Product (Re)forms Part II Improving Performance

Kristy Gardner Computer Science Department Amherst College Rhonda Righter IEOR Department UC Berkeley

INFORMS APS, July 4, 2019



BERKELEY THE DESCRIPTION OF THE PROPERTY OF TH

Tyler Maxey



Tyler Maxey

Outline – Improving Performance

- 1. Improving performance by collaborating (or not)
- 2. Improving performance by scheduling
- 3. Improving performance by increasing flexibility
- 4. Improving performance for whom?
- 5. Improving performance by how much?

Joint with Adan, Harchol-Balter, Hyytiä, Kleiner, Weiss More details in talks by Gardner, Weiss

| Collaboration Scheduling | Flexibility | Fairness | Convexity |
|--------------------------|-------------|----------|-----------|
|--------------------------|-------------|----------|-----------|

Does collaboration help?

Note: Even though we have product forms for steady-state distributions for both the collaborative and noncollaborative models, a direct comparison, even for steady-state means, is not possible except in simple cases.

We'll explore the question both for steady-state means and along sample paths.

| Collaboration Scheduling Flexibility Fairness Convexi | y |
|---|---|
|---|---|

Improving Performance by Collaborating?

Does collaboration help? Sometimes!

Collaboration > noncollaboration (C > NC) to minimize the number in system for:

- M/M/K: One class, heterogeneous servers
 Collaboration ⇔ M/M/1 with fast server
- 2. Fully symmetric system:

Servers and job classes stochastically identical

$$\lambda_i = \lambda_j, \mu_k = \mu_l, \forall i, j, k, l$$

 S_i symmetric: *d* randomly selected (power of d)

3. Symmetric W model: $\lambda_1 = \lambda_2$ $\mu_1 = \mu_2$

Intuition: Collaboration keeps servers busy

| Collaboration Scheduling Flexibility Fairness Convexity | 3 |
|---|---|
|---|---|

Improving Performance by Collaborating?

Does collaboration help? Sometimes!

Collaboration > noncollaboration (C > NC) to minimize the number in system for:

- M/M/K: One class, heterogeneous servers
 Collaboration ⇔ M/M/1 with fast server
- 2. Fully symmetric system:

Servers and job classes stochastically identical

$$\lambda_i = \lambda_j, \mu_k = \mu_l, \forall i, j, k, l$$

S_i symmetric: d randomly selected (power of d)

3. Symmetric W model: $\lambda_1 = \lambda_2$ $\mu_1 = \mu_2$

But sometimes collaboration can be much worse ...

| Collaboration S | Scheduling | Flexibility | Fairness | Convexity | 3 |
|-----------------|------------|-------------|----------|-----------|---|
|-----------------|------------|-------------|----------|-----------|---|

Thm: For the N model,

 $N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$ all can be arbitrarily large. (*F* for flexible)

Proof: By an example sample path



Thm: For the N model,

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

Proof: By an example sample path



Thm: For the N model,

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

Proof: By an example sample path



State:

(2,2,F,2,2,2)

 $(2, F_1, 2, 2, 2)$

 λ_2

Sample path: Service on 2,2,1, Arrival F, 2,2,2

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 4 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Thm: For the N model,

State:

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

Proof: By an example sample path



Sample path: Service on 2,2,1, Arrival *F*, 2,2,2, Service on 2,2,2,1

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 4 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Thm: For the N model,

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

State:

Proof: By an example sample path



Sample path: Continuing – arbitrary build up of class 2 in the C system

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | |
|---------------|------------|-------------|----------|-----------|--|
|---------------|------------|-------------|----------|-----------|--|

Thm: For the N model,

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

Proof: By an example sample path



State:

(2,2,2,2,2,F,F,F,F,F)

Sample path: Arrival F, F, F, F

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 4 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Thm: For the N model,

$$N^{C}(t) - N^{NC}(t), N_{2}^{C}(t) - N_{2}^{NC}(t), \text{ and } N_{F}^{C}(t) - N_{F}^{NC}(t)$$

all can be arbitrarily large.

