Routing and Staffing in Large-Scale Service Systems with Heterogeneous Servers and Impatient Customers

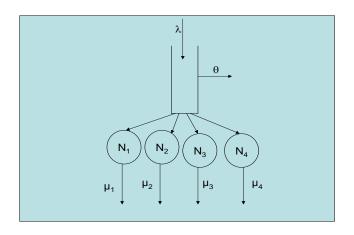
Mor Armony New York University

Joint Work with Avi Mandelbaum

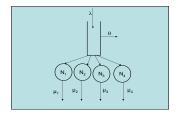
Motivation: Call Centers



The Inverted-V Model with Abandonment



The Inverted-V Model with Abandonment: Motivation



- ▶ Heterogeneous server population
- Learning Effects
- Various server skills

The Model

- ▶ Single customer class Poisson arrival process with rate λ .
- \blacktriangleright K server pools $(N_1, N_2, ..., N_K \text{ servers})$
- Exponential non-preemptive service times with rates

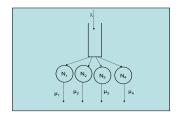
$$\mu_1 < \mu_2 < \dots < \mu_K$$
.

Exponential time to abandonment with rate θ .

Our Focus: Staffing and Routing

- ▶ How many servers of each type are needed?
- How to route incoming or waiting customers?

Inverted-V model Without Abandonment (Armony '05)



Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $P(W > 0) \le \alpha$, for some routing policy

4 D > 4 D > 4 E > 4 E > E 9 Q Q

Even <u>little abandonment</u> can have a significant effect on performance:

▶ An unstable M/M/N system $(\rho > 1)$ becomes stable with abandonment (M/M/N + M).

- ▶ An unstable M/M/N system $(\rho > 1)$ becomes stable with abandonment (M/M/N + M).
- ▶ Example: (Mandelbaum & Zeltyn '06) Consider $\lambda=2000$, $\mu=20$. Service level target: "80% of customers should be served within 30 second".

- ▶ An unstable M/M/N system $(\rho > 1)$ becomes stable with abandonment (M/M/N + M).
- ▶ Example: (Mandelbaum & Zeltyn '06) Consider $\lambda=2000$, $\mu=20$. Service level target: "80% of customers should be served within 30 second".
 - ▶ 106 agents $(\theta = 0)$.

- ▶ An unstable M/M/N system $(\rho > 1)$ becomes stable with abandonment (M/M/N + M).
- ▶ Example: (Mandelbaum & Zeltyn '06) Consider $\lambda=2000$, $\mu=20$. Service level target: "80% of customers should be served within 30 second".
 - ▶ 106 agents $(\theta = 0)$.
 - ▶ 95 agents ($\theta = 20$ (avg. patience of 3 minutes), P(ab) = 6.9%)

- ▶ An unstable M/M/N system $(\rho > 1)$ becomes stable with abandonment (M/M/N + M).
- ▶ Example: (Mandelbaum & Zeltyn '06) Consider $\lambda=2000,~\mu=20.$ Service level target: "80% of customers should be served within 30 second".
 - ▶ 106 agents ($\theta = 0$).
 - ▶ 95 agents ($\theta = 20$ (avg. patience of 3 minutes), P(ab) = 6.9%)
 - ▶ 84 agents ($\theta = 60$ (avg. patience of 1 minute), P(ab) = 16.8%)

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $P(W > T) \le \alpha_T$, for some routing policy $EW \le \bar{W}$, $P(ab) \le \Delta$.

Issues related to formulation:

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $P(W > T) \le \alpha_T$, for some routing policy $EW \le \bar{W}$, $P(ab) \le \Delta$.

Issues related to formulation:

FCFS: Natural but not necessarily optimal.

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $P(W > T) \le \alpha_T$, for some routing policy $EW \le \bar{W}$, $P(ab) \le \Delta$.

Issues related to formulation:

- ▶ FCFS: Natural but not necessarily optimal.
- ▶ Intensional Idling can improve service level.

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $P(W > T) \le \alpha_T$, for some routing policy $EW \le \bar{W}$, $P(ab) \le \Delta$.

Issues related to formulation:

- ▶ FCFS: Natural but not necessarily optimal.
- ▶ Intensional Idling can improve service level.
- ▶ Not all are Customer-subjective measurements.

