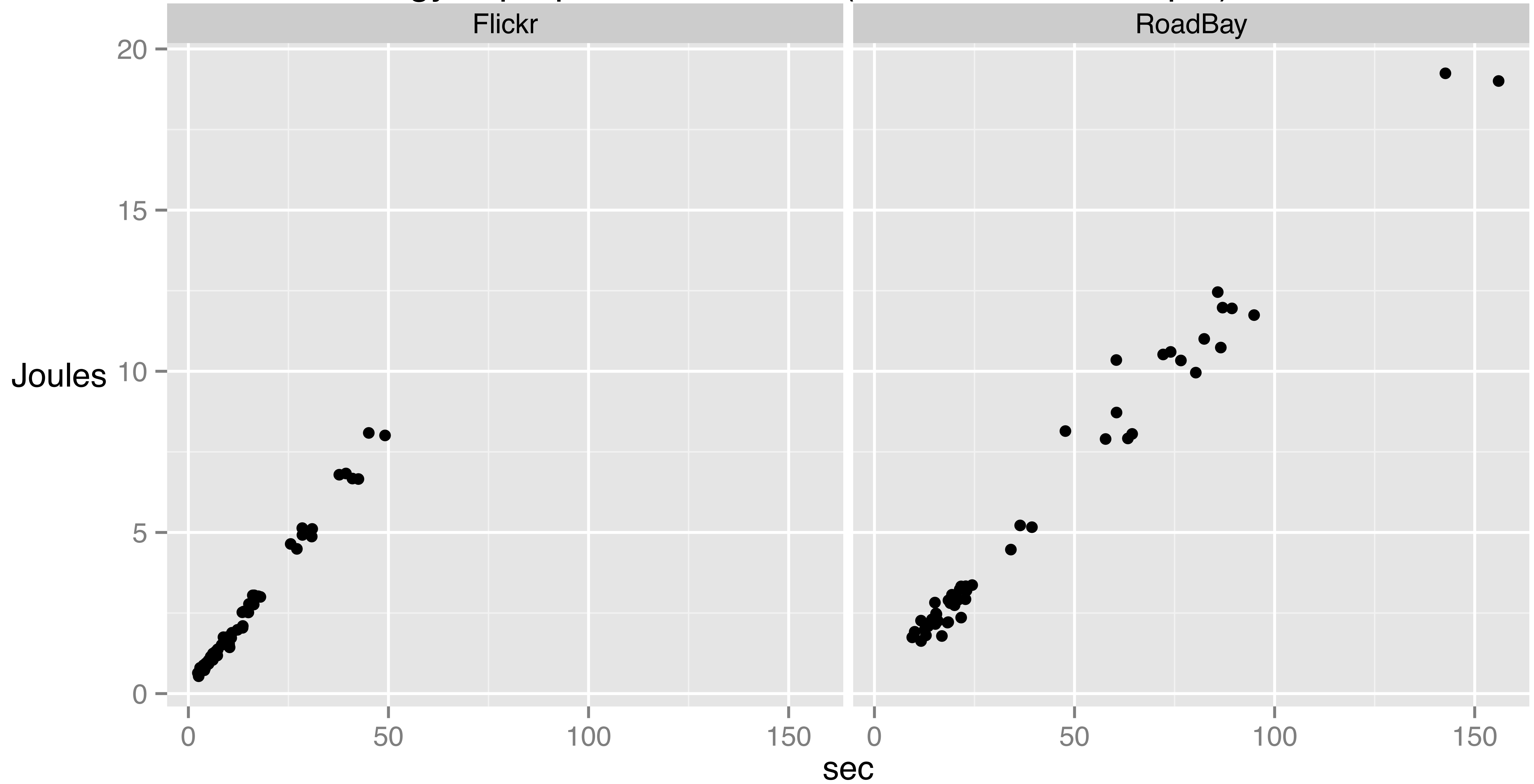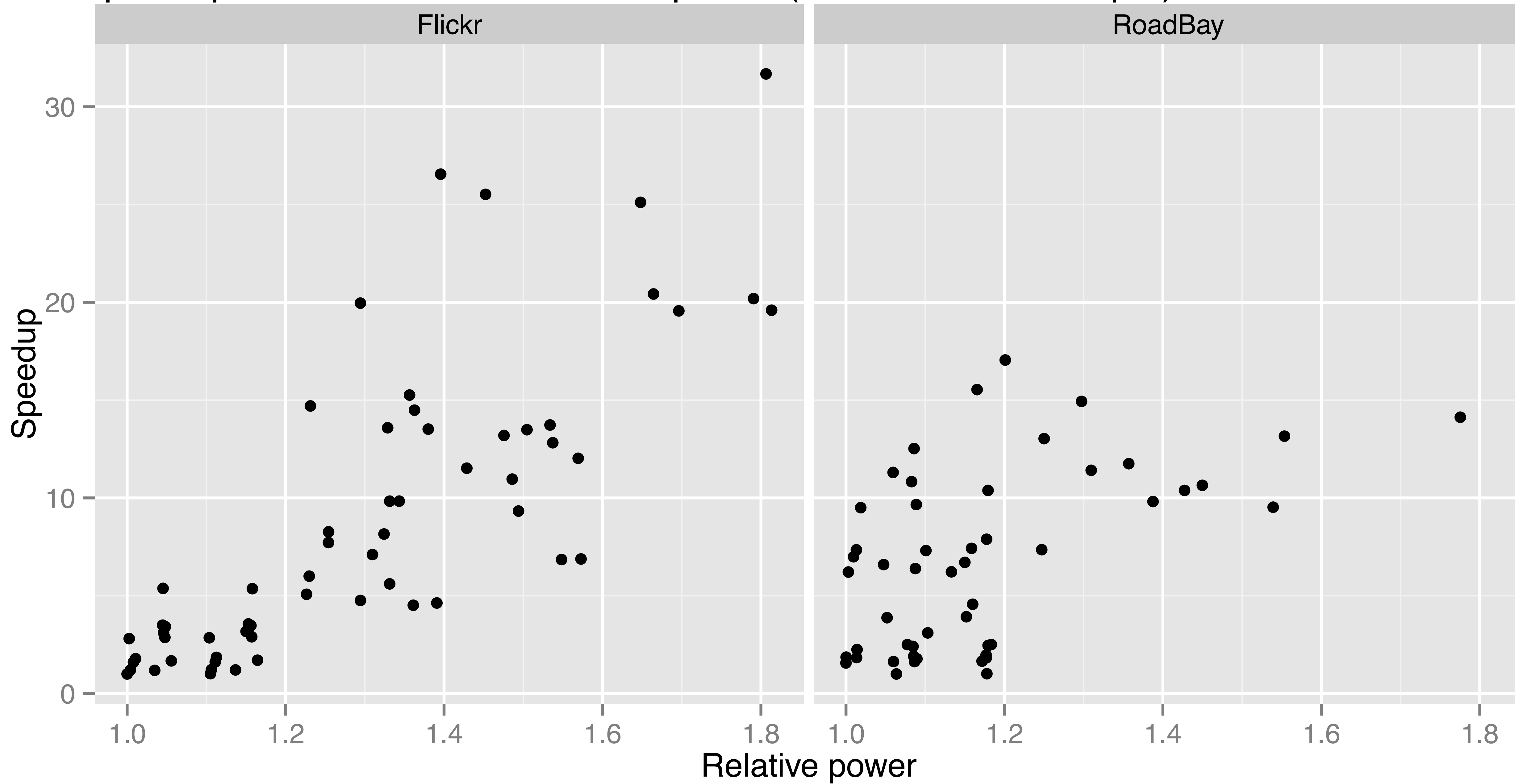*Time ~ (# of operations) / (number of processors)*

*Energy ~ (# of operations)*

*Power = Energy / Time ~* **(number of processors) = Speedup**

Conclusion:
To save time & energy: *Must reduce work (# ops or cost/op)*
To save power: *Must* **slow down** *(e.g., use fewer cores)*

25

# Speedup increases with additional power (SSSP+GPU example)

*Main question of this talk:*

# Can you **design an algorithm** in a way that you can **control its power**?

*We are interested in "algorithmic" methods that **complement** techniques available in hardware, like DVFS, and systems software or middleware.*

**Yes!**
Example: A **power-tunable** graph algorithm to compute single-source shortest paths **(SSSP)**.

*Sara Karamati (Ph.D. student), Dr. Jeff Young (research faculty), R. Vuduc — new,* **unpublished** *work*

- Baseline: Fastest, work-efficient **"delta-stepping-like"** method*

  - **Tunable work-parallelism tradeoff**

- Tuned for a GPU and run on an **NVIDIA Jetson TK1**, which has tunable core frequencies (**10x**) and memory frequencies (**3x**)
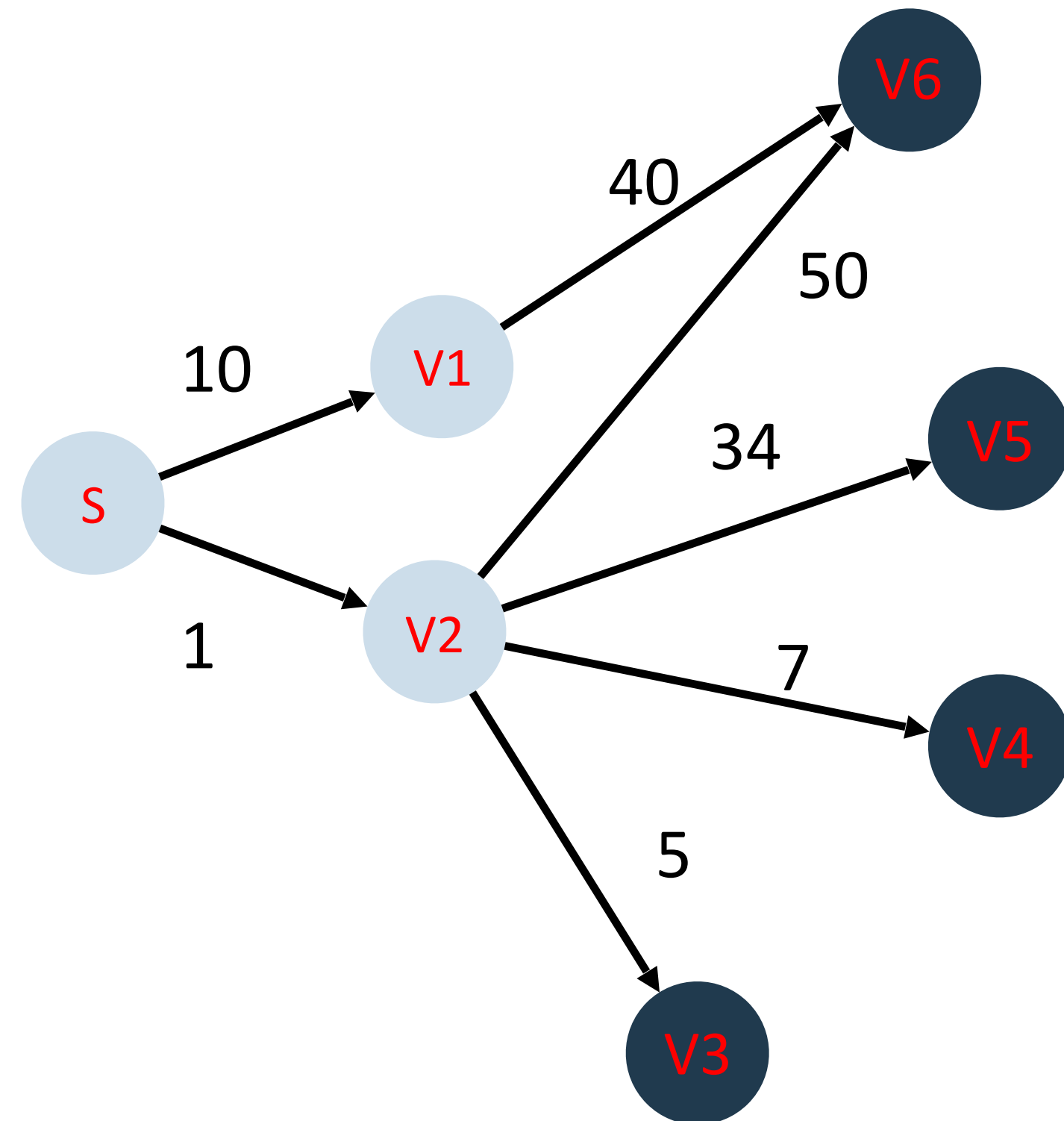
- *No preprocessing shortcuts, a la PHAST***

