

On Pooled versus Dedicated Service in the Presence of Triage Errors



Yonit Barron, Opher Baron

2024 Sandra Rotman Centre for Health Sector Strategy

Thx: Ren Yi

Access to Care- Long Lines, Long Waits

Waiting [is] a defining characteristic of Canadian health care. ... A median waiting time of **27.7 weeks** [between referral and treatment in 2023] , 27.4 weeks [in 2022], 9.3 weeks [1993].¹

FRASER
INSTITUTE

The median wait time far **exceeds two hours** in some states. The rate of ... visits has increased significantly... admitted patients in the nation's capital wait a median of **286 minutes** for their room in the hospital.

U.S. News & WORLD REPORT

CNN Money

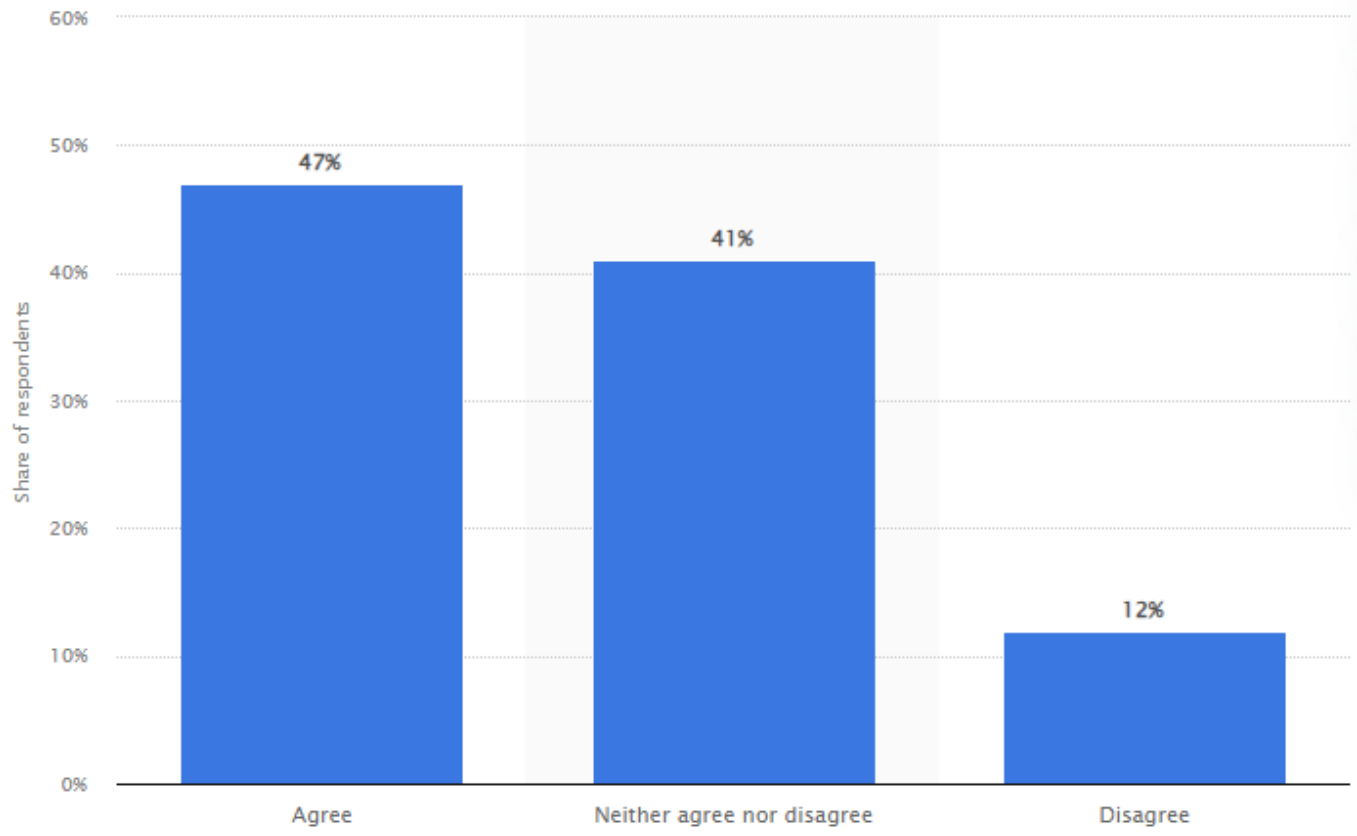
The average time to see (any) dermatologist is **72 days** in Boston, **56 days** in Minneapolis, and only **14 days** in San Diego.