Background: Garnett, Mandelbaum & Reiman '02

In a sequence of M/M/N+M models, N=1,2,3,..., with $R=\lambda/\mu$, the following are equivalent:

$$N \approx R + \beta \sqrt{R}, -\infty < \beta < \infty;$$

Background: Garnett, Mandelbaum & Reiman '02

In a sequence of M/M/N+M models, N=1,2,3,..., with $R=\lambda/\mu$, the following are equivalent:

- $ightharpoonup N pprox R + \beta \sqrt{R}, -\infty < \beta < \infty;$

Background: Garnett, Mandelbaum & Reiman '02

In a sequence of M/M/N+M models, N=1,2,3,..., with $R=\lambda/\mu$, the following are equivalent:

- $N \approx R + \beta \sqrt{R}, -\infty < \beta < \infty;$
- $| \lim_{N \to \infty} P_N \{W > 0\} = \alpha, \quad 0 < \alpha < 1;$
- ▶ $\lim_{N\to\infty} \sqrt{N} P_N \{ab\} = \Delta$, $0 < \Delta < \infty$;

Here
$$\alpha = w(-\beta, \sqrt{\mu/\theta})$$
, $\Delta = \left[\sqrt{\theta/\mu} \cdot h(\beta\sqrt{\mu/\theta}) - \beta\right] \alpha$, $w(x, y) = \left[1 + \frac{h(-xy)}{yh(x)}\right]^{-1}$, $h(x) = \frac{\phi(x)}{1 - \Phi(x)}$.

- 1. QED regime: $N = R + \beta \sqrt{R}$
 - $\sqrt{\lambda}P\{ab\} \leq \Delta$, $\sqrt{\lambda}EW \leq \bar{w}$

- 1. QED regime: $N = R + \beta \sqrt{R}$
 - $\sqrt{\lambda}P\{ab\} \leq \Delta$, $\sqrt{\lambda}EW \leq \bar{w}$
- 2. ED regime: $N = (1 \gamma) \cdot R$, $\gamma > 0$
 - ▶ $P{ab} \le \Delta$, $EW \le \bar{w}$

- 1. QED regime: $N = R + \beta \sqrt{R}$
 - $\sqrt{\lambda}P\{ab\} \leq \Delta$, $\sqrt{\lambda}EW \leq \bar{w}$
- 2. ED regime: $N = (1 \gamma) \cdot R$, $\gamma > 0$
 - ▶ $P{ab} \le \Delta$, $EW \le \bar{w}$
- 3. ED + QED regime: $N = (1 \gamma)R + \delta\sqrt{R}$, $\gamma > 0$.
 - ▶ $P\{W > T\} \le \alpha$.



Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $\sqrt{\lambda}P(ab) \leq \Delta$, for some routing policy.

Challenges:

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $\sqrt{\lambda}P(ab) \leq \Delta$, for some routing policy.

Challenges:

Offered load not well-defined

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

Subject to $\sqrt{\lambda}P(ab) \leq \Delta$, for some routing policy.

Challenges:

- Offered load not well-defined
- Optimal routing unknown

Minimize
$$C_1(N_1) + C_2(N_2) + ... + C_K(N_K)$$

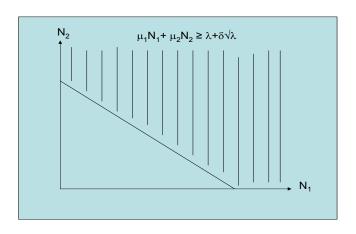
Subject to $\sqrt{\lambda}P(ab) \leq \Delta$, for some routing policy.

Challenges:

- Offered load not well-defined
- ▶ Optimal routing unknown

$$\mu_1 N_1 + \mu_2 N_2 + \dots + \mu_K N_K = \lambda + \delta \sqrt{\lambda}, \quad -\infty < \delta < \infty$$

Asymptotic Feasible Region



The Asymptotic Feasible Region

Theorem (Asymptotic Feasible Region): Consider a sequence of systems indexed by $\lambda\uparrow\infty$. Suppose that $\liminf_{\lambda\to\infty}N_1/N>0$. Then there exists a non-preemptive policy under which

$$\limsup_{\lambda\to\infty}\sqrt{\lambda}P_{\lambda}(ab)\leq\Delta,$$

if and only if

$$\mu_1 N_1 + \mu_2 N_2 + ... + \mu_K N_K \ge \lambda + \delta \sqrt{\lambda} + o(\sqrt{\lambda}), \quad -\infty < \delta < \infty.$$

Exact Optimal Routing: Unknown (De Vericourt and Zhou '06)

Exact Optimal Routing: Unknown (De Vericourt and Zhou '06)

Proposed Routing: Route to Faster Servers First (FSF).