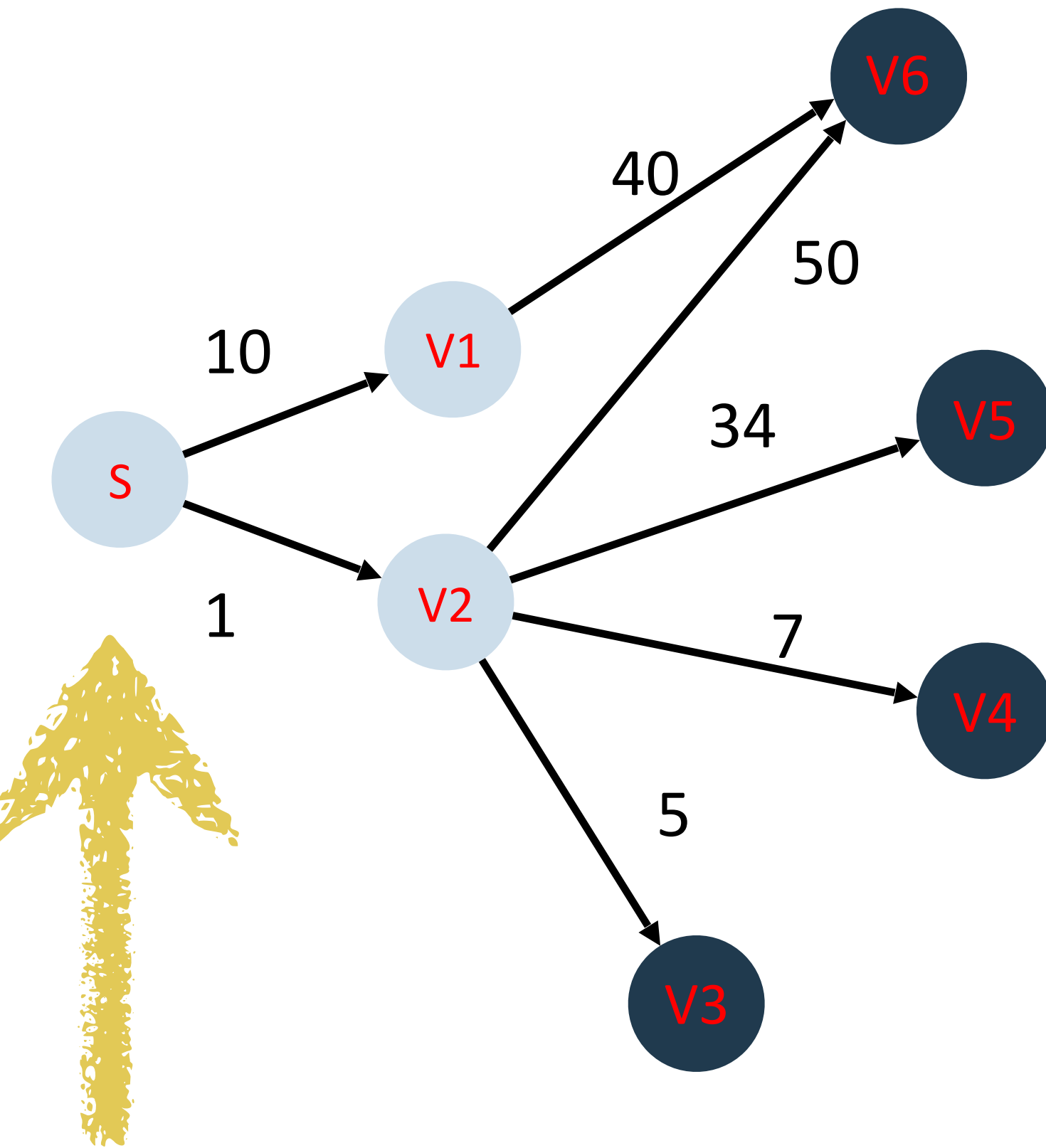| | Jetson TK1 |
|---|---|
| CPU | ARM A15 (*32-bit, 2.3 GHz, 4+1 cores*) |
| GPU | 192 core Kepler, 326 GF/s (*peak*) |
| Memory | 2 GB LPDDR3 |
| Storage | 16 GB eMMC |
| Networking | Ethernet |
| Form Factor | Dev board |
| I/O | USB, HDMI, Serial |
| Release Date | 2014 |

*\* Based on GunRock implementations of Davidson, Baxter, Garland, and Owens (IPDPS'14)*
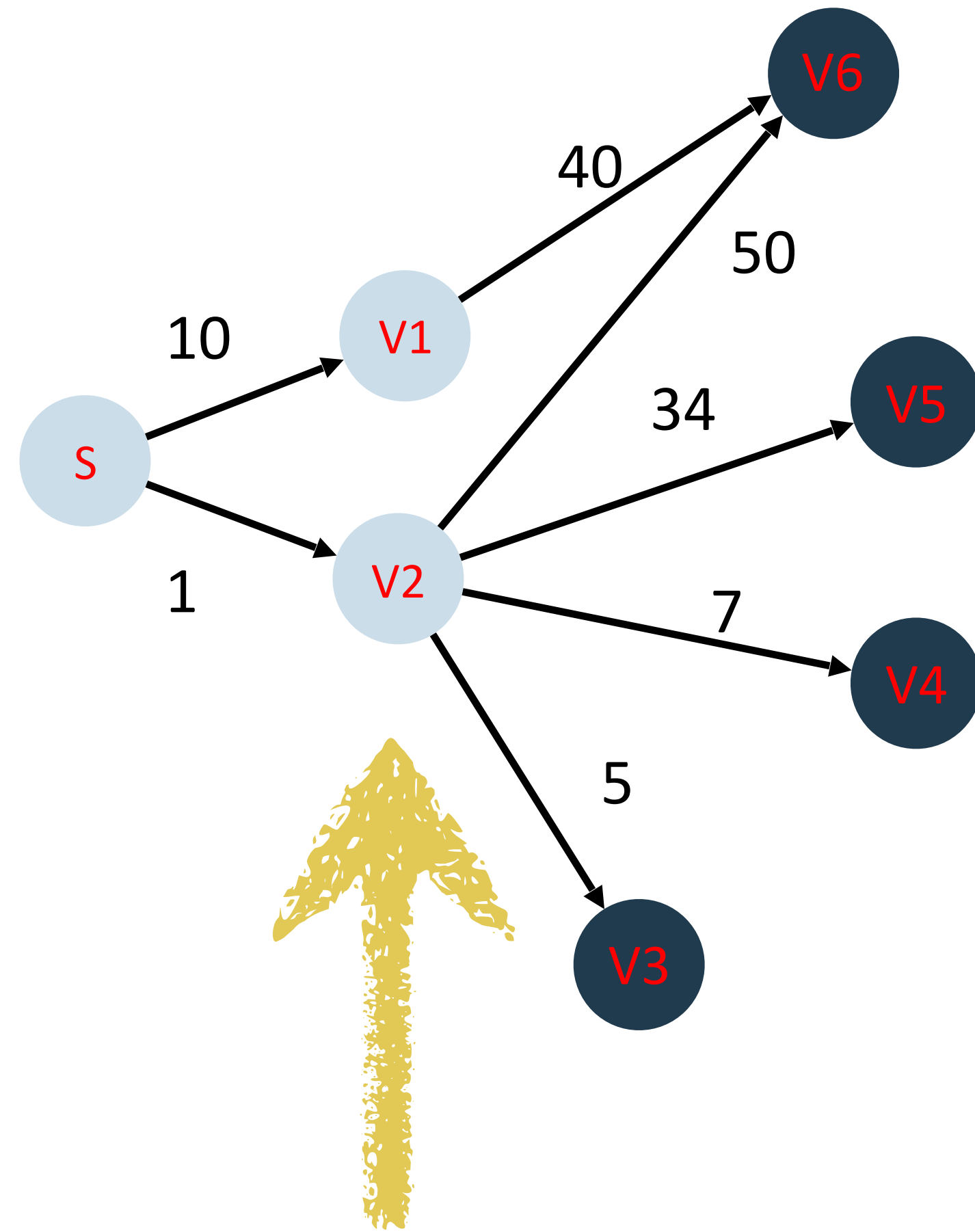*\*\* Delling et al. "PHAST: Hardware-accelerated shortest path trees" (JPDC'10)*

# Baseline: Gunrock's **"Near+Far"** algorithm

# Baseline: Gunrock's "**Near+Far**" algorithm

# Baseline: Gunrock's "**Near+Far**" algorithm

# Baseline: Gunrock's "**Near+Far**" algorithm

Splitting Distance = 20

| | V3 | V3, d=6 < 20 | | |
|---|---|---|---|---|
| V1 | V4 | V4, d=8 < 20 | V3 | V5 |
| V2 | V5 | V5, d=35 > 20 | V4 | V6 |
| | V6 | V6, d=51 > 20 | | |
| | V6 | | | |

Advance → Filter → Bisect

Frontier    Frontier    Frontier

One Iteration

Far Pile

Near Pile = Frontier

31

# Baseline: Gunrock's "**Near+Far**" algorithm

Splitting Distance = 20

| | | | | |
|---|---|---|---|---|
| V1<br>V2 | **Advance** → | V3<br>V4<br>V5<br>V6<br>V6 | **Filter** → | V3, d=6 < 20<br>V4, d=8 < 20<br>V5, d=35 > 20<br>V6, d=51 > 20 |

Frontier | Frontier | Frontier

**Bisect** →

V5<br>V6 — Far Pile

V3<br>V4 — Near Pile = Frontier

One Iteration

**|Frontier| ~ parallelism**

31

# Baseline: Gunrock's "**Near+Far**" algorithm

Splitting Distance = 20

| | |
|---|---|
| 40 | |
| 50 | |
| 10 | |
| 34 | |
| 1 | |
| 7 | |
| 5 | |

**Frontier**: V1, V2 → **Advance** → **Frontier**: V3, V4, V5, V6, V6 → **Filter** → **Frontier**: V3, d=6 < 20 / V4, d=8 < 20 / V5, d=35 > 20 / V6, d=51 > 20 → **Bisect** → V3, V4

Far Pile: V5, V6

Near Pile = Frontier: V3, V4

One Iteration

31

# Baseline: Gunrock's "**Near+Far**" algorithm

Splitting Distance = 20

**Frontier:** V1, V2

Advance →

**Frontier:** V3, V4, V5, V6, V6

Filter →

**Frontier:**
V3, d=6 < 20
V4, d=8 < 20
V5, d=35 > 20
V6, d=51 > 20

Bisect →

Near Pile = Frontier: V3, V4

Far Pile: V5, V6

One Iteration

31

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ ⬇**

Splitting Distance = 20

| V6 |
|----|

40

50

10

| V1 |

34 → V5

| S |

1

| V2 |

7 → V4

5

| V3 |

**Far Pile**

V5
V6

**Near Pile = Frontier**

| V1 V2 | → Advance → | V3 V4 V5 V6 V6 | → Filter → | V3, d=6 < 20 V4, d=8 < 20 V5, d=35 > 20 V6, d=51 > 20 | → Bisect → | V3 V4 |

Frontier     Frontier     Frontier

One Iteration

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ**

Splitting Distance = 20

Frontier: V1, V2 → Advance → Frontier: V3, V4, V5, V6, V6 → Filter → Frontier:
V3, d=6 < 20
V4, d=8 < 20
V5, d=35 > 20
V6, d=51 > 20
→ Bisect → Near Pile = Frontier: V3, V4

Far Pile: V5, V6

One Iteration

31

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ**

Splitting Distance = 20

V1
V2

**Frontier**

Advance →

V3
V4
V5
V6
V6

**Frontier**

Filter →

V3, d=6 < 20
V4, d=8 < 20
V5, d=35 > 20
V6, d=51 > 20

**Frontier**

Bisect →

V3
V4

Far Pile

V5
V6

Near Pile = Frontier

One Iteration

Graph edges:
S → V1: 10
S → V2: 1
V1 → V6: 40
V2 → V6: 50
V2 → V5: 34
V2 → V4: 7
V2 → V3: 5

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ**

Splitting Distance = 20

Graph edges:
- S → V1: 10
- S → V2: 1
- V1 → V6: 40
- V2 → V6: 50
- V2 → V5: 34
- V2 → V4: 7
- V2 → V3: 5

**Frontier:** V1, V2

**Advance →**

**Frontier:** V3, V4, V5, V6, V6

**Filter →**

**Frontier:**
- V3, d=6 < 20
- V4, d=8 < 20
- V5, d=35 > 20
- V6, d=51 > 20

**Bisect →**

**Far Pile:** V5, V6

**Near Pile = Frontier:** V3, V4

One Iteration

31

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ**

Splitting Distance = 20

Graph edges:
- S → V1: 10
- S → V2: 1
- V1 → V6: 40
- V2 → V6: 50
- V2 → V5: 34
- V2 → V4: 7
- V2 → V3: 5

Frontier (V1, V2) — **Advance** → Frontier (V3, V4, V5, V6, V6) — **Filter** →

Frontier:
- V3, d=6 < 20
- V4, d=8 < 20
- V5, d=35 > 20
- V6, d=51 > 20

— **Bisect** →

Near Pile = Frontier: V3, V4

Far Pile: V5, V6

# Baseline: Gunrock's "**Near+Far**" algorithm

Graph edges:
- S → V1: 10
- S → V2: 1
- V1 → V6: 40
- V2 → V6: 50
- V2 → V5: 34
- V2 → V4: 7
- V2 → V3: 5