How long will you wait to see a doctor?

Here's the average wait time for a new patient to see a doctor for a non-emergency issue.

	Wait Time <small>What does this show me?</small>	Expected Length of Stay <small>What does this show me?</small>	Status <small>What does this show me?</small>
Vancouver & Area			
★ Vancouver General Hospital <small>Patients of ages 17 and older seen</small>	06:24	07:15	✓
★ St. Paul's Hospital <small>Patients of all ages seen</small>	00:40	05:12	✓
★ Mount Saint Joseph Hospital <small>Patients of all ages seen</small>		Currently closed	
★ UBC Hospital (UBCH) <small>Patients of all ages seen</small>		Currently closed	
★ City Centre Urgent & Primary Care Centre <small>Patients of all ages seen UPCC is for mild to moderate illness</small>		Currently closed	
★ REACH Urgent and Primary Care Centre <small>Patients of all ages seen (lab & x-ray offsite) UPCC is for mild to moderate illness</small>		Currently closed	
★ Northeast Urgent and Primary Care Centre <small>Patients of all ages seen (lab & x-ray offsite) UPCC is for mild to moderate illness</small>		Currently closed	
★ Southeast Urgent and Primary Care Centre <small>Patients of all ages seen (lab & x-ray offsite) UPCC is for mild to moderate illness</small>		Currently closed	
★ BC Children's Hospital <small>Patients seen up to age 16</small>	07:38	06:35	✓

Share of individuals perceiving the waiting times to get an appointment with doctors as too long in Japan in 2018



- ★
- 🔔
- ⚙️
- 🔗
- “ ”
- 🖨️

DOWNLOAD

📄 PDF +
📊 XLS +
🖼️ PNG +
📄 PPT +

Source

- [Show sources information](#)
- [Show publisher information](#)
- [Use Ask Statista Research Service](#)

Release date

July 2018

Region

Japan

Survey time period

May 25 to June 8, 2018

Number of respondents

Approx. 1,000

Age group

16-64 years

[Additional Information](#)

© Statista 2024 [Show source](#)

Access to Medical Care in Emergency Departments in America, Asia, and, well, Globally

- EDs serve different types of patients
- Triage: emergency ($\leq 3\%$), acute ($\sim 20\%$), non acute ($\sim 80\%$)
- Wait targets for patients of different acuity level
 - CTAS: TPIA for each patients' type, measures of LOS
 - Chinese govt. guideline on triage motivate EDs to provide a high level of service to patients
- Fundamental queueing insight: pooling is effective [Smith & Whitt, 1981]
 - Ignores multiple customer types with different targets and priorities
 - Pooling may not be helpful in EDs [Song et al 2019]
- Importantly, in EDs this insight ignores triage and **triage errors**
 - Nurses from 4 Swiss hospitals triage only 59.6% of the patients correctly [Jordi, et al 2015]
 - For elderly patients, 117 out of 519 cases were assigned to a lower type [Grossmann 2012]

