

# THE COST OF IMPATIENCE IN DYNAMIC MATCHING: SCALING LAWS AND OPERATING REGIMES

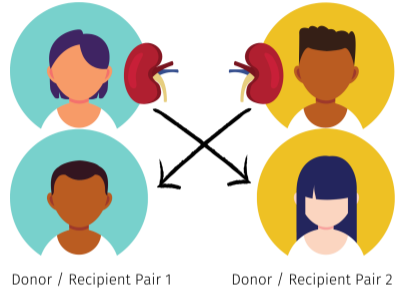
---

Angela Kohlenberg and Itai Gurvich

## BLOOD ALLOCATION



## PAIRED KIDNEY EXCHANGE



## ■ HOW DOES IMPATIENCE IMPACT MATCH RATE?



THE RELATIONSHIP BETWEEN PARAMETERS

## ■ NEED TO CONSIDER: TIME TO ABANDON & TIME TO MATCH

# HOW DOES IMPATIENCE IMPACT MATCH RATE?

1. **CLASSIFY SETTINGS** based on how impatience impacts match loss  
[OPERATING REGIMES]
2. **IDENTIFY KEY DETERMINANTS** of match loss from impatience  
[SCALING LAWS]



## RELATED LITERATURE (A SAMPLE): MATCHING AND SINGLE-SIDED QUEUES WITH ABANDONMENT

### EXACT RESULTS

- Transient and steady-state performance of a two-sided queue (Conolly et al. 2002, Afèche et al. 2014, Diamant and Baron 2019)
- Performance analysis / stability region under specific policy (Castro et al. 2020b, Zubeldia et al. 2022)

ALL PARAMETERS ARE EQUALLY  
IMPORTANT

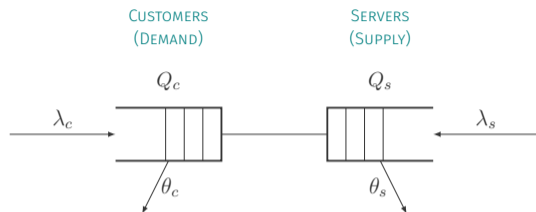
### APPROXIMATIONS (MOSTLY HEAVY TRAFFIC)

- Performance analysis and control of a two-sided queue (Liu et al. 2015, Büke and Chen 2017, Chen and Hu 2020, Aveklouris et al. 2023)
- Performance analysis and control of a single-sided queue (Ward and Glynn 2003, Lee and Ward 2019)
- Control policies for a matching network (Aveklouris et al. 2021, Collina et al. 2020, Aouad and Saritaç 2022, Castro et al. 2020a, Wang et al. 2022)

IMPOSES SPECIFIC RELATIONSHIP  
ON PARAMETERS

OUR FOCUS: SIMPLE MODEL, UNIVERSAL RESULTS (IN PARAMETERS) & GENERAL CHARACTERIZATION

# THE SIMPLEST MATCHING MODEL: TWO-SIDED QUEUE WITH EXPONENTIAL DISTRIBUTIONS



MATCH RATE WITHOUT ABANDONMENT (UPPER BOUND):  
 $\min \{ \lambda_c, \lambda_s \} = \lambda_c$

**LABEL:**  $\lambda_c \leq \lambda_s$

**UTILIZATION (MEASURE OF EXCESS CAPACITY):**

$$\rho = \lambda_c / \lambda_s$$

**ACTUAL MATCH RATE:**

$$d = \lim_{t \uparrow \infty} \frac{1}{t} \mathbb{E}[D(t)]$$

**ARRIVALS = MATCHES + ABANDONMENTS:**

$$\lambda_c = d + \theta_c \mathbb{E}[Q_c]$$

## Definition (Cost-of-Impatience, Col).

$$\text{Col} = \text{No-abandonment match rate } (\lambda_c) - \text{Actual match rate } (d) = \theta_c \mathbb{E}[Q_c] = \theta_s \mathbb{E}[Q_s] - (\lambda_s - \lambda_c)$$

