



#### I. INTRODUCTION

- A Two-tier Service System: Two service providers (SPs) free but long wait time v.s. toll but short wait time.
  - Public hospital v.s. private hospital
  - Regular service v.s. expediting service with a fee
- Strategic Customers:
  - Delay/cost sensitive
  - Heterogeneous
  - Selecting an SP to obtain the service based on "congestion information"
- Information Scenarios
  - Real-time queue length at an arrival instant.
  - Long term statistics (or no real-time information)



#### **Research Issues**

- How does the real-time queue length information affect the system performance?
- Is there any significant difference in system performance behavior between with real-time information and without real-time information scenarios? (long-asked question by both practitioners and researchers in academia)

#### **Related Literature**

- Noar (1969) a single server queueing system with identical customers influenced by either real-time information or the long-term statistics.
- Schroeter (1982) heterogeneous customers with uniformly distributed unittime waiting costs.
- Hassin and Haviv (2003) a book "To Q or Not to Q" a survey on more research in this area
- Armony and Maglaras (2004a, b) service systems where an arriving customer knowing some delay information chooses to balk, wait, or leave a message, in which case the SP calls back within a guaranteed time.
- Guo and Zipkin (2007) a single server queueing system with heterogeneous customers who either join or balk and three levels of information.
- Guo and Zhang (2012a, b) investigate two-tier service systems motivated by healthcare and border-crossing systems where customers can choose to join a free system or a toll system. Either long-terms statistics with two queues or real-time information with one queue.

## **Our Approach**

- Consider more realistic factors
- Use exact analysis supplemented by simulation studies









The system state is defined as  $(X_f(t), X_c(t))$  on the state space

 $\Omega = \{(n,m): n = 0, 1, ...; m = 0, 1, ..., M\}$ 

• Under the steady-state condition, the system reaches steadystate

 $\lim_{t\to\infty} P\{X_f(t)=n, X_c(t)=m)=p_{nm}.$ 

• The equilibrium arrival rates to the two systems:

$$\begin{array}{lll} \lambda_f(n,m) &=& \Lambda F_{\theta}(\frac{p}{n/\mu_f - m/\mu_c}), \\ \lambda_c(n,m) &=& \Lambda \left(1 - F_{\theta}(\frac{p}{n/\mu_f - m/\mu_c})\right) \end{array}$$

for states with m < M.

# Customer self-interest choice behavior => Structure of the QBD generator

• There exists a threshold n<sub>0</sub> for the free queue length

such that  $\lambda_f(n,m) < \varepsilon$  whenever  $n > n_0$  and m < M, where  $\varepsilon$  is a small positive value. The lower bound  $n_0$  can be determined by

$$n_{\mathbf{0}} > \frac{p\mu_f}{F_{\boldsymbol{\theta}}^{-1}(\varepsilon/\Lambda)} + \frac{\mu_f}{\mu_c}M$$

$$n_{\mathbf{0}} = \operatorname{int} \left( \frac{p}{U} \frac{\Lambda \mu_{f}}{\varepsilon} + (M-1) \frac{\mu_{f}}{\mu_{c}} \right)$$

- Using this property, we can develop a level independent QBD process.
- M determines the number of phases
- n<sub>0</sub> determine the number of boundary states of the QBD process.







$$A_{10} = \begin{bmatrix} \mu_f \\ \ddots \\ \mu_f \end{bmatrix} = \mu_f I,$$

$$B_{11} = \begin{bmatrix} -(\mu_f + \Lambda) & \lambda_c(1,0) \\ \mu_c & -(\Lambda + \mu_c + \mu_f) \\ & \ddots & \ddots \\ & & \mu_c & -(\Lambda + \mu_c + \mu_f) \end{bmatrix},$$

$$C_{12} = \begin{bmatrix} \lambda_f(1,0) \\ \Lambda \\ & \ddots \\ & & \Lambda \end{bmatrix}$$

$$A_{21} = \begin{bmatrix} \mu_{f} & & \\ & \ddots & \\ & & \mu_{f} \end{bmatrix} = \mu_{f}I,$$

$$B_{22} = \begin{bmatrix} -(\mu_{f} + \Lambda) & \lambda_{c}(2, 0) & & \\ & \mu_{c} & -(\Lambda + \mu_{c} + \mu_{f}) & \lambda_{c}(2, 1) & \\ & & \mu_{c} & -(\Lambda + \mu_{c} + \mu_{f}) \\ & & & \ddots & \ddots \\ & & & & \mu_{c} & -(\Lambda + \mu_{c} + \mu_{f}) \end{bmatrix},$$

$$C_{23} = \begin{bmatrix} \lambda_{f}(2, 0) & & \\ & \lambda_{f}(2, 1) & & \\ & & & \ddots & \\ & & & & \Lambda \end{bmatrix}$$

$$A_{M,M-1} = \begin{bmatrix} \mu_{f} \\ & \ddots \\ & \mu_{f} \end{bmatrix} = \mu_{f}I,$$

$$B_{M,M} = \begin{bmatrix} -(\mu_{f}+\Lambda) & \lambda_{c}(5,0) \\ & \mu_{c} & -(\Lambda+\mu_{c}+\mu_{f}) & \lambda_{c}(5,1) \\ & & \ddots & \ddots & \ddots \\ & & \mu_{c} & -(\Lambda+\mu_{c}+\mu_{f}) & \lambda_{c}(M,M-1) \\ & & & \mu_{c} & -(\Lambda+\mu_{c}+\mu_{f}) \end{bmatrix},$$

$$C_{M,M+1} = \begin{bmatrix} \lambda_{f}(5,0) & & & \\ & \ddots & & \\ & \lambda_{f}(M,M-1) & & \\ & & & \Lambda \end{bmatrix}.$$