Model Description

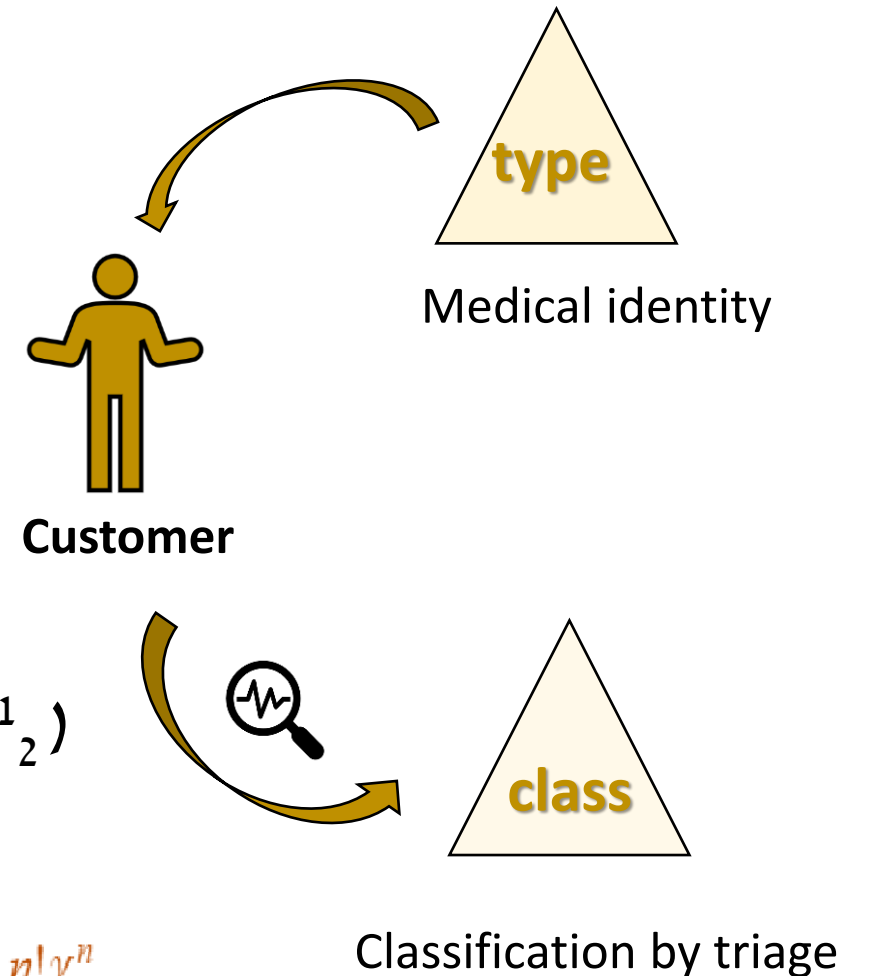
- ED system
- Acute (type 1), non-acute (type 2) patients
- Poisson arrival, rate $\lambda_i, i=1,2$.
- Workload brought $S_i \sim \exp(1/\gamma_i)$
- Triage: p_{ij} = type **i** is classified as type **j**
- System's moments (total workload $L^1 = L^1_1 + L^1_2$)

$$L^r_j := p_{1j} \lambda_1 (\gamma_1)^r + p_{2j} \lambda_2 (\gamma_2)^r$$

- Capacity $\mu_i, i=0$ (pooled), 1,2.

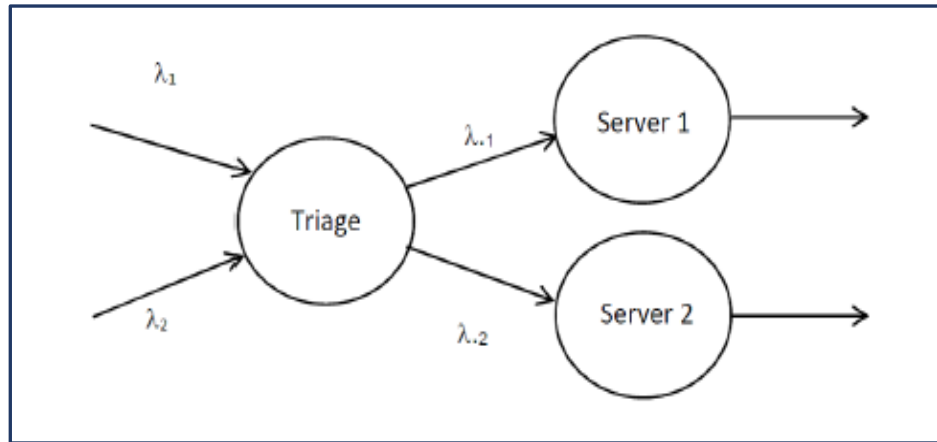
$$Y \sim \exp(\mu/\gamma),$$

$$E[Y^n] = \frac{E[S^n]}{\mu^n} = \frac{n! \gamma^n}{\mu^n}.$$



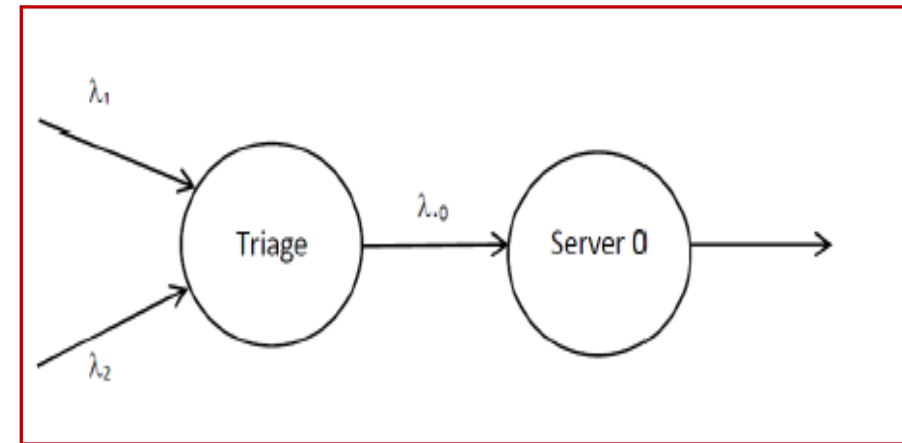
Two M/M/1 like models:

Dedicated system



FCFS

Pooled system



FCFS

Priority (PR)

(non-preemptive)

Objective functions and SL Constraints

w_i - target for expected waiting time of type i customer

W_i - (realized) expected waiting time of type i customer

c_j - capacity cost for one unit of workload of server j

Dedicated Cost

$$\min C_d = c^1 \mu^1 + c^2 \mu^2$$

$$\text{s.t. } W_i \leq w_i, \forall i,$$

Pooled Cost FCFS

$$\min C_p = c^0 \mu^0$$

$$\text{s.t. } W \leq w_1.$$

Pooled Cost PR

$$\min C_p = c^0 \mu^0$$

$$\text{s.t. } W_1 \leq w_1.$$

$$W_2 \leq w_2$$

Our *study*

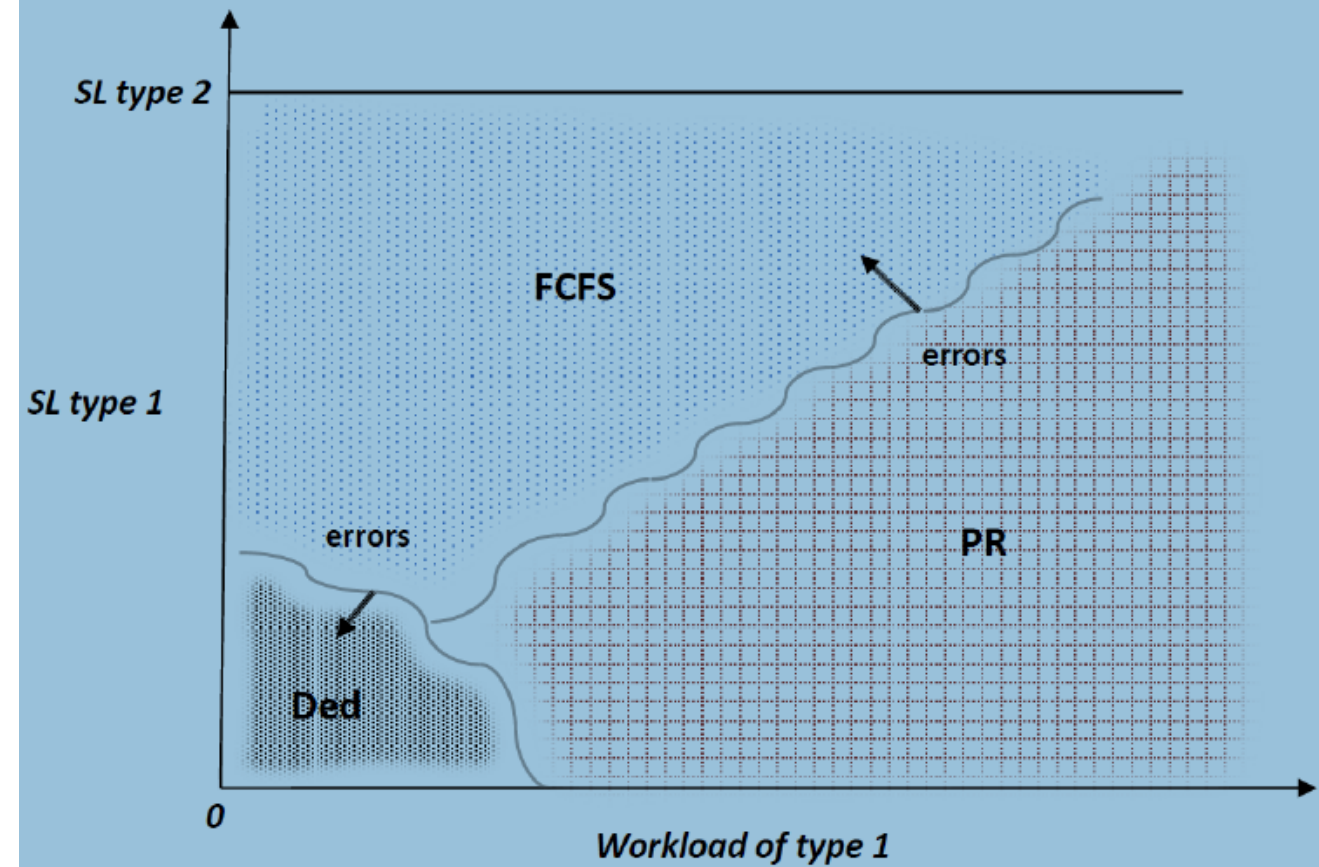
- Deriving the optimal capacities μ_0, μ_1, μ_2 *under SL constraints*
- *Compare dedicated system vs. pooled systems (FCFS and PR policies)*
- The impact of triage errors