# HOW DOES IMPATIENCE IMPACT MATCH RATE?

1. CLASSIFY SETTINGS based on how impatience impacts match loss

[OPERATING REGIMES]

2. IDENTIFY KEY DETERMINANTS of match loss from impatience

[SCALING LAWS]



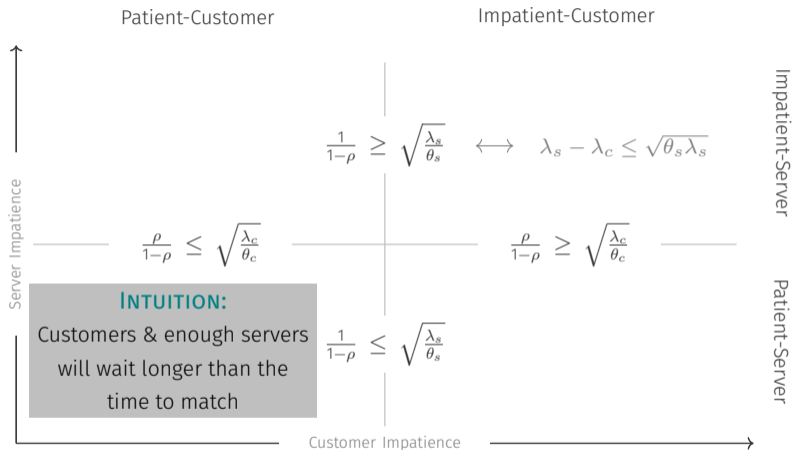
# OPERATING REGIMES: A CLASSIFICATION BY RELATIVE IMPATIENCE



■ **PATIENT VS. IMPATIENT:** MEASURE OF MEAN PATIENCE RELATIVE TO EXCESS CAPACITY ( $\rho = \lambda_c / \lambda_s$ )

■ THE COI BEHAVES DIFFERENT IN EACH OPERATING REGIME

# OPERATING REGIMES: A CLASSIFICATION BY RELATIVE IMPATIENCE

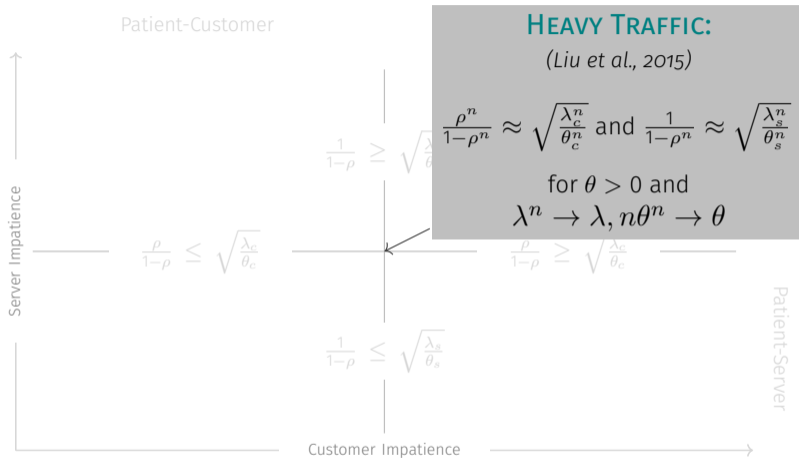


■ **PATIENT VS. IMPATIENT:** MEASURE OF MEAN PATIENCE RELATIVE TO EXCESS CAPACITY ( $\rho = \lambda_c/\lambda_s$ )

■ THE COI BEHAVES DIFFERENT IN EACH OPERATING REGIME



# OPERATING REGIMES: A CLASSIFICATION BY RELATIVE IMPATIENCE

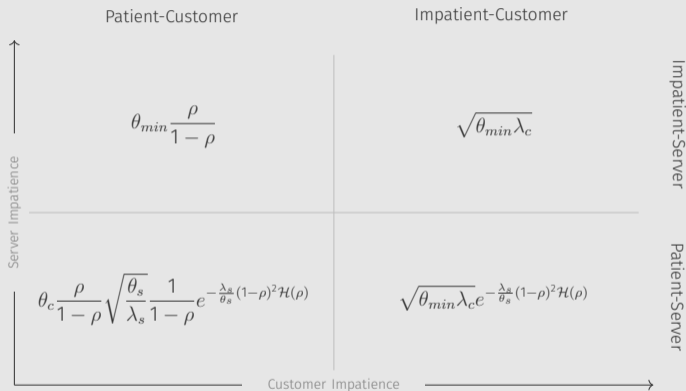


■ **PATIENT VS. IMPATIENT:** MEASURE OF MEAN PATIENCE RELATIVE TO EXCESS CAPACITY ( $\rho = \lambda_c/\lambda_s$ )