Proof: By an example sample path



Sample path: Arrival F, F, F, F, Service on 2,2,2,2

State:

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 4 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Collaboration is better in symmetric systems

- Collaboration can be worse than noncollaboration:
- We've seen that

$$N^{C}(t) - N^{NC}(t)$$

can be arbitrarily large on a sample path

More flexible jobs can block less flexible jobs in the C (collaborative) system

Are these sample paths sufficiently rare so that the mean is bounded?

| Collaboration Scheduling Flexibility Fairness Convexity | Flexibility Fairness Convex | |
|---|-----------------------------|--|
|---|-----------------------------|--|

Improving Performance by Collaborating?

No! **Collaboration can be much worse**, even on average, for asymmetric systems.



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 6 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

We've seen that the **cost of collaboration**, in terms of # of jobs, can be arbitrarily large:

 $N^{C}(t) - N^{NC}(t)$ can be arbitrarily large

 $E[N^{C}] - E[N^{NC}]$ can be arbitrarily large

Can the **benefit of collaboration** also be arbitrarily large?

No! The benefit of collaboration is bounded by M = # of servers

 $N^{NC}(t) - N^{C}(t) \le M$ on coupled sample paths

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 7 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Collaboration Benefit Upper Bound

<u>Thm</u>: We can couple the C and NC systems so that $\{N^{NC}(t)\} \leq \{N^{C}(t)\} + M$ with probability 1 and $\{N_{i}^{NC}(t)\} \leq \{N_{i}^{C}(t)\} + |S_{i}|$ w.p. 1, where $|S_{i}| = \#$ servers compatible with class *i*

Proof Idea: Recall from Part I that $\{N_Q^{NC}(t) \mid \text{keep all servers busy}\} =_{st} \{N^C(t)\}$

| Collaboration Scheduling Flexibility | y Fairness Convexity | 8 |
|--------------------------------------|----------------------|---|
|--------------------------------------|----------------------|---|

Collaborative vs. Noncollaborative



has the same stationary distribution as the collaborative system

| Sustam | Detaile | d states | Polatod syste | Partial ag | gregation | Per-class aggregation | 0 |
|--------|---------|-----------|---------------|------------|-----------|-----------------------|---|
| System | Collab | Noncollab | Related systs | Collab | Noncollab | Collab | 9 |

Collaborative vs. Noncollaborative



Collaboration Benefit Upper Bound

Thm: We can couple the C and NC systems so that

$$\{N^{NC}(t)\} \leq \{N^{C}(t)\} + M$$
 with probability 1 and
 $\{N_{i}^{NC}(t)\} \leq \{N_{i}^{C}(t)\} + |S_{i}|$ w.p. 1,
where $|S_{i}| = \#$ servers compatible with class *i*

Proof Idea: Recall from Part I that

$$\{N_Q^{NC}(t) \mid \text{keep all servers busy} \} =_{st} \{N^C(t)\}$$

$$\{N^{NC}(t)\} \leq_{st} \{N_Q^{NC}(t) \mid \text{all servers busy} \} + M \quad (\text{easy})$$

$$=_{st} \{N^C(t)\} + M \quad (\text{from } *)$$

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 10 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Tighter Bound for the W Model

<u>Thm</u>: For the W model, we can couple the NC and C systems so that, w.p.1, $\{N^{NC}(t)\} \le \{N^{C}(t)\} + 1$ $\{N^{NC}_{i}(t)\} \le \{N^{C}_{i}(t)\} + 1, i = 1,2$ $\{N^{NC}_{i}(t) + N^{NC}_{3}\} \le \{N^{C}_{i}(t) + N^{C}_{3}\} + 1, i = 1,2$



→ NC has at most one more job than C (call it the tagged job)
 And if it does, then among all non-tagged jobs:
 NC has fewer jobs of each class than C (nonstrictly)
 The theorem also holds for W models embedded in larger nested systems.