Three points of view:

- I. Servers' point of view (capacities, utilizations, waiting time at server)*
 - II. Customers' point of view (waiting time observed by customers)*
 - III. System's point of view (total cost)*
- Many comparisons...

Qualitative Areas where Each Policy is Cost Minimizing System's POV

Triage errors significantly impact the optimal dedicated system, slightly impact the optimal PR pooled system, and have no impact on the FCFS pooled system => arrows



The dedicated system w. errors may be better **jointly** on all 3 POV:

1. Servers: lower utilization (one capacity increases)
2. Customers: are served faster (one type strictly faster)
3. System: costs are lower

Analysis with no triage errors

Optimal Dedicated system

Safety stock of servers 1 and 2

$$\begin{aligned} C_d^* &= c_1 \left(\frac{\lambda_1 \gamma_1}{2} + \sqrt{\frac{\lambda_1 \gamma_1^2}{w_1} + \left(\frac{\lambda_1 \gamma_1}{2} \right)^2} \right) + c_2 \left(\frac{\lambda_2 \gamma_2}{2} + \sqrt{\frac{\lambda_2 \gamma_2^2}{w_2} + \left(\frac{\lambda_2 \gamma_2}{2} \right)^2} \right) \\ &= c_1 \left(\lambda_1 \gamma_1 + \frac{\lambda_1 \gamma_1^2}{w_1 \left(\sqrt{\frac{\lambda_1 \gamma_1^2}{w_1} + \left(\frac{\lambda_1 \gamma_1}{2} \right)^2} + \frac{\lambda_1 \gamma_1}{2} \right)} \right) + c_2 \left(\lambda_2 \gamma_2 + \frac{\lambda_2 \gamma_2^2}{w_2 \left(\sqrt{\frac{\lambda_2 \gamma_2^2}{w_2} + \left(\frac{\lambda_2 \gamma_2}{2} \right)^2} + \frac{\lambda_2 \gamma_2}{2} \right)} \right) \end{aligned}$$

Analysis with no triage errors

Optimal Pooled FCFS

$$C_{FCFS}^* = c_0 \left(\frac{L}{2} + \sqrt{\frac{L^{(2)}}{w_1} + \frac{L^2}{4}} \right) = c_0 \left(L + \frac{L^{(2)}}{w_1 \left(\frac{L}{2} + \sqrt{\frac{L^{(2)}}{w_1} + \frac{L^2}{4}} \right)} \right)$$

Safety stock

$$W_1 = W_2 = W_{FCFS} \leq w_1$$

Fixed $w_1 \longrightarrow C_{FCFS}^*$ independent of w_2

Analysis with no triage errors

Optimal Pooled PR

$$C_{PR}^* = \begin{cases} (i) \ c_0 \left(\frac{\lambda_1 \gamma_1}{2} + \sqrt{\frac{L^{(2)}}{w_1} + \left(\frac{\lambda_1 \gamma_1}{2} \right)^2} \right) = c_0 \left(\lambda_1 \gamma_1 + \frac{L^{(2)}}{w_1 \left(\sqrt{\frac{L^{(2)}}{w_1} + \left(\frac{\lambda_1 \gamma_1}{2} \right)^2} + \frac{\lambda_1 \gamma_1}{2} \right)} \right), & w_1 \leq w_1^*, \\ (ii) \ c_0 \left(\lambda_1 \gamma_1 + \frac{\lambda_2 \gamma_2}{2} + \sqrt{\frac{L^{(2)}}{w_2} + \left(\frac{\lambda_2 \gamma_2}{2} \right)^2} \right) = c_0 \left(L + \frac{L^{(2)}}{w_2 \left(\sqrt{\frac{L^{(2)}}{w_2} + \left(\frac{\lambda_2 \gamma_2}{2} \right)^2} + \frac{\lambda_2 \gamma_2}{2} \right)} \right), & w_1 > w_1^*, \end{cases}$$