**delta: δ**

Splitting Distance = 20

**Frontier**
V1
V2

→ Advance →

**Frontier**
V3
V4
V5
V6
V6

→ Filter →

**Frontier**
V3, d=6 < 20
V4, d=8 < 20
V5, d=35 > 20
V6, d=51 > 20

→ Bisect →

Near Pile = Frontier
V3
V4

Far Pile
V5
V6

One Iteration

# Power ~ Parallelism ~ Queue sizes

32

# Baseline: Gunrock's "**Near+Far**" algorithm

**delta: δ**

Splitting Distance = 20

| | |
|---|---|
| 40 | |
| 50 | V6 |
| 10 | V1 |
| 34 | V5 |
| 1 | V2 |
| 7 | V4 |
| 5 | V3 |

S

**Frontier**
- V1
- V2

→ **Advance** →

**Frontier**
- V3
- V4
- V5
- V6
- V6

→ **Filter** →

**Frontier**
- V3, d=6 < 20
- V4, d=8 < 20
- V5, d=35 > 20
- V6, d=51 > 20

→ **Bisect** →

**Near Pile = Frontier**
- V3
- V4

**Far Pile**
- V5
- V6

One Iteration

# Power ~ Parallelism ~ Queue sizes

32

# What is the effect of **delta (δ)**?

# What is the effect of **delta (δ)?**

Delta = 1e6

*(road network)*

34

# What is the effect of **delta (δ)?**

Delta = 1e6

Delta = 1e5

**Time**

*(road network)*

**Mean power**

**Max power**
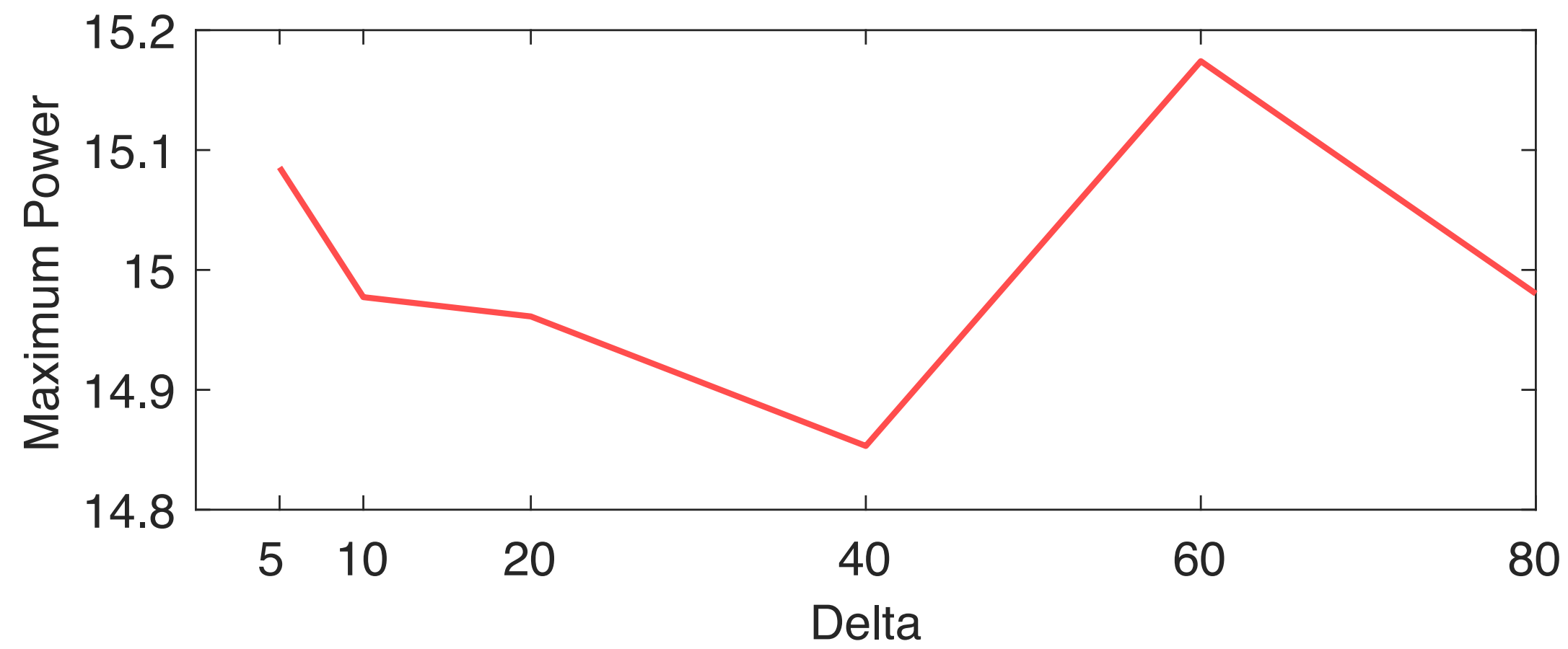
**delta: δ**

36

**Time**  *(scale-free)*

**Mean power**

**Max power**

37

*Observation:*

Delta ($\delta$) is a **tuning parameter** that can be used to control power-time tradeoffs.

*But* **how** *to choose it? It is input-dependent. And, in the literature, it is always treated as a fixed* a priori *parameter with little guidance on its ideal value.*
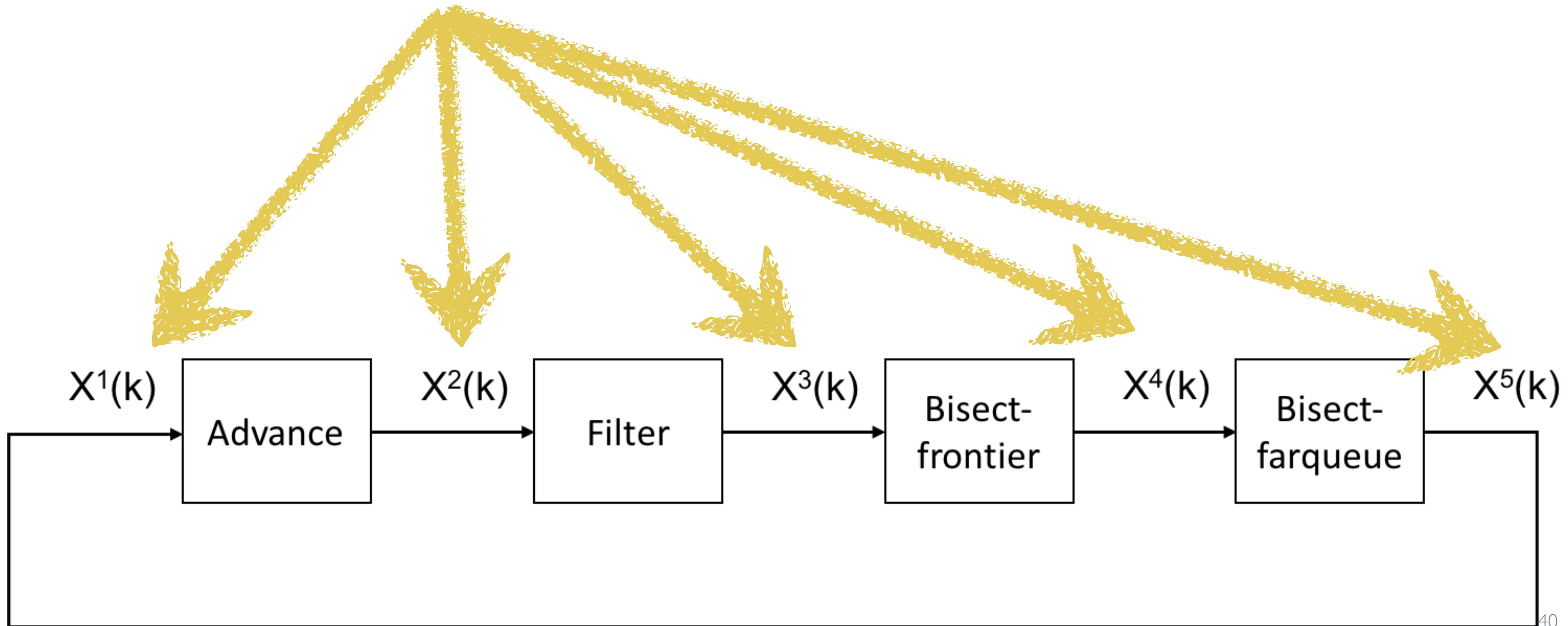
*Sara's insight:*

Treat $\delta$ as a parameter to be **learned** and **controlled**, **dynamically**.
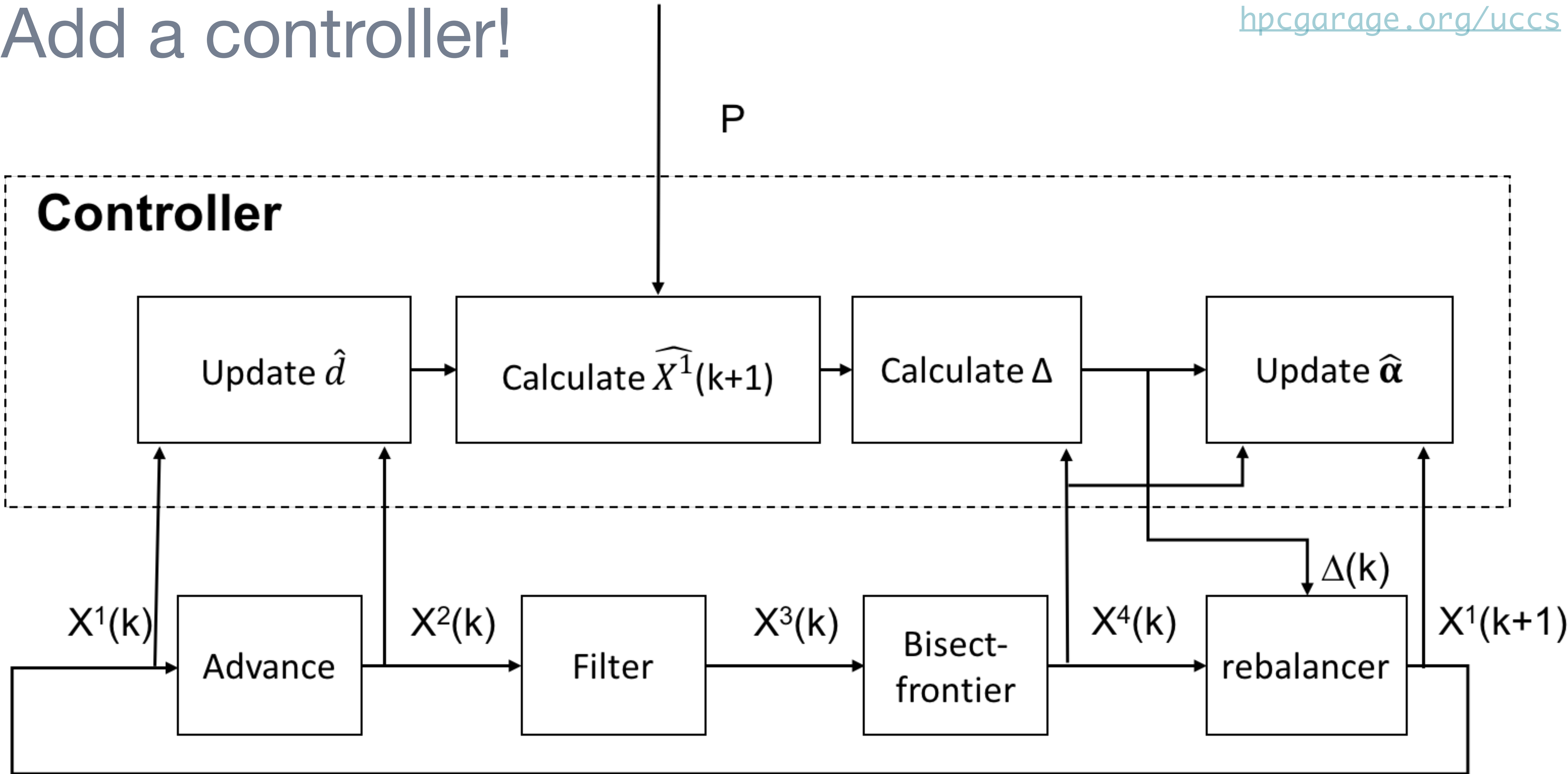
# Recall: Near+Far == stages.

$X^1(k)$ → **Advance** → $X^2(k)$ → **Filter** → $X^3(k)$ → **Bisect-frontier** → $X^4(k)$ → **Bisect-farqueue** → $X^5(k)$

# Recall: Near+Far == stages.

## Intermediate frontier (queue) sizes



$X^1(k)$ → **Advance** → $X^2(k)$ → **Filter** → $X^3(k)$ → **Bisect-frontier** → $X^4(k)$ → **Bisect-farqueue** → $X^5(k)$

# Add a controller!

P



**Controller**

| Update $\hat{d}$ | Calculate $\widehat{X^1}(k+1)$ | Calculate $\Delta$ | Update $\widehat{\boldsymbol{\alpha}}$ |

$X^1(k)$ → Advance → $X^2(k)$ → Filter → $X^3(k)$ → Bisect-frontier → $X^4(k)$ → rebalancer → $X^1(k+1)$

$\Delta(k)$

41

# Simple models between stages...

**Controller**

Update $\hat{d}$

Calculate $\widehat{X^1}(k$

$\hat{\alpha}$

**Try to estimate ("learn")**

$X^1(k)$    Advance    $X^2(k)$

Filter    $X^3(k)$    Bisect-frontier    $X^4(k)$    $\downarrow \Delta(k)$    rebalancer    $X^1(k+1)$