■ THE COI BEHAVES DIFFERENT IN EACH OPERATING REGIME

# KEY DETERMINANTS OF MATCH LOSS FROM IMPATIENCE: SNEAK PEAK

## Theorem (Cost-of-Impatience Scaling by Operating Regime).



■ SCALING LAW,  $\mathcal{S} \sim \text{COI}$ : CHARACTERIZE HOW THE COI CHANGES AS A FUNCTION OF PARAMETERS

# HOW DOES IMPATIENCE IMPACT MATCH RATE?

1. CLASSIFY SETTINGS based on how impatience impacts match loss  
[OPERATING REGIMES]
2. IDENTIFY KEY DETERMINANTS of match loss from impatience  
[SCALING LAWS]



## SCALING LAW DEFINITION: A UNIVERSAL APPROXIMATION

Definition (Scaling Law).

$$\text{Col} = \lambda_c - d \sim \mathcal{S}$$

when, for “all” parameter combinations,

$$\frac{1}{\Gamma} \leq \frac{\text{Col}(\lambda, \theta)}{\mathcal{S}(\lambda, \theta)} \leq \Gamma$$

for some function  $\mathcal{S}(\lambda, \theta)$  and a constant  $\Gamma \geq 1$  that does not depend on  $\lambda = (\lambda_c, \lambda_s), \theta = (\theta_c, \theta_s)$ . Recall that  $d$  denotes the actual match rate.

- **SCALING LAW:** CHARACTERIZE HOW THE COI CHANGES AS A FUNCTION OF PARAMETERS
- **GOAL:** IDENTIFY THE **RELATIVE IMPORTANCE** OF EACH PARAMETER ON MATCH LOSS

## Theorem (Universal Cost-of-Impatience Scaling).

Col  $\sim \mathcal{S}$

$$= \theta_c \min \left\{ \frac{\rho}{1-\rho}, \sqrt{\frac{\lambda_c}{\theta_c}} \right\} \left( 1 + \left[ 1 + \frac{\rho}{1-\rho} \frac{\theta_c}{\lambda_c} \min \left\{ \frac{\rho}{1-\rho}, \sqrt{\frac{\lambda_c}{\theta_c}} \right\} \right] \sqrt{\frac{\lambda_s}{\theta_s}} (1-\rho) e^{\frac{\lambda_s}{\theta_s} (1-\rho)^2 \mathcal{H}(\rho)} \right)^{-1}$$

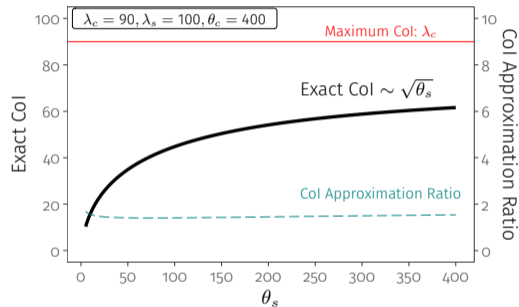
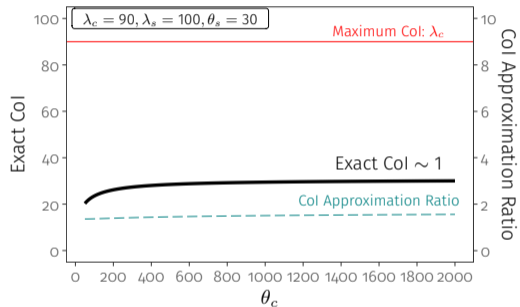
where  $\rho = \lambda_c/\lambda_s$  and  $\mathcal{H}(\rho) = \sum_{n=1}^{\infty} \frac{1}{n(n+1)} (1-\rho)^{n-1}$ .