Type 1 dominates

where

$$w_1^* = \frac{(d-1)^2 L^{(2)}}{d(\lambda_1 \gamma_1 + d\lambda_2 \gamma_2)L}$$

$$d = w_1/w_2, \ (d=1 \Rightarrow w_1^*=0)$$

Both customers impact

Systems' comparisons-no triage errors



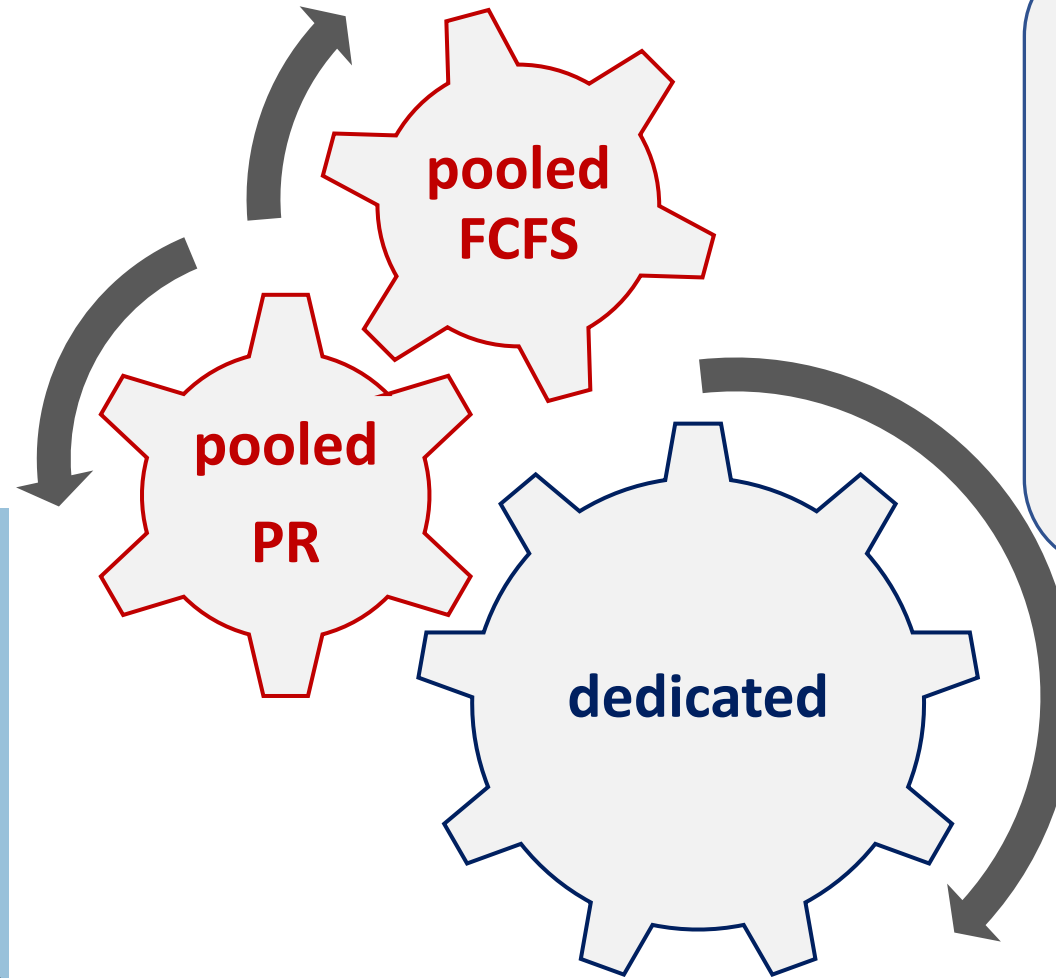
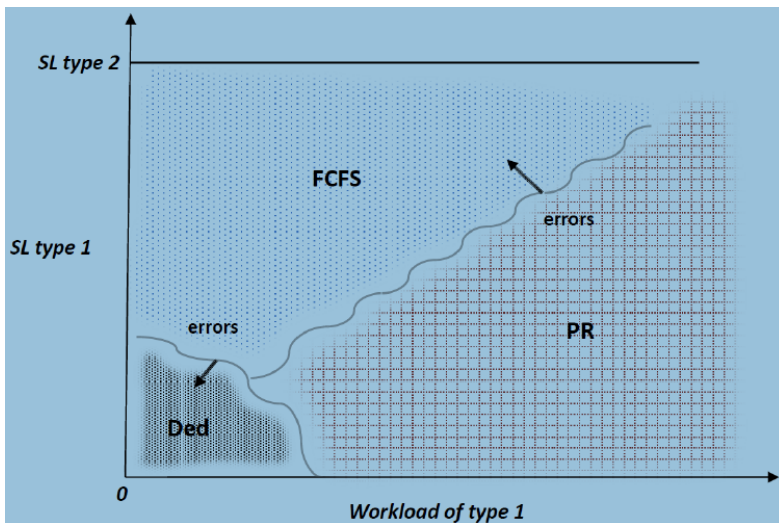
Conclusions:

$w_1 = w_2$ \Rightarrow **Pooled FCFS**

$w_1 \ll w_2$ \Rightarrow **Pooled PR**

$w_1 \rightarrow \infty$ \Rightarrow **Similar**

or $w_1 = 0$ \Rightarrow **Similar**



w_1 is strict (low)
workload of type 1 small
 \downarrow
 μ_0 increases
The benefit from pooling
decreases

Analysis With Triage Errors

Better or worse?

Triage errors highlight **two** effects:

- Servers' workload patters
- Customers service pattern

$$L_j = p_{1j} \lambda_1 \gamma_1 + p_{2j} \lambda_2 \gamma_2$$

Workload of server ***j***
(dedicated systems)

$$W_i = p_{ii} W_{.i} + p_{ij} W_{.j}$$

Service pattern => waits of
customer type ***i***

SL is for **Realized** waits

Analysis With Triage Errors

Optimal Dedicated system

$$\min_{\mu_1, \mu_2} \{C_d = c_1\mu_1 + c_2\mu_2\}$$

s.t.

$$W_i = p_{ii} \frac{L_i^{(2)}}{(\mu_i - L_i)\mu_i} + p_{ij} \frac{L_j^{(2)}}{(\mu_j - L_j)\mu_j} \leq w_i, \quad i = 1, 2$$

$$\mu_j > L_j, \quad j = 1, 2.$$

$$\min_{W_1, W_2} C_d = c_i \left(\frac{L_i}{2} + \sqrt{\left(\frac{L_i}{2}\right)^2 + \frac{L_i^{(2)}}{W_i}} \right) + c_j \left(\frac{L_j}{2} + \sqrt{\left(\frac{L_j}{2}\right)^2 + \frac{L_j^{(2)}}{W_j}} \right)$$

s.t.

$$(a) W_i = p_{ii}W_i + p_{ij}W_j \leq w_i,$$

$$(b) W_j = p_{ji}W_i + p_{jj}W_j \leq w_j.$$

W_i, W_j decision variables

The optimal waits have
at least one tight SL

$$W_i = w_i$$

or

$$W_j = w_j$$

Solution steps (dedicated system...)

- ✓ Transform the problem to a single decision variable w
- ✓ Under some constraints, the objective function $C_d^*(w)$ is convex in w

↓
 $C_d^*(w)$ has a unique minimizer \bar{w}_i

Proposition:

(a) If $\bar{w}_i \geq \hat{W}(w_i)$ for both $i = 1$ and $i = 2$, the optimal solution is $W_{.i}^* = \hat{W}(w_i), i = 1, 2$.

$$W_{.i}^* = w_i, W_{.j}^* = w_j$$

(b) Otherwise,

$$W_{.i}^* = \bar{w}_i, W_{.j}^* = \frac{w_j - p_{ji}\bar{w}_i}{p_{jj}};$$

$$W_{.i}^* < w_i, W_{.j}^* = w_j$$

Conclusions (dedicated system)



The interplay between the changes of the *workload & service patterns* effects the performance of the dedicated system



The system w. errors may be better **jointly** on all 3 POV:

1. Servers: lower utilization (one capacity increases)
2. Customers: are served faster (one type strictly faster)
3. System: costs are lower



Pondering

Intuitively, this may occur when the cheaper server becomes busier due to triage errors. But, this happens more generally!

Analysis With Triage Errors

Pooled FCFS: no analysis 😊

Pooled PR: (no closed form)

There exist a unique $\mu_i > L$ that solves $W_i = w_i$, $i=1,2$. The minimized cost for the PR w. triage errors is

$$C_{PR} = c_0 \mu_0^*, \text{ where } \mu_0^* = \max\{\mu_i\}.$$

Some insights (pooled system...)

Comparing FCFS and PR pooled systems w. errors:

- (i) When $p_{ij} = 0.5$ for $i=1,2$ (and all other parameters are the same), FCFS and PR have the same capacity and cost, i.e., $\mu_0^* = \mu_{FCFS}^*$;
- (ii) ow., FCFS may perform strictly better only when the SL of type 2 customers is tight under PR, i.e., when under PR $W_2^* = w_2$ holds.