$X^2(k) \sim$ **(degree)** $* X^1(k)$

42

$$\hat{X}_k^{(2)} = d \cdot X_k^{(1)}$$

# Estimator

$$\hat{X}_k^{(2)} = d \cdot X_k^{(1)}$$

# Estimator

$$\hat{X}_k^{(2)} = d \cdot X_k^{(1)}$$

**Parameter**

**Estimator**

**Loss**

$$\hat{X}_k^{(2)} = d \cdot X_k^{(1)}$$

$$\min_d \sum_k \left( X_k^{(2)} - \hat{X}_k^{(2)} \right)^2$$

**Parameter**

**Estimator**

**Loss**

$$\hat{X}_k^{(2)} = d \cdot X_k^{(1)}$$

**Parameter**

$$\min_d \sum_k \left( X_k^{(2)} - \hat{X}_k^{(2)} \right)^2$$

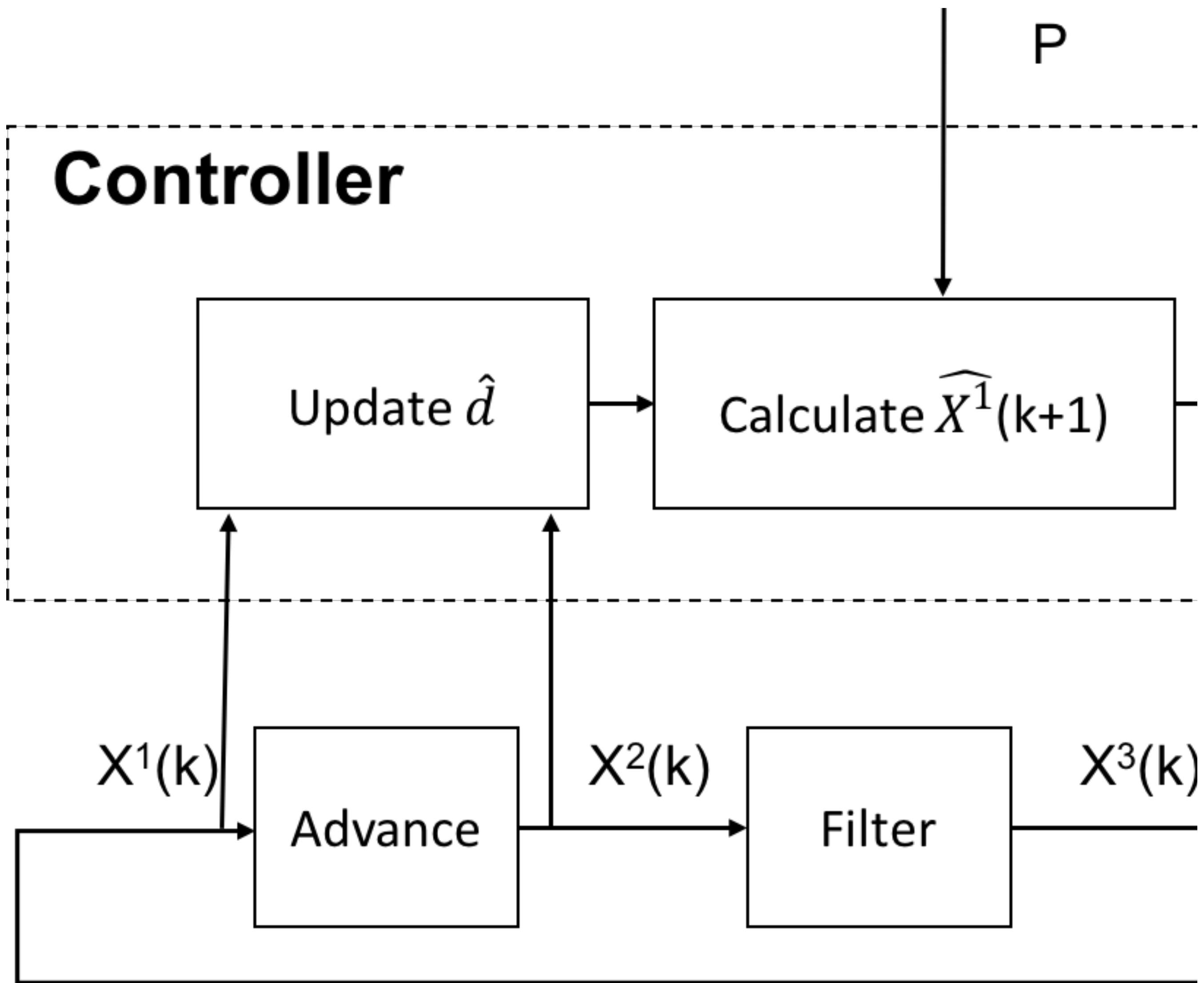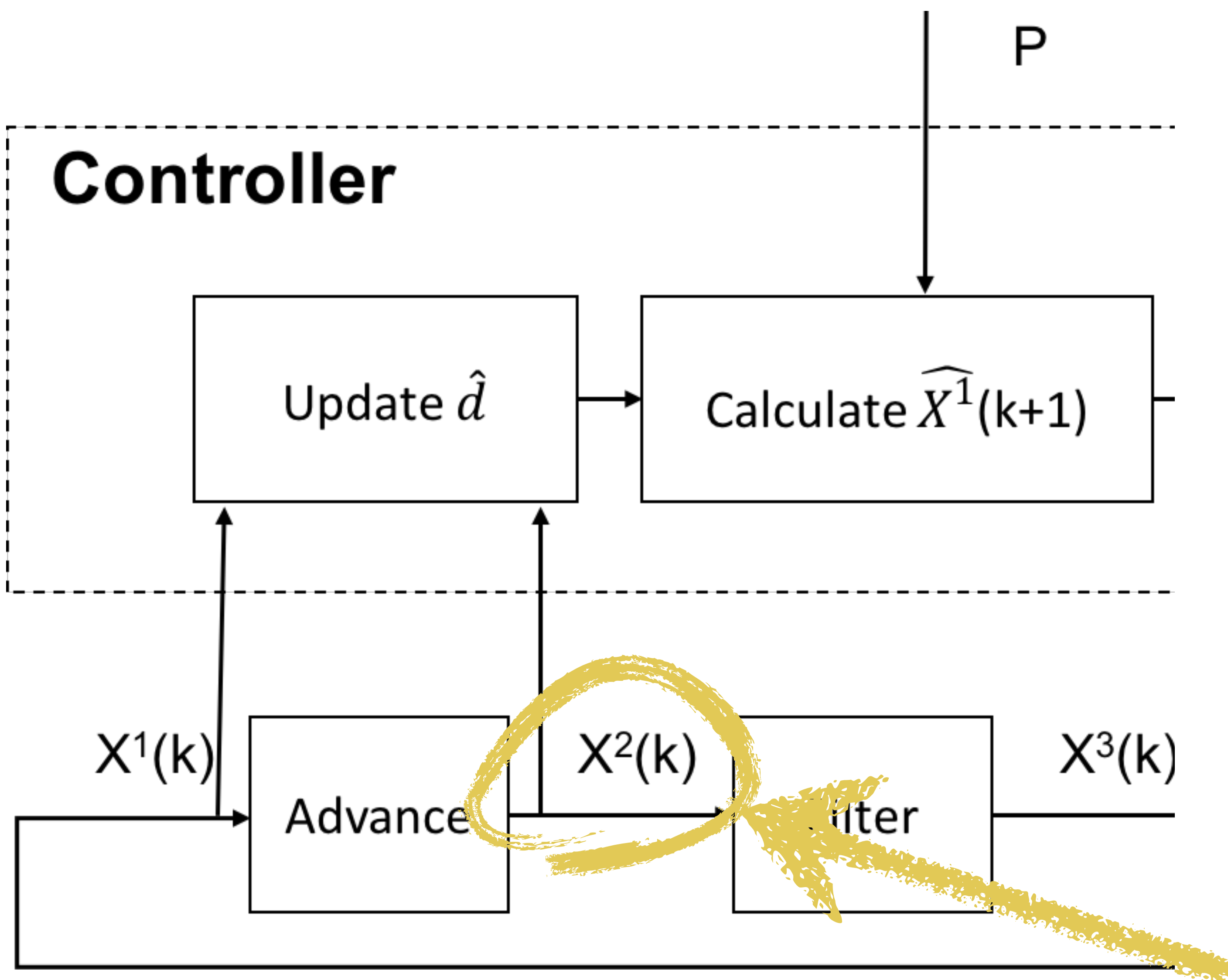$$= \min_d \sum_k \left( X_k^{(2)} - d \cdot X_k^{(1)} \right)^2$$

# **Loss**

$$\min_d \sum_k \left( X_k^{(2)} - d \cdot X_k^{(1)} \right)^2 - \lambda \ln d$$

Regularize to stabilize the estimator during the early iterations and use online fitting method (e.g., stochastic gradient descent)
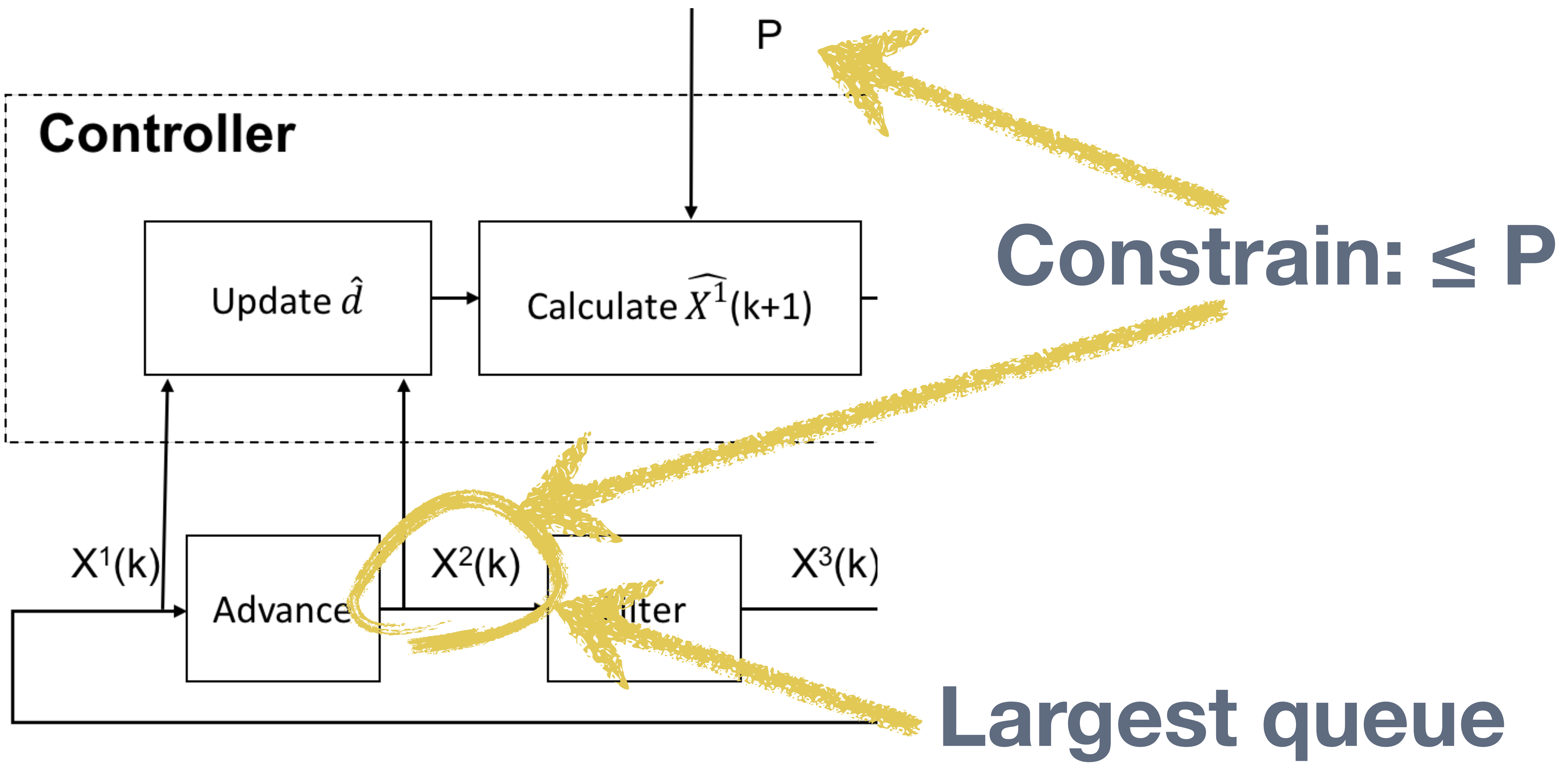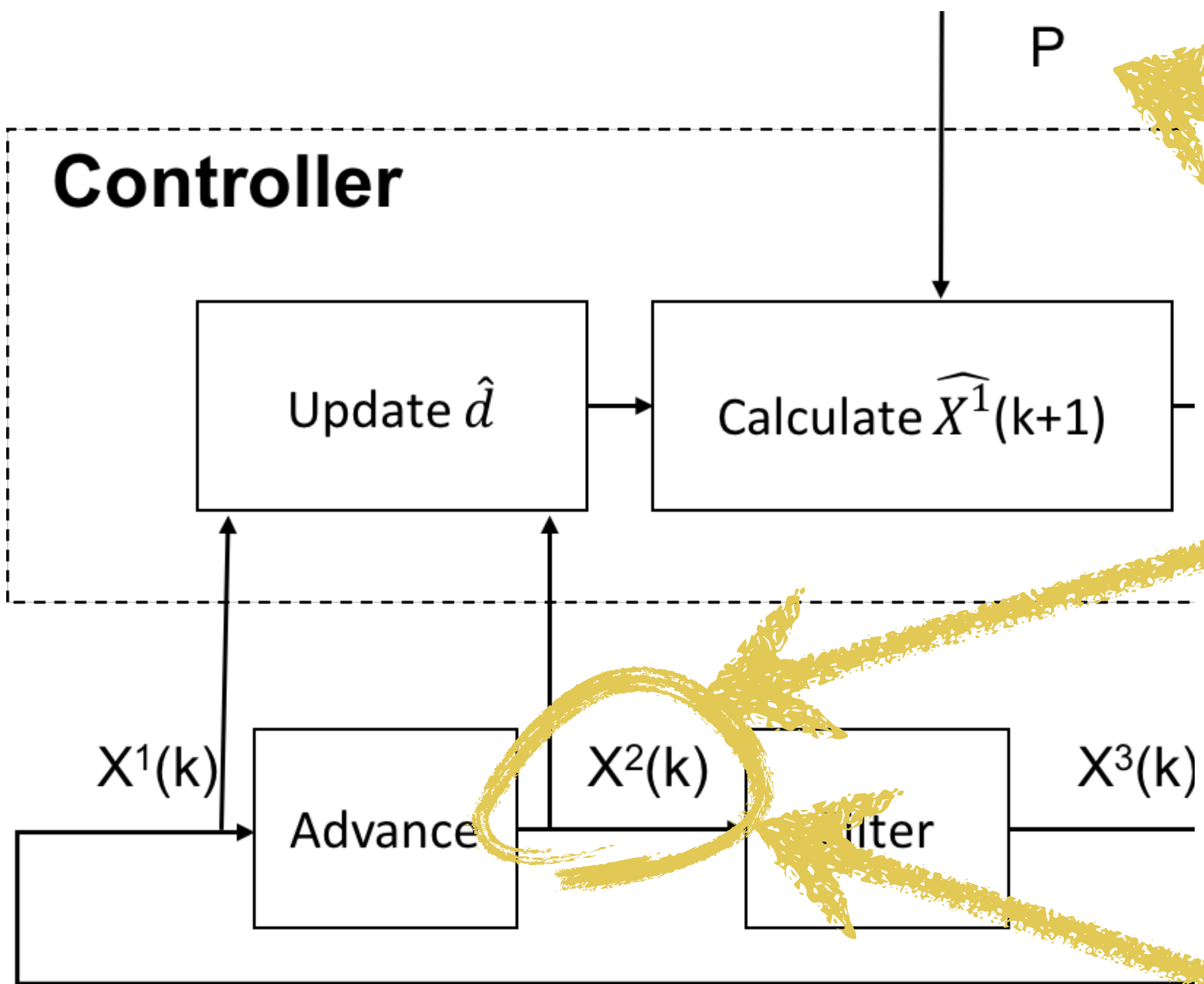
# Constrain parallelism…