**EXACT COI:**  $\text{Col} = \theta_c \sum_{n=1}^{\infty} n \prod_{i=1}^n \frac{\lambda_c}{\lambda_s + i\theta_c} \left( 1 + \sum_{n=1}^{\infty} \prod_{i=1}^n \frac{\lambda_c + i\theta_s}{\lambda_s} + \sum_{n=1}^{\infty} \prod_{i=1}^n \frac{\lambda_c}{\lambda_s + i\theta_c} \right)^{-1}$

**PROOF CONCEPT:** COUPLING ARGUMENTS AND EXPANSIONS OF EXPLICIT EXPRESSIONS TO UPPER- AND LOWER-BOUND STEADY-STATE DISTRIBUTIONS

# UNIVERSAL COI SCALING LAW: AN EXAMPLE

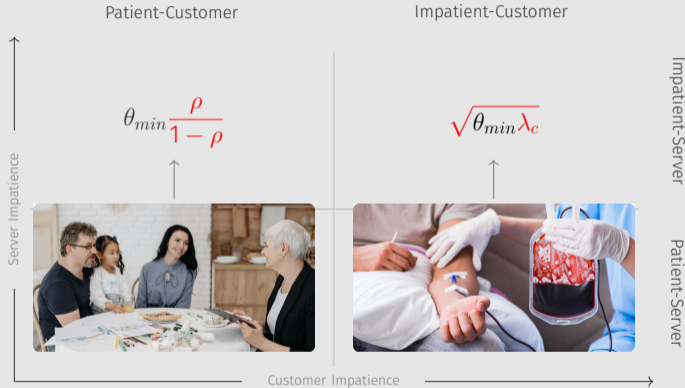
## IMPATIENT-CUSTOMER, IMPATIENT-SERVER REGIME:



- SCALING LAW,  $\mathcal{S}$ : CHARACTERIZE HOW THE COI CHANGES AS A FUNCTION OF PARAMETERS
- COI APPROXIMATION RATIO =  $\mathcal{S} / \text{Exact Col}$

# “WINNER-TAKE-ALL” COMPETITION BETWEEN EXCESS CAPACITY AND IMPATIENCE

## Theorem (Cost-of-Impatience Scaling by Operating Regime).



# “WINNER-TAKE-ALL” COMPETITION BETWEEN EXCESS CAPACITY AND IMPATIENCE

$$\mathbb{E}[Q_c | Q_s = 0] \sim \min \left\{ \frac{\rho}{1-\rho}, \sqrt{\frac{\lambda_c}{\theta_c}} \right\}$$

EXCESS CAPACITY “WINS”: M/M/1 QUEUE

IGNORES ABANDONMENT



IMPATIENCE “WINS”: M/M/1+M QUEUE

CRITICALLY-LOADED, IGNORES EXCESS CAPACITY

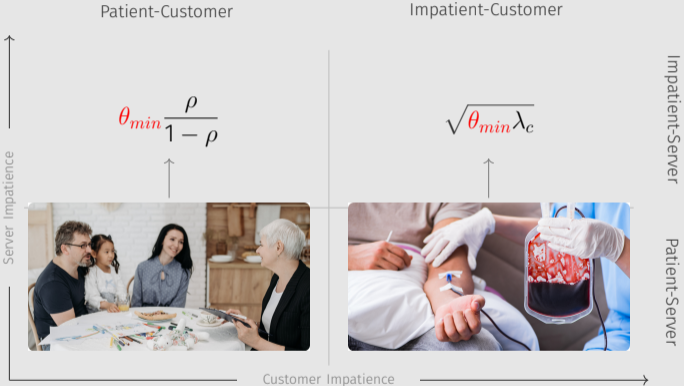


- IN HEAVY TRAFFIC,  $M/M/1 + M \sim M/M/1$  WHEN  $\sqrt{\theta} \ll 1 - \rho$  (WARD AND GLYNN, 2003)
- ONLY EXCESS CAPACITY OR IMPATIENCE MATTERS - NOT BOTH
- PATIENT CUSTOMERS: COI SENSITIVE TO SMALL CHANGES IN EXCESS CAPACITY,  $\rho$