Exact Optimal Routing: Unknown (De Vericourt and Zhou '06)

Proposed Routing: Route to Faster Servers First (FSF).

Potential Problem: Preemption may lead to excessive idling of fast servers.

Exact Optimal Routing: Unknown (De Vericourt and Zhou '06)

Proposed Routing: Route to Faster Servers First (FSF).

Potential Problem: Preemption may lead to excessive idling of fast servers.

Proposition (Optimal Preemptive Routing): FSF_P is optimal in the sense that it stochastically minimizes the cumulative number of abandoning customers.

Exact Optimal Routing: Unknown (De Vericourt and Zhou '06)

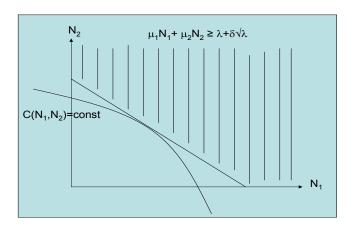
Proposed Routing: Route to Faster Servers First (FSF).

Potential Problem: Preemption may lead to excessive idling of fast servers.

Proposition (Optimal Preemptive Routing): FSF_P is optimal in the sense that it stochastically minimizes the cumulative number of abandoning customers.

Proposition (Asymptotically Optimal Routing): FSF is asymptotically optimal in the sense that in the limit FSF and FSF $_P$ have the same performance. (**Proof:** State-space collapse - Faster servers are always busy)

Asymptotically Optimal Staffing



Asymptotically Optimal Staffing: Example

Problem:

$$\mbox{Minimize} \hspace{0.5cm} \mathcal{C}_1 \textit{N}_1^{\textit{p}} + \mathcal{C}_2 \textit{N}_2^{\textit{p}} + ... + \mathcal{C}_{\textit{K}} \textit{N}_{\textit{K}}^{\textit{p}}, \hspace{0.2cm} \textit{p} > 1$$

Subject to
$$\sqrt{\lambda}P\{ab\} \leq \Delta$$
.

Solution:

Minimize
$$C_1 N_1^p + C_2 N_2^p + ... + C_K N_K^p, p > 1$$

Subject to
$$\mu_1 N_1 + \mu_2 N_2 + ... \mu_K N_K \ge \lambda + \delta \sqrt{\lambda}$$
.

To get:
$$rac{N_k}{N_j} = \left(rac{\mu_k/C_k}{\mu_j/C_j}
ight)$$
 . (Note: $N_1/N > 0!!!$)

Outperforming the Homogeneous Server System

Consider an M/M/N+M system with $\mu=\sum_{k=1}^K q_k\mu_k$. Then $\sqrt{\lambda}P\{ab\}\leq \Delta$ if and only if:

$$\mu \mathsf{N} \ge \lambda + \beta \sqrt{\mu} \sqrt{\lambda}.$$

Compared to:

$$\mu_1 \mathbf{N}_1 + \mu_2 \mathbf{N}_2 + \ldots + \mu_K \mathbf{N}_K \ge \lambda + \beta \sqrt{\mu_1} \sqrt{\lambda}.$$

► Simple Solution to a Difficult Problem:

- ► Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.

- Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.

- Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.
 - Linear boundary for feasible region.

- ▶ Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.
 - Linear boundary for feasible region.
 - FSF is asymptotically optimal (although not exactly optimal)

- Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.
 - Linear boundary for feasible region.
 - ► FSF is asymptotically optimal (although not exactly optimal)
 - Preemptive and non-preemptive policies are asymptotically equivalent (General result: Atar '03)

- ▶ Simple Solution to a Difficult Problem:
 - ► Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.
 - Linear boundary for feasible region.
 - ► FSF is asymptotically optimal (although not exactly optimal)
 - Preemptive and non-preemptive policies are asymptotically equivalent (General result: Atar '03)
 - Asymptotic cost minimization made simple easy to include other constraints.

- Simple Solution to a Difficult Problem:
 - Separation of Staffing and Control.
 - Square-root Safety Staffing is asymptotically optimal.
 - Linear boundary for feasible region.
 - ► FSF is asymptotically optimal (although not exactly optimal)
 - Preemptive and non-preemptive policies are asymptotically equivalent (General result: Atar '03)
 - Asymptotic cost minimization made simple easy to include other constraints.
- ▶ The inverted-V system outperforms its homogeneous server