P

**Controller**

| Update $\hat{d}$ | Calculate $\widehat{X^1}(k+1)$ |

$X^1(k)$

Advance

$X^2(k)$

Filter

$X^3(k)$

# Constrain parallelism…

P

**Controller**

Update $\hat{d}$ → Calculate $\widehat{X^1}(k+1)$

$X^1(k)$

Advance

$X^2(k)$

...lter

$X^3(k)$

**Largest queue**

# Constrain parallelism...

**Controller**

Update $\hat{d}$ → Calculate $\widehat{X^1}(k+1)$

P

**Constrain: ≤ P**

$X^1(k)$    Advance    $X^2(k)$    ...ter    $X^3(k)$

**Largest queue**

45

# Constrain parallelism…

P

**Controller**

| Update $\hat{d}$ | | Calculate $\widehat{X^1}(k+1)$ |

**Constrain: ≤ P**

$X^1(k)$ | Advance | $X^2(k)$ | ...lter | $X^3(k)$

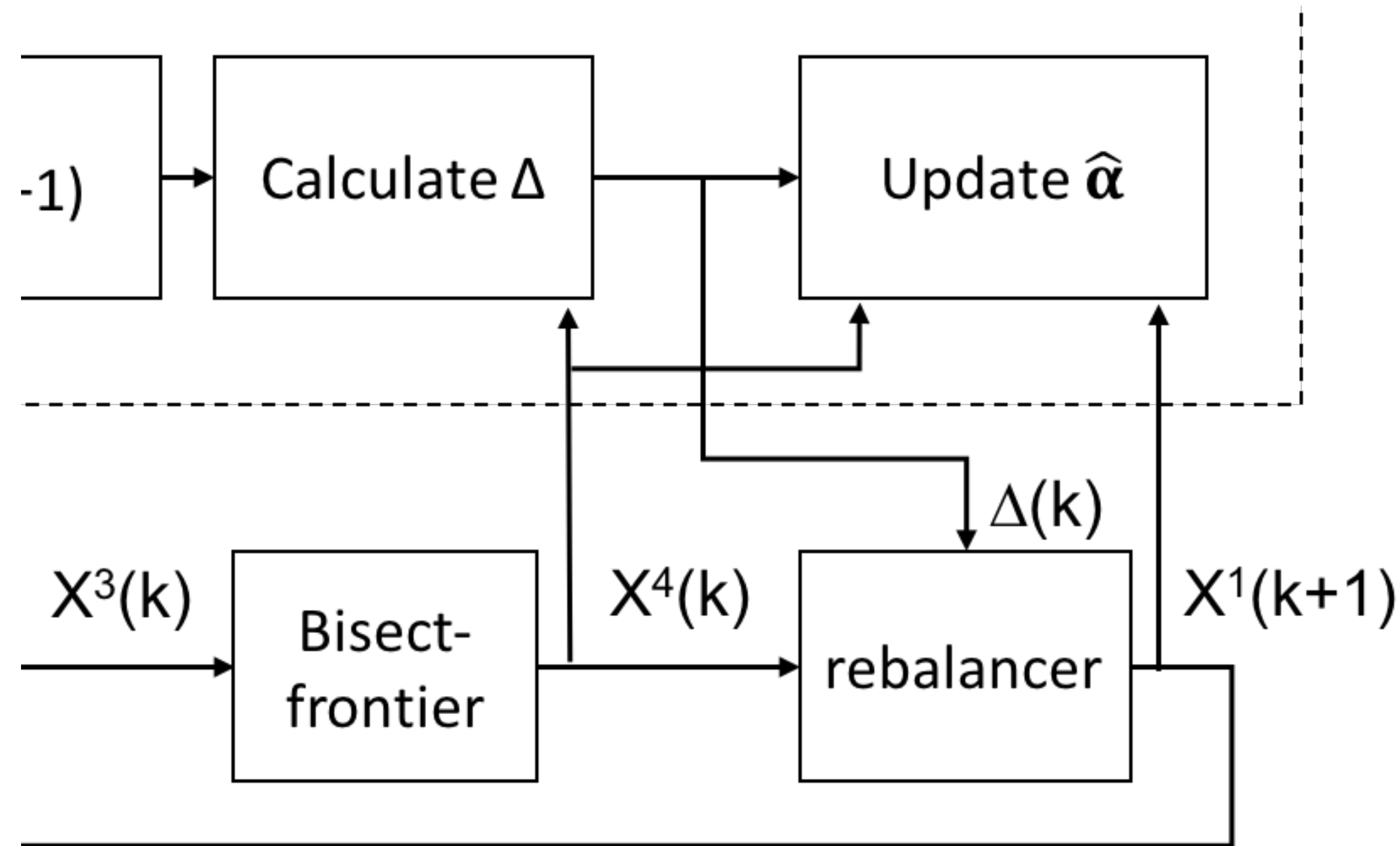$$\hat{X}^{(1)}_{k+1} \leq \frac{P}{d}$$

**Largest queue**

45

# Estimate the effect of a *change* in $\boldsymbol{\delta}$

$$\hat{X}_{k+1}^{(1)} = X_k^{(4)} + \alpha \cdot \Delta\delta_k$$

# Estimate the effect of a *change* in **δ**

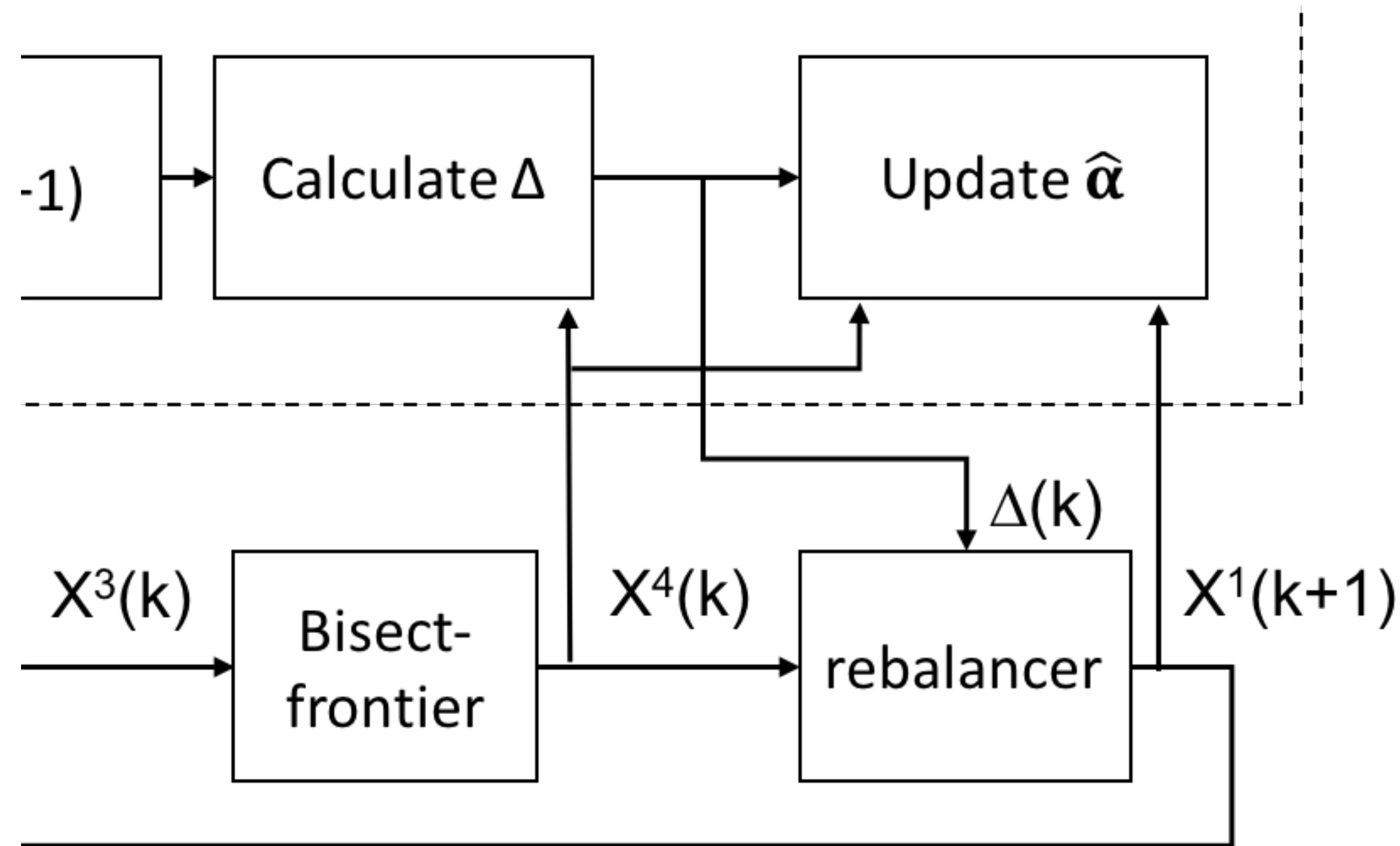$$\hat{X}^{(1)}_{k+1} = X^{(4)}_k + \alpha \cdot \Delta\delta_k$$

**Estimator**

| | Calculate Δ | | Update $\widehat{\boldsymbol{\alpha}}$ |
|---|---|---|---|
| -1) | | | |

$\Delta(k)$

$X^3(k)$ → Bisect-frontier → $X^4(k)$ → rebalancer → $X^1(k+1)$

# Estimate the effect of a *change* in $\boldsymbol{\delta}$

$$\hat{X}_{k+1}^{(1)} = X_k^{(4)} + \alpha \cdot \Delta\delta_k$$

**Estimator**

**Parameter**

| | Calculate Δ | | Update $\widehat{\boldsymbol{\alpha}}$ |

-1) → Calculate Δ → Update $\widehat{\boldsymbol{\alpha}}$

$X^3(k)$ → Bisect-frontier → $X^4(k)$ → rebalancer

$\Delta(k)$

$X^1(k+1)$

*Cool feature:*

# User sets $P$, the max parallelism, not $\delta$!

*The controller chooses the hard-to-set $\delta$ fully automatically, adjusting it as frequently as **every iteration**. <u>Not discussed</u>: One can prove that the controller is bounded-input bounded-output (BIBO) stable, meaning frontier sizes will be "well-behaved."*

*Next question:*
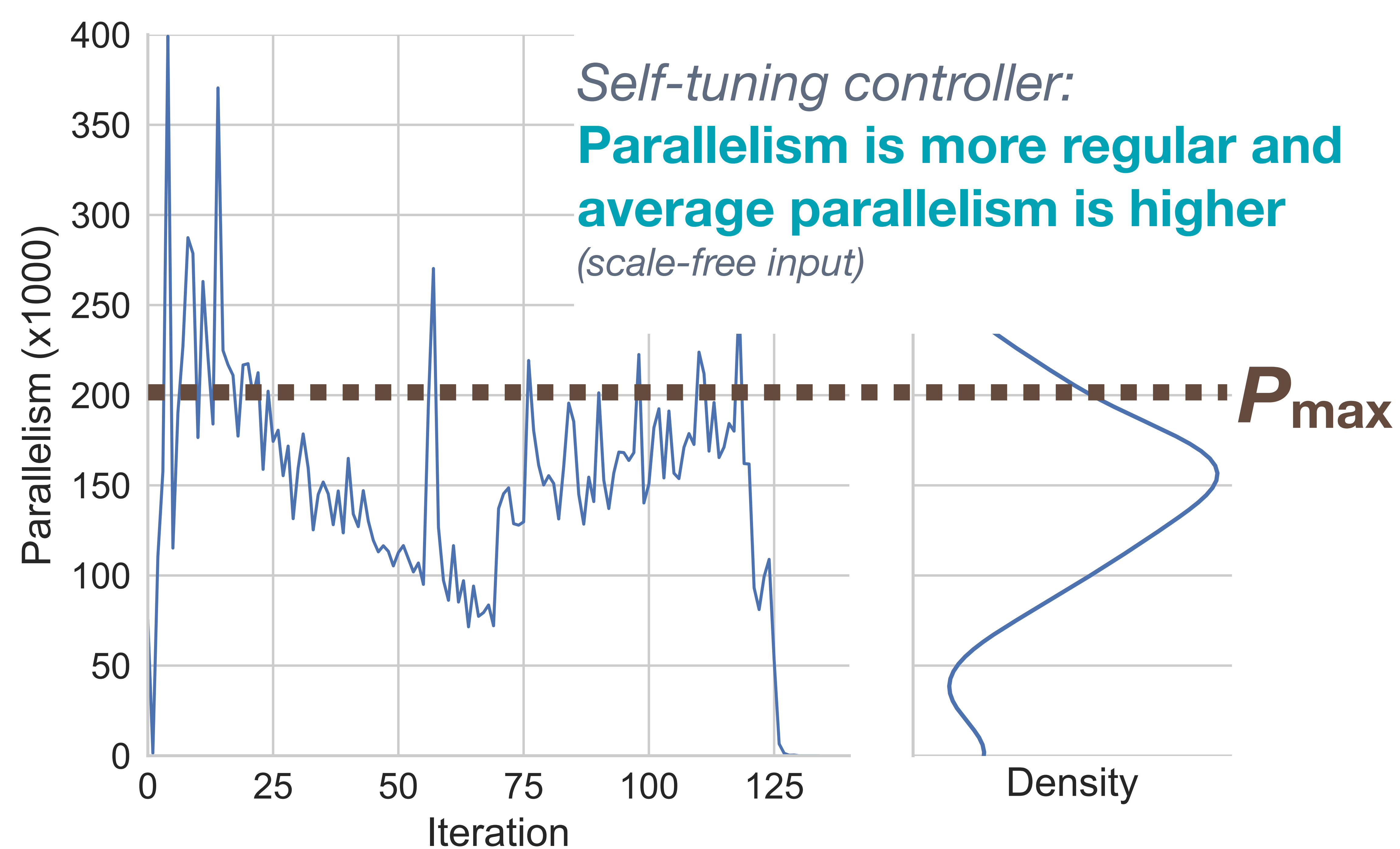
# Does it work?