W Model Bound – Proof Idea

- In C and NC: Class *i* jobs FIFO, *i* = 1,2 (they leave in arrival order)
- In C (only!): Class 3 jobs FIFO
- In C: Class i jobs + class 3 jobs FIFO (Class i jobs can't "pass" class 3 jobs)
- If NC has an extra tagged job T, and otherwise fewer of each class than C
- T must be in service in NC,
- The "next" job (A) on the other server must be in service in both C and NC







| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 12 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Outline – Improving Performance

- 1. Improving performance by collaborating (or not)
- 2. Improving performance by scheduling
- 3. Improving performance by increasing flexibility
- 4. Improving performance for whom?
- 5. Improving performance by how much?

Scheduling in Nested Systems



For every pair of job classes *i* and *j*, either:



•
$$S_j \subset S_i$$

•
$$S_i \cap S_j = \emptyset$$



 $R(i) = \{j: S_j \subseteq S_j\}$ = the set of classes that **R**equire at least one of the servers in S_i = class *i* plus all of the classes within its subsystem

 $K_i(t) = \sum_{j \in R(i)} N_j(t)$ = the total number of jobs in class *i*'s subsystem

E.g.,
$$K_3(t) = N_3 + N_4 + N_5$$

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 14 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Scheduling in Nested Systems

<u>Thm</u>: For nested (collaborative or noncollaborative) systems, preemptive non-idling LFF (least-flexible-first) minimizes $\{\vec{K}(t)\}$ w.p. 1 for coupled sample paths, where $\vec{K}(t) = (K_1(t), K_2(t), \dots, K_J(t))$ and $K_i(t) =$ the total number of jobs in class *i*'s subsystem



LFF: 4 > 3 > 6 on server 4, etc.

Intuition: Save more flexible jobs to run when servers might otherwise be idle

Product form no longer holds, though least flexible jobs experience an M/M/1 queue

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 15 |
|---------------|------------|-------------|----------|-----------|----|
| | | | | | |

LFF Optimality Proof for the W model

<u>Thm</u>: For the W model, LFF minimizes, w.p. 1, $\{(N_1(t), N_2(t), N_1(t) + N_2(t) + N_3(t))\}$



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 16 |
|---------------|------------|-------------|----------|-----------|----|
| | | | | | |

LFF Optimality Proof for the W model

<u>Thm</u>: For the W model, LFF minimizes, w.p. 1, $\{(N_1(t), N_2(t), N_1(t) + N_2(t) + N_3(t))\}$

Proof:

At time 0, π serves a class 3 job (A) on server 1 Let π' serve the class 1 job (B) on server 1

Case 1. The next event is an arrival or server 2 completion: Let $\pi' = \pi$ thereafter, so $N'_i(t) = N_i(t)$ for all *i*,*t*





Collaboration

LFF Optimality Proof for the W model

<u>Thm</u>: For the W model, LFF minimizes, w.p. 1, $\{(N_1(t), N_2(t), N_1(t) + N_2(t) + N_3(t))\}$

<u>Proof</u>:

Case 2. Next event a server 1 completion: Let π' serve A whenever π serves B; π otherwise let $\pi' = \pi$ (OK because of **nested structure**, may **idle**) All job completions at same times, but A, B interchanged: A leaves earlier under π and **B** leaves earlier under π' π $N_1'(t) \le N_1(t), N_2'(t) \le N_2(t),$ $N_1'(t) + N_2'(t) + N_3'(t) \le N_1(t) + N_2'(t) + N_3(t)$





8

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 1 |
|---------------|------------|-------------|----------|-----------|---|
|---------------|------------|-------------|----------|-----------|---|

Non-nested Systems: LFF is not optimal



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 19 |
|---------------|------------|-------------|----------|-----------|----|
| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 19 |

Is Collaboration a Good Thing?