# ABILITY TO ACCUMULATE “INVENTORY” OF SERVERS

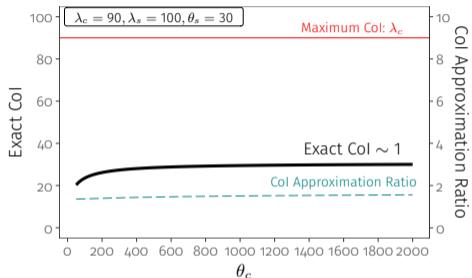
## Theorem (Cost-of-Impatience Scaling by Operating Regime).



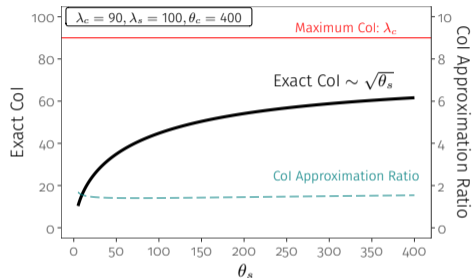
■ PATIENT CUSTOMERS: NOT JUST REDUCTION TO SINGLE-SIDED CUSTOMER QUEUE

# ONLY THE MOST PATIENT TYPE MATTERS ( $\min\{\theta_c, \theta_s\}$ )

## MAXIMUM ABANDONMENT RATE ( $\theta_c$ ):



## MINIMUM ABANDONMENT RATE ( $\theta_s$ ):



- IT ONLY MATTERS THAT ONE TYPE IS “PATIENT ENOUGH”
- LOWER MATCH LOSS  $\implies$  FOCUS ON MOST PATIENT TYPE
- WAITING SERVERS OFFSET MATCH LOSS



## Theorem (Cost-of-Impatience Scaling by Operating Regime).



### WHAT IS THE OPTIMAL CAPACITY LEVEL?

# OPTIMAL CAPACITY: LOGARITHMIC SCALING IN ABANDONMENT COST

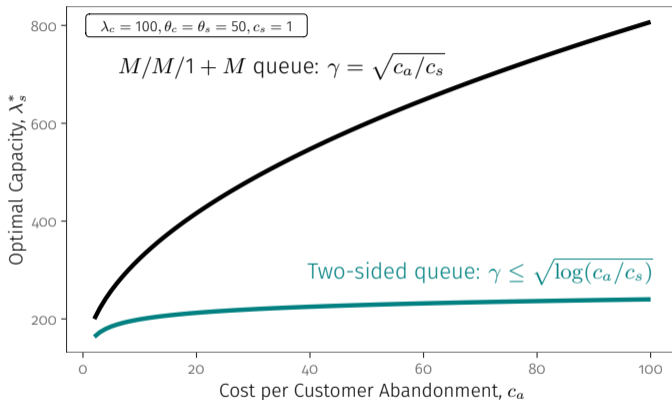
## Lemma (Optimal Capacity Scaling).

The optimal safety capacity,  $\lambda_s^* - \lambda_c$ , that balances abandonment and capacity costs,  $c_a$  and  $c_s$ , is:

$$\lambda_s^* - \lambda_c \sim \gamma \sqrt{\theta_{\min} \lambda_c}$$

where

$$\lambda_s^* = \operatorname{argmin}_{\lambda_s \geq \lambda_c} \{c_a \theta_c \mathbb{E}[Q_c] + c_s \lambda_s\}.$$



■ ABILITY TO HOLD INVENTORY OF SERVERS  $\implies$  SLOWER SCALING OF CAPACITY

# ONGOING WORK: TRADE-OFF BETWEEN MATCH QUALITY AND EFFICIENCY

## SUPPLY



TYPE O- (UNIVERSAL DONOR)