*(Experimental results)*

*Baseline Near+Far:*
**Parallelism is highly irregular, average parallelism is low**
*(scale-free input)*

Baseline Near+Far:
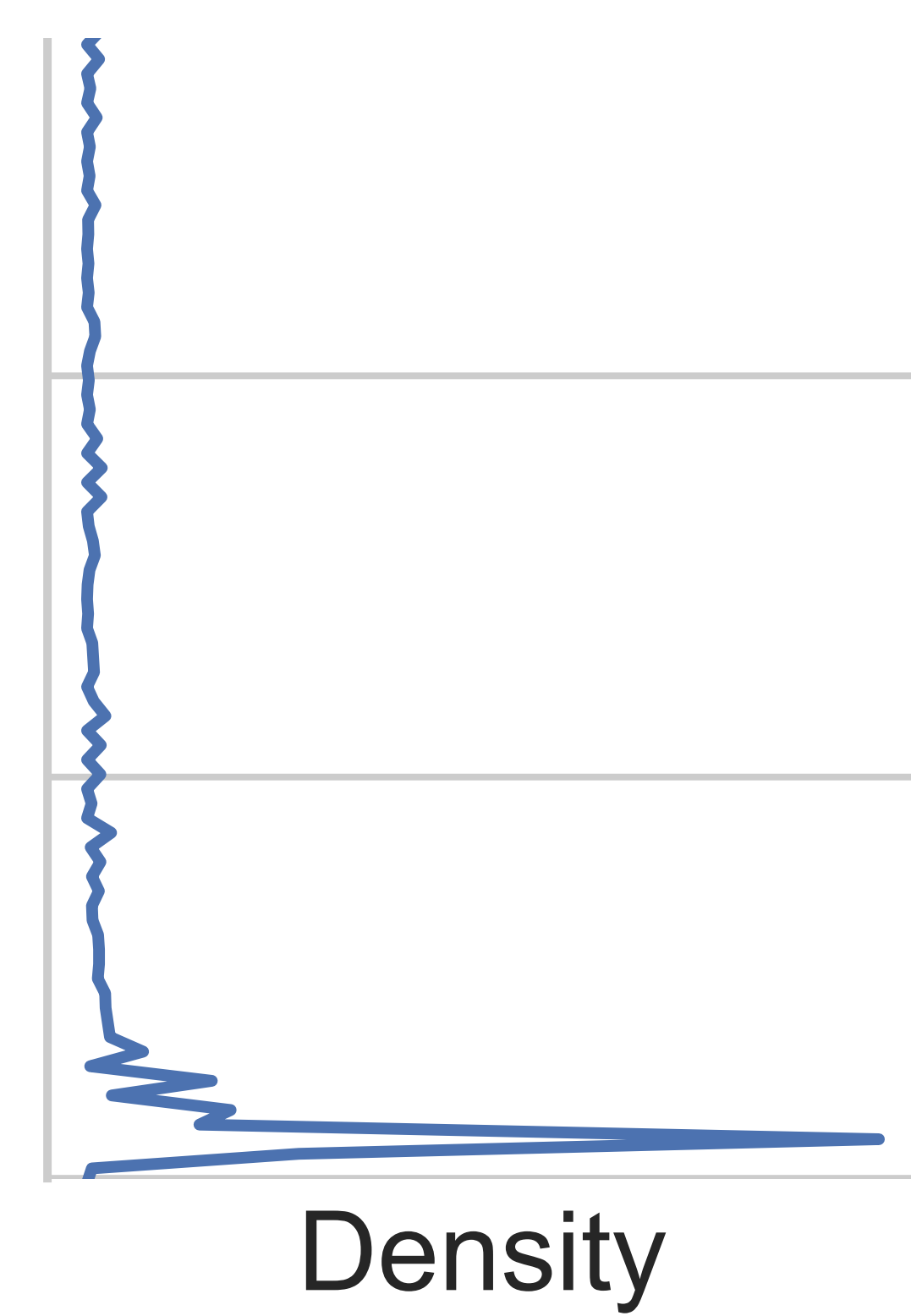**Parallelism is highly irregular, average parallelism is low**
*(scale-free input)*

$P_{max}$

Parallelism (x1000)

Iteration

Density

*Self-tuning controller:*
**Parallelism is more regular and average parallelism is higher**
*(scale-free input)*

$P_{max}$

*Self-tuning controller:*
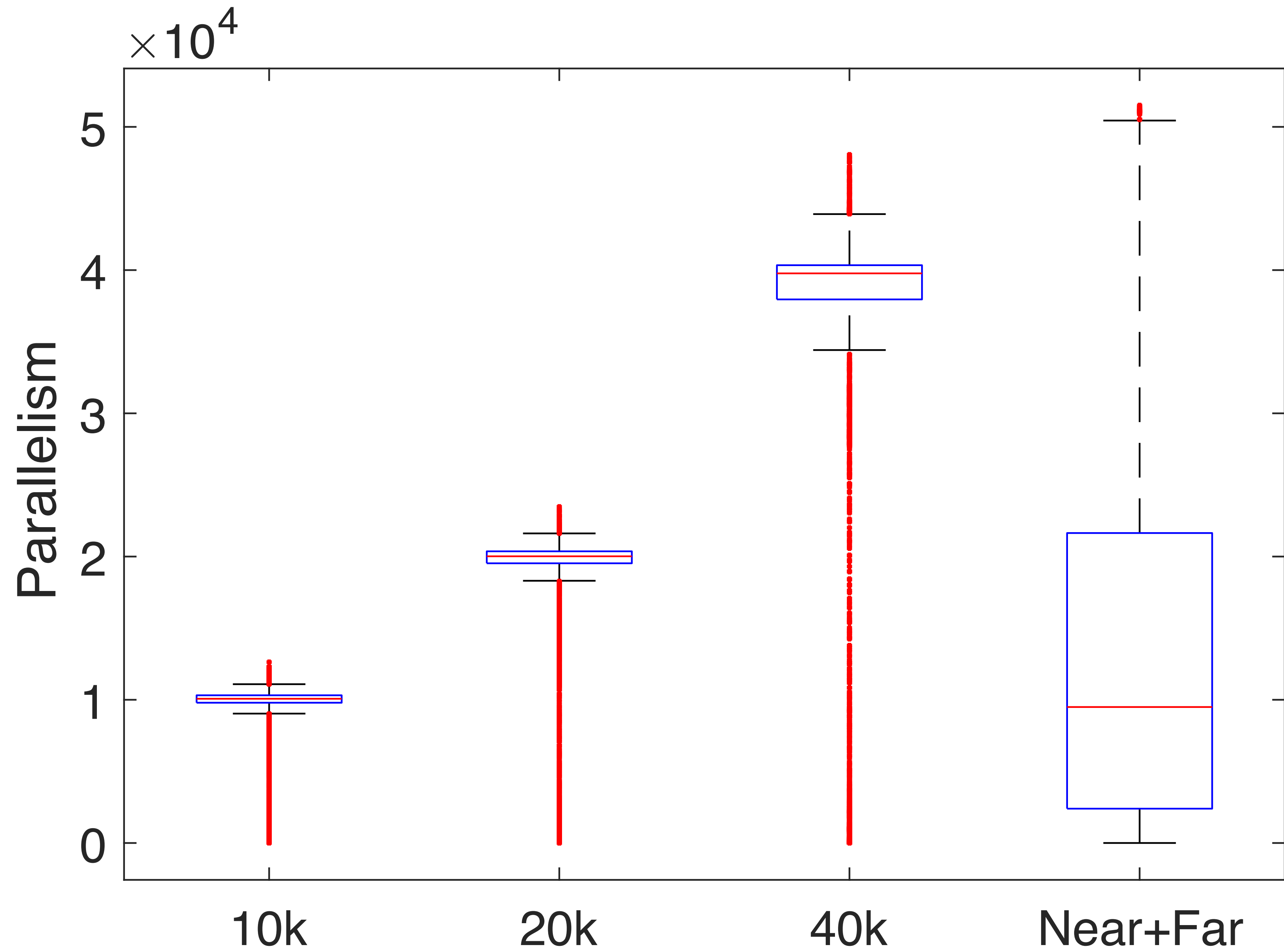**Parallelism is more regular and average parallelism is higher**
*(scale-free input)*

Baseline Near+Far:
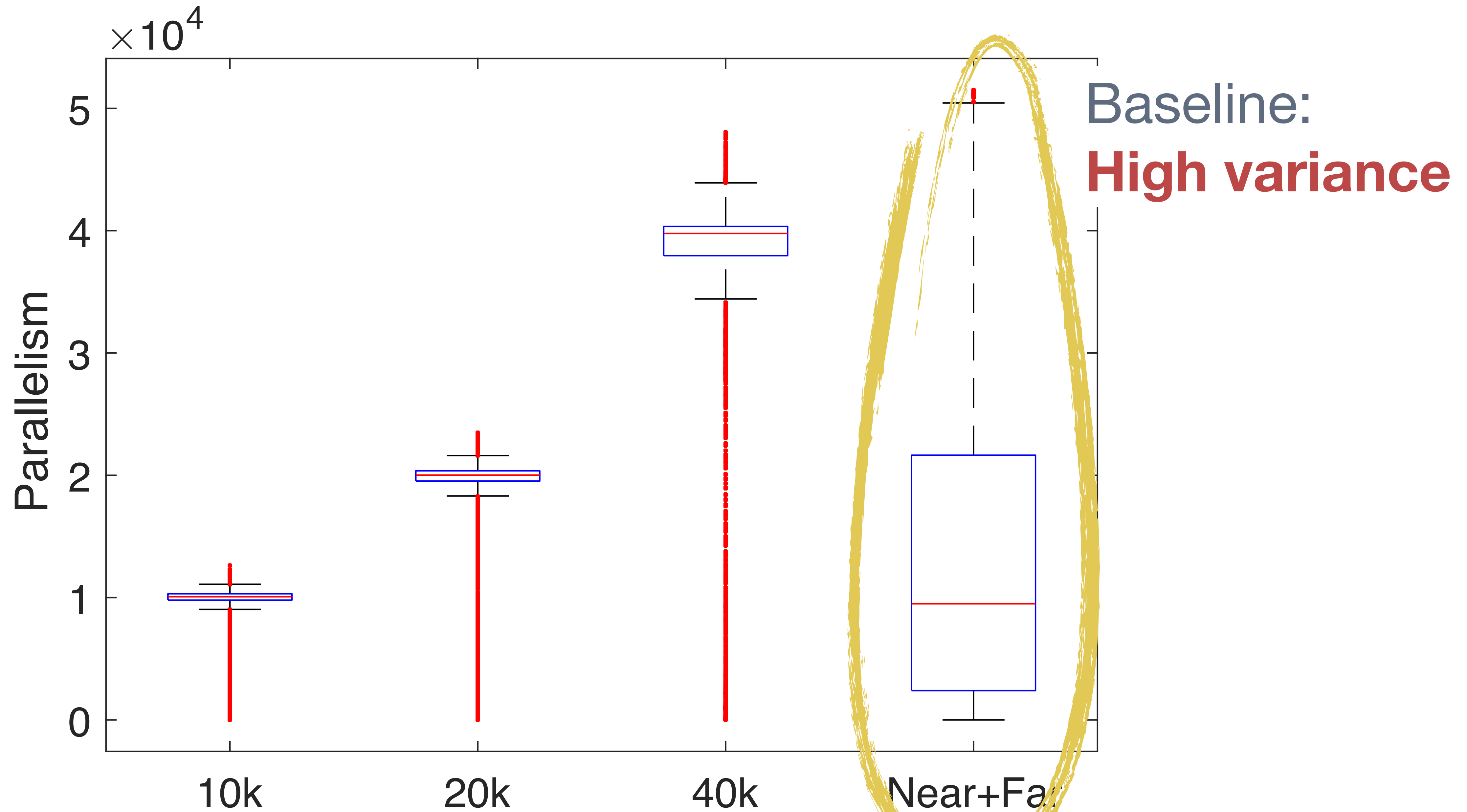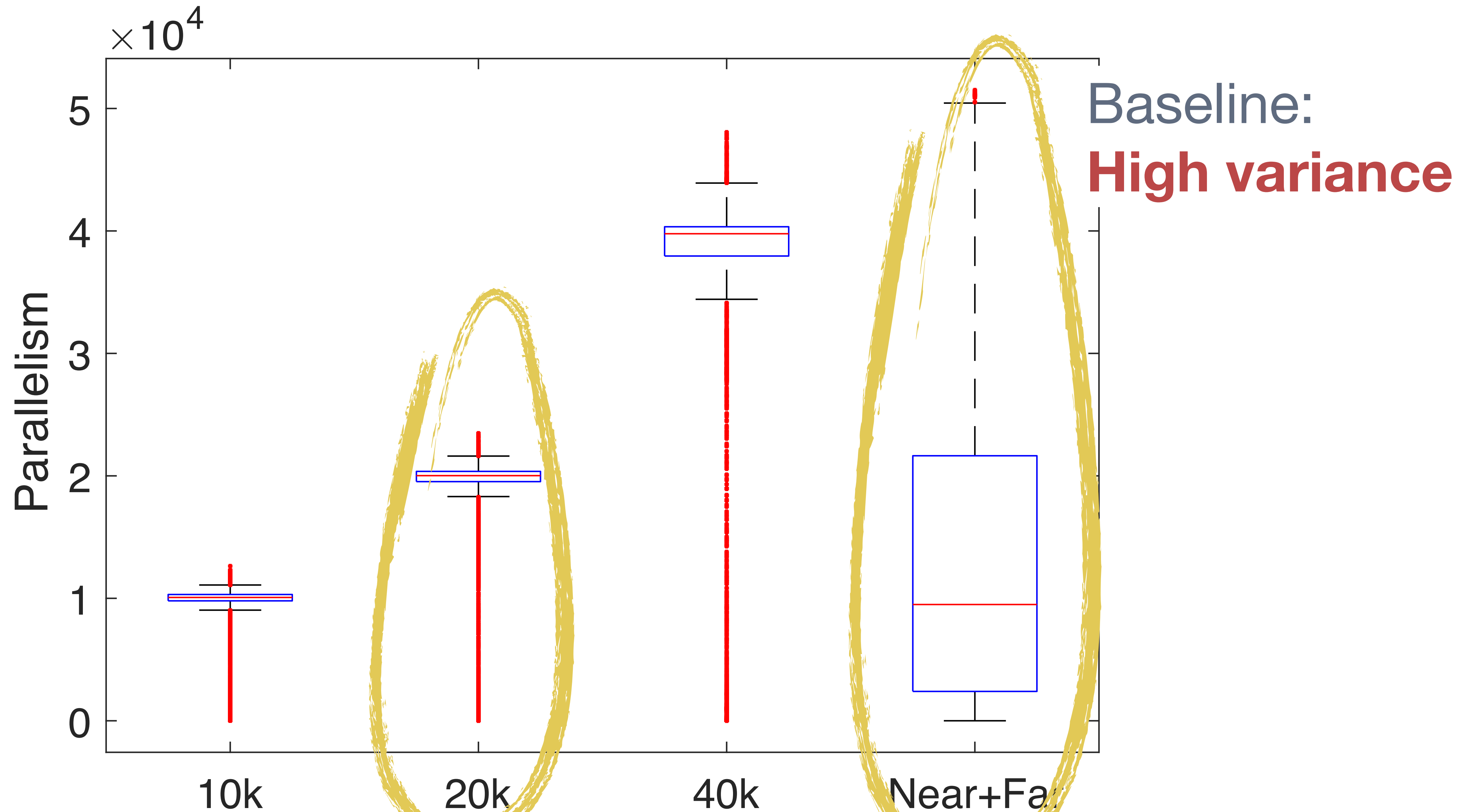**Low average parallelism ⇒ many iterations**