For FCFS – Sometimes!

For LFF in nested systems – Always!

For example, in the W model, class 3 jobs are only served on server *i* when no class *i* jobs are present, *i* = 1,2 – when using both servers can only help.



| \frown | | I | | · | |
|----------|-----|----|------|---|--|
| | 112 | no | rati | n | |
| | l a | | ιαι | | |
| | | | | | |

Outline – Improving Performance

- 1. Improving performance by collaborating (or not)
- 2. Improving performance by scheduling
- 3. Improving performance by increasing flexibility
- 4. Improving performance for whom?
- 5. Improving performance by how much?

Improving Performance with Flexibility?

Under FCFS, more flexibility helps in symmetric systems, when flexibility added symmetrically

1. Fully symmetric system Servers and job classes stochastically identical $\lambda_i = \lambda_j, \mu_k = \mu_l, \forall i, j, k, l$ S_i symmetric: *d* randomly selected (power of d) 2. Symmetric W model



Response time decreasing in *d*

Response time decreasing in Δ

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 22 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

W Model: Adding Flexibility Symmetrically

Under FCFS (and LFF) more flexibility helps in symmetric systems



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 23 |
|---------------|------------|-------------|----------|-----------|----|
| | | | | 4 | |

Improving Performance with Flexibility?

Under FCFS, more flexibility can hurt (asymmetric) Under LFF and nested, more flexibility always helps



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 24 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Improving Performance with Flexibility?

LFF: A little bit (of flexibility) goes a long way

FCFS: A little bit more goes the wrong way!



| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 24 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Outline – Improving Performance

- 1. Improving performance by collaborating (or not)
- 2. Improving performance by scheduling
- 3. Improving performance by increasing flexibility
- 4. Improving performance for whom?
- 5. Improving performance by how much?

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 25 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Improving Performance for Whom?

Collaboration, scheduling, flexibility are all means for (possibly) improving **overall** performance

But even when overall performance improves, some classes of jobs may be hurt

FCFS → LFF hurts most flexible jobs
 FCFS → LFF helps least flexible jobs

W model (collaborative) example:

$$\mu_1 = \mu_2 = 1, \lambda = 1.6, \lambda_2 = 0.8$$

 $\lambda_1 = 0.8 - \Delta, \lambda_3 = \Delta$



 Δ = shift from class 1 to class 3

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 26 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

W Model Response Times – Shift 1 to 3



Improving Performance for Whom?

- 1. Going from FCFS → LFF hurts some job classes
- 2. Under FCFS
 - even when more flexibility helps overall, it may hurt some job classes
 - *collaboration* may hurt some job classes
- 3. Under LFF, in nested systems
 - *more flexibility* helps all classes (win-win!)
 - *collaboration* helps all classes

How can we improve performance, relative to FCFS, without hurting any job class?

4. Going from FCFS to the PF (primaries first) policy ensures more flexibility helps all classes

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 28 |
|---------------|------------|-------------|----------|-----------|----|
| | | | | | 4 |

Primaries First Policy

Objective: Improve from status quo Starting with an initial FCFS system, call it System *P* (primary)

Create a new system through scheduling and increased flexibility that provides improved performance for all job classes.

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 29 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Primaries First Policy

System *P* (primary): Initial FCFS system with $S_i \subset S_j$. System *P*+*S* (primary plus secondary): Shift arrival rates from λ_i and λ_j to $\lambda_i - \Delta$ and $\lambda_j + \Delta$ (still stable!)