## DEMAND



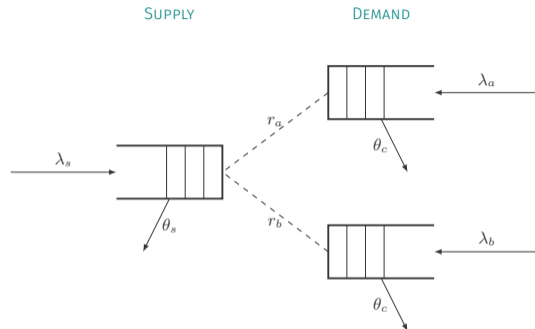
TYPE O- (HIGH REWARD)



TYPE AB+ (LOW REWARD)

■ WHAT IS THE OPTIMAL TIMING OF MATCHES?

# THE SIMPLEST MATCHING MODEL INVOLVING A MATCH DECISION



- **GOAL:** IDENTIFY NEARLY OPTIMAL **POLICIES** (SPECIFIC POLICY FOR ANY COMBINATION OF PARAMETERS)
- **SCALING LAWS:** LOWER BOUND ON MATCH LOSS
- **OUR FOCUS:** SIMPLE MODEL, UNIVERSAL RESULTS (IN PARAMETERS) & GENERAL CHARACTERIZATION

THANK YOU!

---

### REFERENCES

---

- Afèche, P., Diamant, A., and Milner, J. (2014). Double-sided batch queues with abandonment: Modeling crossing networks. *Operations Research*, 62(5):1179–1201.
- Aouad, A. and Saritaç, Ö. (2022). Dynamic Stochastic Matching under Limited Time. *Operations Research*, 70(4):2349–2383.
- Aveklouris, A., DeValve, L., and Ward, A. R. (2021). Matching Impatient and Heterogeneous Demand and Supply.
- Aveklouris, A., Puha, A. L., and Ward, A. R. (2023). A fluid approximation for a matching model with general renegeing distributions. *Queueing Systems*.
- Büke, B. and Chen, H. (2017). Fluid and diffusion approximations of probabilistic matching systems. *Queueing Systems*, 86:1–33.



## REFERENCES II

- Castro, F., Frazier, P., Ma, H., Nazerzadeh, H., and Yan, C. (2020a). Matching queues, flexibility and incentives. Available at SSRN.
- Castro, F., Nazerzadeh, H., and Yan, C. (2020b). Matching queues with renegeing: a product form solution. *Queueing Systems*, 96(3-4):359–385.
- Chen, Y. and Hu, M. (2020). Pricing and matching with forward-looking buyers and sellers. *Manufacturing & Service Operations Management*, 22(4):717–734.
- Collina, N., Immorlica, N., Leyton-Brown, K., Lucier, B., and Newman, N. (2020). Dynamic Weighted Matching with Heterogeneous Arrival and Departure Rates. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 12495 LNCS.
- Conolly, B., Parthasarathy, P., and Selvaraju, N. (2002). Double-ended queues with impatience. *Computers & Operations Research*, 29(14):2053–2072.
- Diamant, A. and Baron, O. (2019). Double-sided matching queues: Priority and impatient customers. *Operations Research Letters*, 47(3):219–224.

## REFERENCES III

- Lee, C. and Ward, A. R. (2019). Pricing and capacity sizing of a service facility: Customer abandonment effects. *Production and Operations Management*, 28(8):2031–2043.
- Liu, X., Gong, Q., and Kulkarni, V. G. (2015). Diffusion Models for Double-Ended Queues with Renewal Arrival Processes. *Stochastic Systems*, 5(1):1–61.
- Wang, G., Zhang, H., and Zhang, J. (2022). On-demand ride-matching in a spatial model with abandonment and cancellation. *Operations Research*, page forthcoming.
- Ward, A. R. and Glynn, P. W. (2003). A Diffusion Approximation for a Markovian Queue with Reneging. *Queueing Systems*, 43(1-2).
- Zubeldia, M., Jhunjhunwala, P. R., and Maguluri, S. T. (2022). Matching queues with abandonments in quantum switches: Stability and throughput analysis. arXiv preprint arXiv:2209.12324.