*(scale-free input)*

# *(road network)*

*(road network)*

Baseline:
**High variance**

*(road network)*

Baseline:
**High variance**

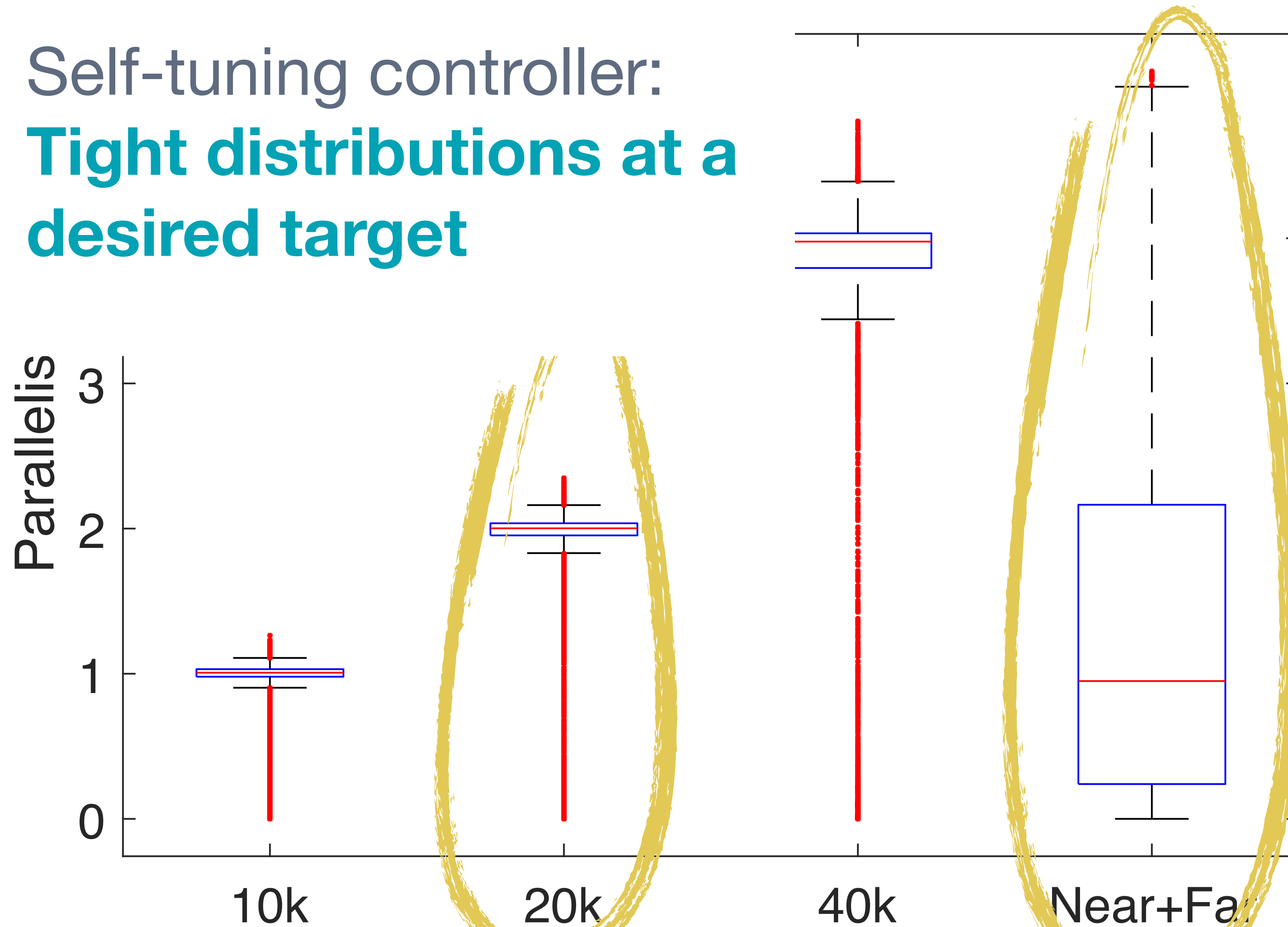*(road network)*

Self-tuning controller:
**Tight distributions at a desired target**

Baseline:
**High variance**



Parallelis

3

2

1

0

10k    20k    40k    Near+Fa
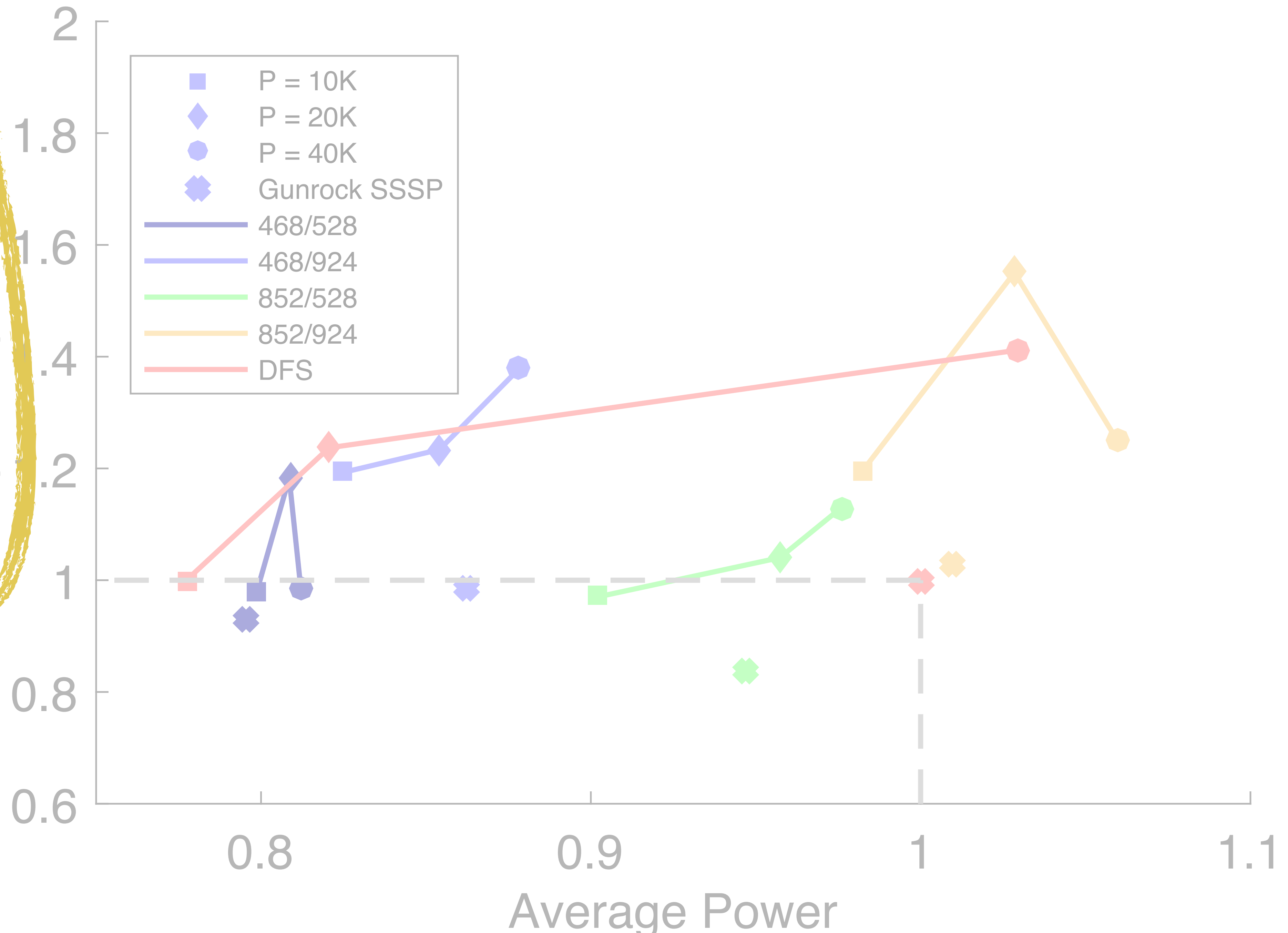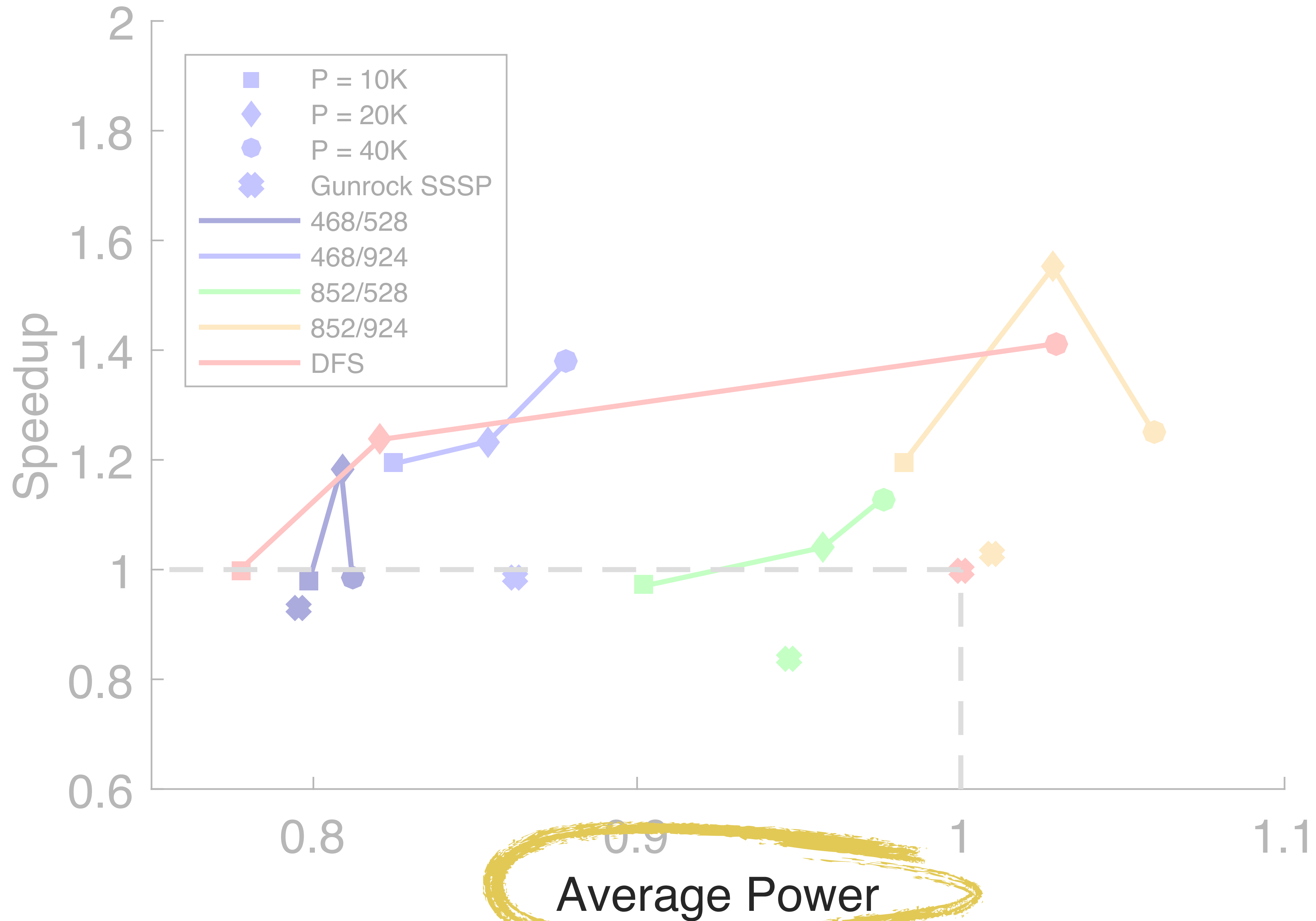
52

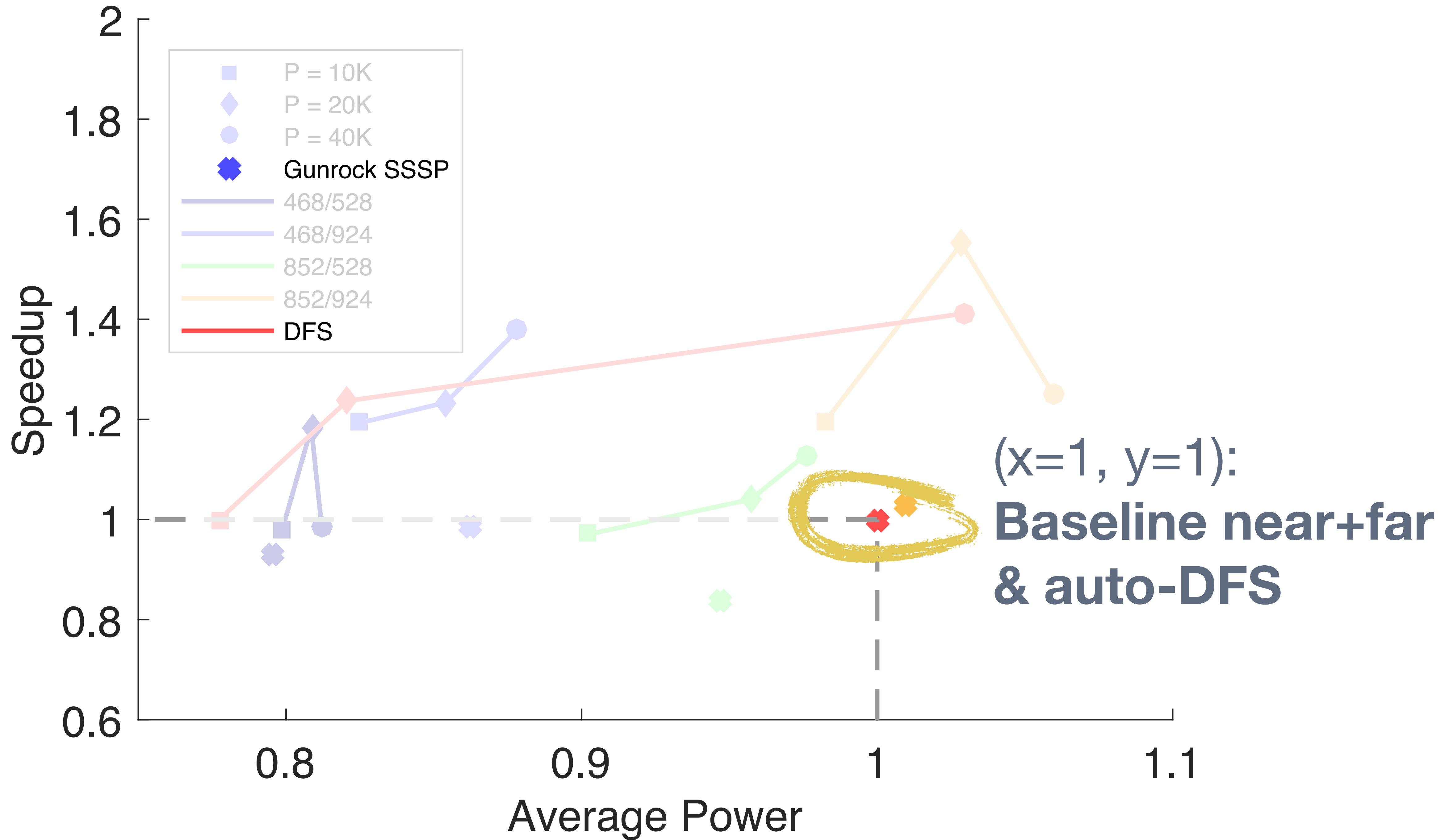*Next question:*

# Can it **save power** or **improve performance**?

*Compare against the baseline with a fixed δ and* **hardware-based dynamic frequency scaling.**