- Original copies are primary
- New copies are secondary
- Primaries have preemptive priority over secondaries
- FCFS within primaries
- FCFS within secondaries

Under PF, all job classes benefit from increasing flexibility

| Collaboration Scheduling Flexibility Fairness Convexity 30 | Collaboration | Scheduling | Flexibility | Fairness | Convexity | 30 |
|--|---------------|------------|-------------|----------|-----------|----|
|--|---------------|------------|-------------|----------|-----------|----|

Outline – Improving Performance

- 1. Improving performance by collaborating (or not)
- 2. Improving performance by scheduling
- 3. Improving performance by increasing flexibility
- 4. Improving performance for whom?
- 5. Improving performance (through increasing flexibility, collaborative)
 by how much?

| Collaboration Scheduling Flexibility Fairness Convexity | Collaboration | Scheduling | Flexibility | Fairness | Convexity | 31 |
|---|---------------|------------|-------------|----------|-----------|----|
|---|---------------|------------|-------------|----------|-----------|----|

FCFS and Convexity

Under **FCFS**, more flexibility **can hurt** overall response time, response time may be neither convex nor concave as a function of flexibility



FCFS and Per Class Convexity

<u>**Thm</u>**: As Δ , the arrival rate shifted from class 1 to class 3, increases</u>

- E[W] may be nonmonotonic, nonconvex, nonconcave
- E[W₁] is decreasing and convex
- E[W₂] is increasing and concave
- $E[W_3]$ is constant



FCFS and Per Class Convexity

<u>**Thm</u>**: As Δ , the arrival rate shifted from class 1 to class 3, increases</u>

- E[W] may be nonmonotonic, nonconvex, nonconcave
- E[W₁] is decreasing and convex
- E[W₂] is increasing and concave
- $E[W_3]$ is constant

The proof for the class-based response times follows from Part I.

<u>Corollary</u>: $W_{i} =_{st} W\left(\lambda(R(S_{i})), \mu(S_{i})\right) + \sum_{j:S_{i} \subseteq S_{j}} W^{Q}\left(\lambda_{j}, \mu(S_{j}) - \lambda(R(S_{j})) + \lambda_{j}\right)$

The theorem can be generalized to arbitrary nested collaborative systems.

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 34 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

LRF and Convexity (Collaborative, Nested)

<u>Thm</u>: Overall mean response time is decreasing and convex in λ_3 , holding $\lambda_1 + \lambda_1 + \lambda_3$ fixed. The benefit of shifting from λ_1 to λ_3 is increasing in the amount shifted from λ_2 to λ_3 .

Mean Overall Response Time **Collaborative W** $\lambda_3 \lambda_2$ λ₁ 15 10 5 0 $\mu_1 = \mu_2 = 1$ $\lambda_1 + \lambda_1 + \lambda_3 = 1.8$ 0.5 0.6 0.4 λ_1 0.2 0 λ_2

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 35 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|

Conclusions – Part I

- Product forms for both collaborative and noncollaborative models can be shown easily by considering detailed states and order independence
- 2. Aggregating states yields partial results for performance measures
- 3. Fully flexible class experiences an M/M/1 system in the collaborative model
- 4. Complete, simple results hold for nested collaborative systems
- 5. Results for the collaborative model hold for the noncollaborative model given all servers busy

| System | Detailed states | | Polatod cysts | Partial aggregation | | Per-class aggregation | |
|--------|-----------------|-----------|---------------|---------------------|-----------|-----------------------|----|
| System | Collab | Noncollab | Related systs | Collab | Noncollab | Collab | 50 |

Conclusions – Part II

- 1. Collaboration and flexibility may not help
 - They help in symmetric FCFS systems
 - They help in nested LFF systems
- 2. Improving overall performance may not be fair
 - Under Primaries First (PF) more flexibility helps all classes
 - Under PF collaboration helps all classes
- 3. Under LFF, in nested systems, flexibility always helps
 - Diminishing marginal returns
 - Increasing cross-derivative effects

| Collaboration | Scheduling | Flexibility | Fairness | Convexity | 37 |
|---------------|------------|-------------|----------|-----------|----|
|---------------|------------|-------------|----------|-----------|----|