Numerical observations:

1. FCFS is better than PR when the proportion of type 1 customers is large

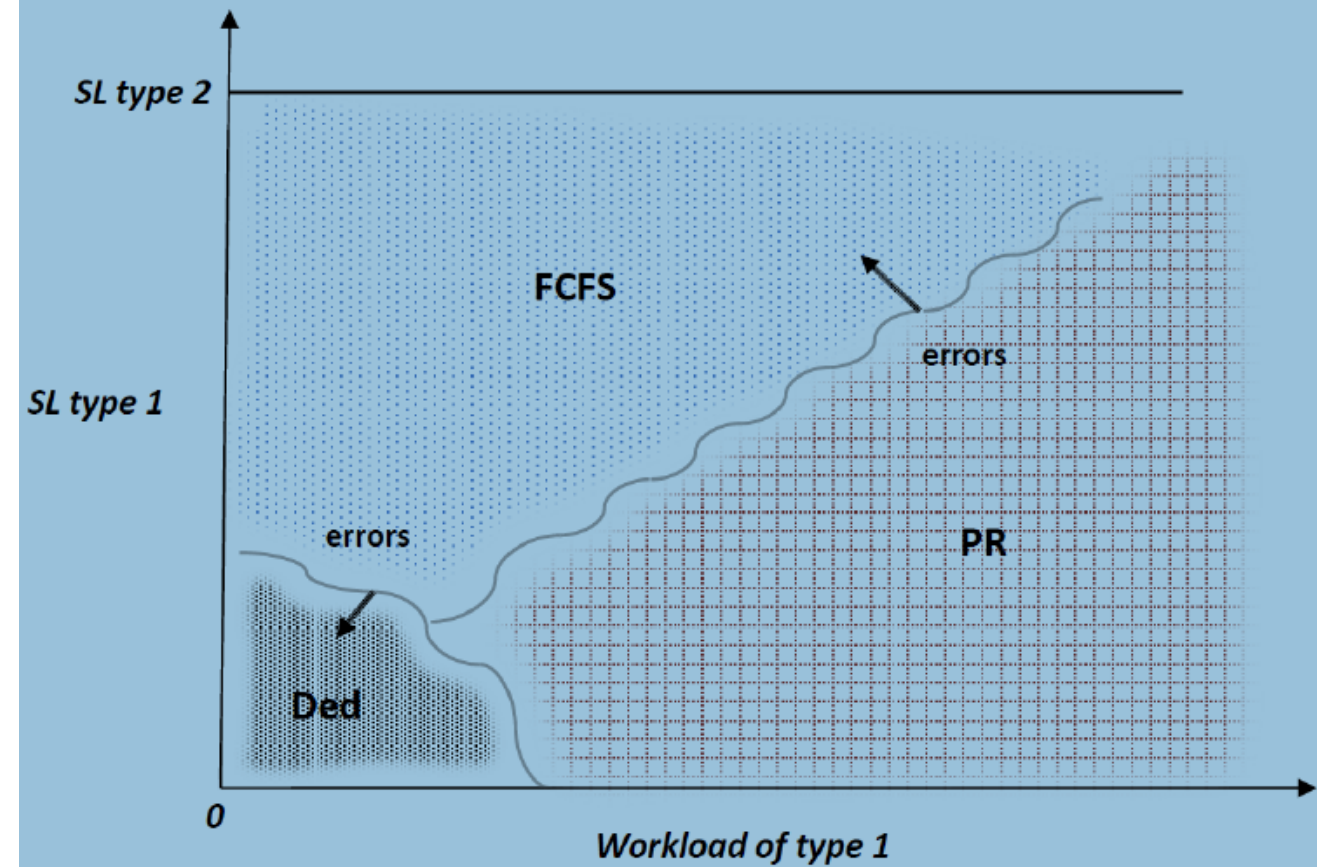
Intuition: Increasing the workload of type 1, increases the delay of type 2. To meet their SL, the PR server increases its capacity...

2. FCFS is better than PR when SL are identical, PR is better when these differ

Intuition: under PR, low priority customers face longer waits=>higher capacity (&cost)

Thx!

Triage errors significantly impact the optimal dedicated system, slightly impact the optimal PR pooled system, and have no impact on the FCFS pooled system => arrows



The dedicated system w. errors may be better **jointly** on all 3 POV:

1. Servers: lower utilization (one capacity increases)
2. Customers: are served faster (one type strictly faster)
3. System: costs are lower