**Legend:**
- P = 10K
- P = 20K
- P = 40K
- Gunrock SSSP
- 468/528
- 468/924
- 852/528
- 852/924
- DFS

(x=1, y=1):
**Baseline near+far & auto-DFS**

Speedup

Average Power

57

*Upper-left quadrant:*

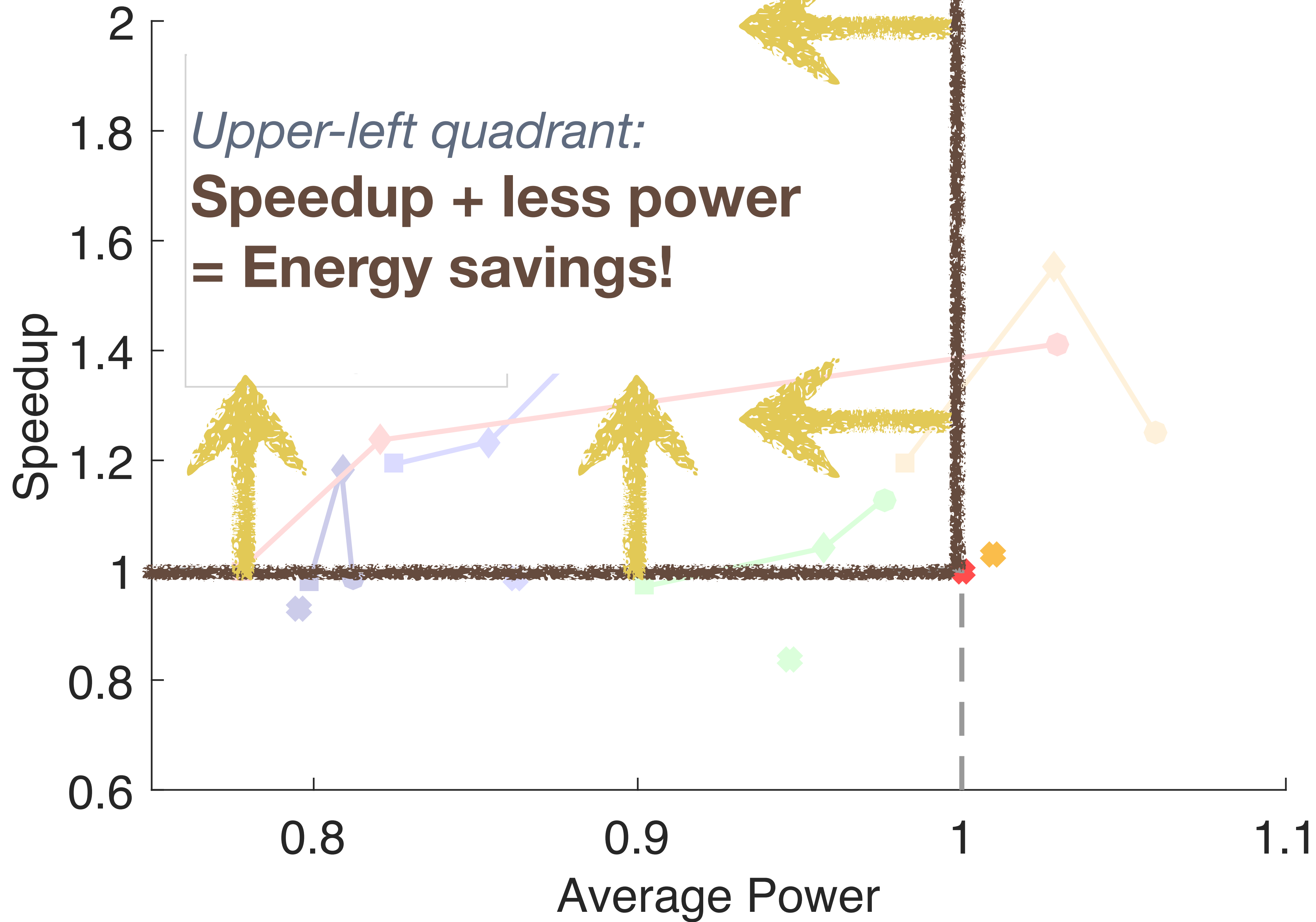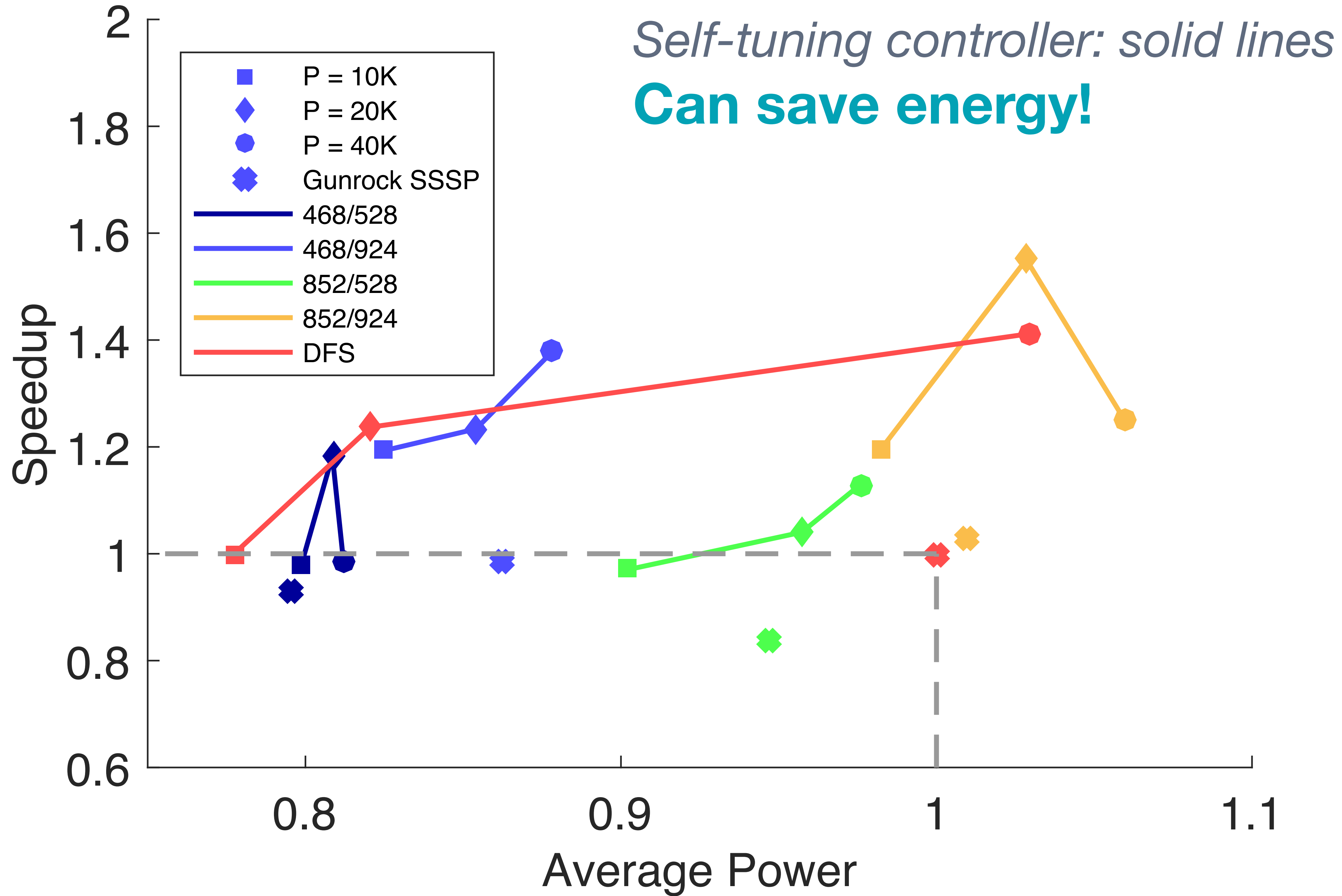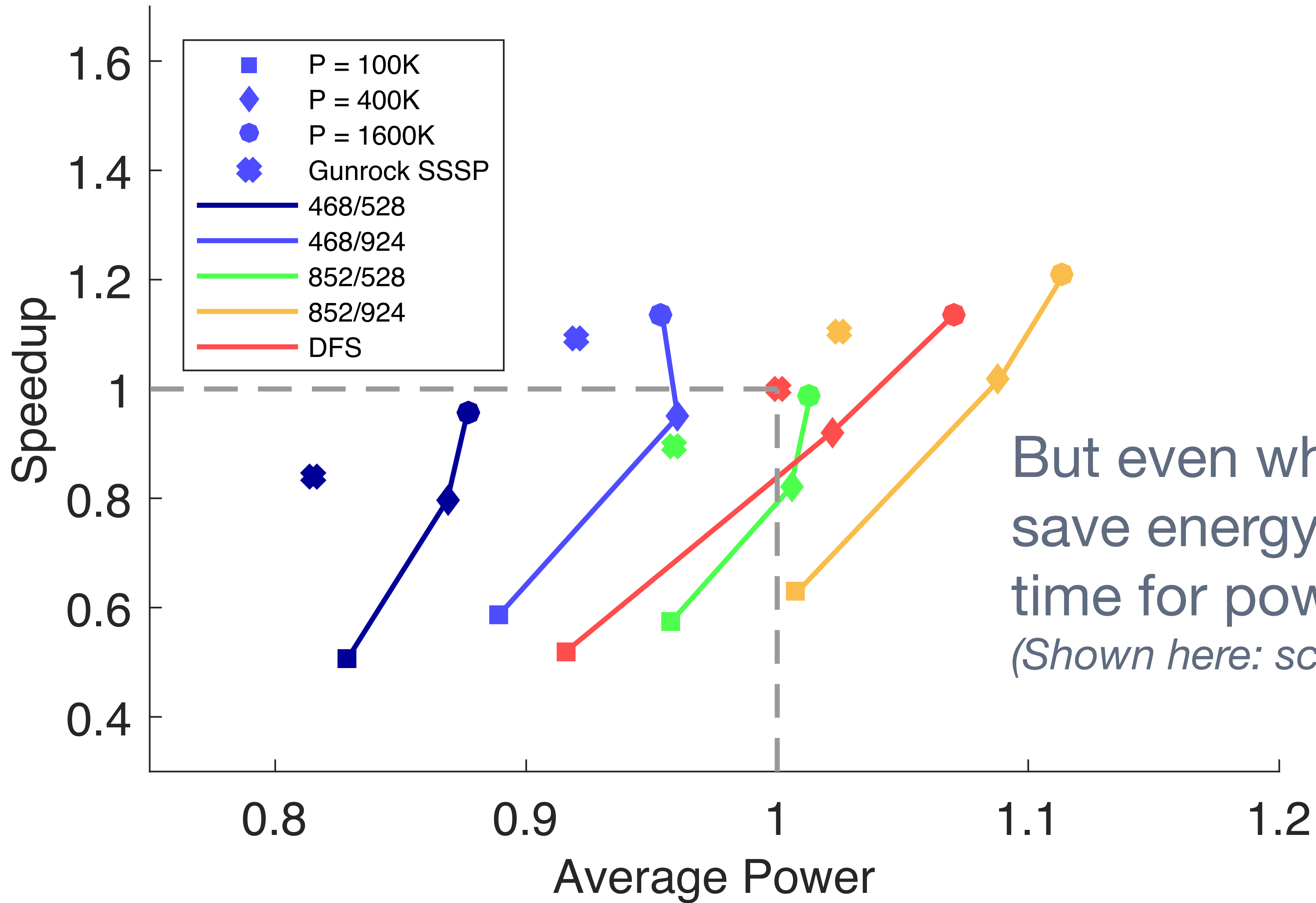**Speedup + less power = Energy savings!**

Speedup

Average Power

But even when it can't save energy, it trades time for power gracefully.
*(Shown here: scale-free network)*

*Limitation 1 (open-question):*

Choosing *P* is **not the same** as asking for max power, which was our motivation.

*There are limits on dynamic power measurement, which is needed to provide feedback to this scheme. But it would likely be easy to incorporate because of the model-based approach.*

*Limitation 2 (observational):*

# Power and energy **savings are not big**.

*We observed ~ 40% speedups and ~ 15% reductions in maximum power consumption over hardware-only DFS. These savings cannot be bigger because the system baseline or "constant" power is high — it is upwards of 50% or more of maximum power, so the amount of dynamically controllable power is small.*

*Conclusions:*

The "dynamic SSSP" algorithm improves on near+far, even when ignoring power. It's easier to choose P than $\delta$, making this SSSP easier to use.

The control-based scheme can be applied to any algorithm that is a "sequence of (filter) banks" (e.g., others in Gunrock), though the models may need specialization.

Energy savings are possible, especially when combined with hardware-only techniques like DFS.

**hpcgarage**

# Rich Vuduc

High-performance computing,
scalable parallel algorithms & software

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