For 
$$n \ge n_0$$
,  

$$A = \begin{bmatrix} \mu_f \\ & \ddots \\ & \mu_f \end{bmatrix} = \mu_f I = A_{i,i-1}, \text{ for } i \ge 1.$$

$$B = \begin{bmatrix} -(\mu_f + \Lambda) & \Lambda \\ & \mu_c & -(\Lambda + \mu_c + \mu_f) & \Lambda \\ & & \ddots & \ddots & \ddots \\ & & & \mu_c & -(\Lambda + \mu_c + \mu_f) \end{bmatrix},$$

$$C = \begin{bmatrix} \ddots \\ & \Lambda \end{bmatrix}.$$

### **Stability Condition**

**Proposition 1** With real-time delay information, the two-tier service system reaches the steady state if

$$\mu_f > \frac{\left(1 - \frac{\Lambda}{\mu_c}\right) \left(\frac{\Lambda}{\mu_c}\right)^M \Lambda}{1 - \left(\frac{\Lambda}{\mu_c}\right)^{M+1}}.$$
(3)

**Corollary 2** If  $\Lambda/\mu_c \leq 1$ , as  $M \longrightarrow \infty$ , (3) becomes  $\mu_f > 0$  or there is no requirement for a positive  $\mu_f$ . If  $\Lambda/\mu_c > 1$ , as  $M \longrightarrow \infty$ , (3) becomes  $\mu_f + \mu_c > \Lambda$ .

Under the stability condition, the stationary probability vector is defined as

$$\mathbf{p}_n = [p_{n0}, p_{n1}, ..., p_{nM}],$$

where  $p_{nm} = \lim_{t\to\infty} P\{X_f(t) = n, X_c(t) = m\}$ . We know that when  $n \ge n_0$ , the matrix geometric solution is given by

 $\mathbf{p}_{n+1} = \mathbf{p}_n \mathbf{R}.$ 

Like any regular QBD process, the rate matrix R should satisfy  $R^2A + RB + C = 0$ 



#### Performance measures

The boundary state vector  $\mathbf{p}_0$  has  $n_0$  components (or M + 1 dimensional vectors) and is the unique solution of the equation system of  $\mathbf{p}_0(B_0 + RA) = 0$  and  $\mathbf{p}_0(I - R)^{-1}\mathbf{1} = 0$ .

• Marginal Probabilities:

$$p_{\bullet j} = \sum_{n=0}^{\infty} p_{nj}$$
  $p_{n\bullet} = \sum_{j=0}^{M} p_{nj}$ 

• Expected Queue Lengths:

$$E(L^c) = \sum_{j=0}^M j p_{\bullet j}$$
 and  $E(L^f) = \sum_{n=0}^\infty n p_{n \bullet}$ 











#### However..

- The challenge is that we do not have any explicit expressions for the equilibrium arrival rates.
- A search algorithm has to be used to numerically compute these arrival rates and the resulting expected waiting times.

# Computing the Nash equilibrium performance measures (Detailed explanations are in the paper)

A Search Algorithm of Computing  $W(\lambda_c, \mu_c)$  and  $W(\lambda_f^{eff}, \mu_f)$ :

- Step 1: Initialization Select an initial small (or large) θ<sub>0</sub> for the right (or left) search, compute λ<sub>f</sub> = ΛF<sub>θ</sub>(θ<sub>0</sub>), λ<sub>c</sub> = Λ(1-F<sub>θ</sub>(θ<sub>0</sub>)), π<sub>M</sub> from (6), W(λ<sub>c</sub>, μ<sub>c</sub>) from (5), and W(λ<sub>f</sub><sup>eff</sup>, μ<sub>f</sub>) using the MMPP/M/1 queue algorithm such that λ<sub>f</sub> + λ<sub>c</sub>π<sub>M</sub> < μ<sub>f</sub> and f(θ<sub>0</sub>) > 0. Use a positive Δθ to search to the right (or a negative Δθ to search to the left). Set n = 0.
- Step 2: Let θ<sub>n+1</sub> = θ<sub>n</sub> + Δθ. Use λ<sub>f</sub> = ΛF<sub>θ</sub>(θ<sub>n+1</sub>) and λ<sub>e</sub> = Λ(1 F<sub>θ</sub>(θ<sub>n+1</sub>)) to compute W(λ<sub>e</sub>, μ<sub>e</sub>) from (6) and check the stability condition for the free system, λ<sub>f</sub> + λ<sub>e</sub>π<sub>M</sub> < μ<sub>f</sub>. If the free system is not stable, increase M until the stability condition is satisfied. Then compute W(λ<sup>eff</sup><sub>f</sub>, μ<sub>f</sub>) based on the MMPP/M/1 algorithm. Calculate f(θ<sub>n+1</sub>) and compared with f(θ<sub>n</sub>).
- Step 3: If the sign of f changes between θ<sub>n</sub> and θ<sub>n+1</sub>, the solution to (9), θ
   , exists on the interval (θ<sub>n</sub>, θ<sub>n+1</sub>), go to next step. If the sign of f does not change between θ<sub>n</sub> and θ<sub>n+1</sub>, let n = n + 1, go to Step 2.
- Step 4: If Δθ > ε, let θ<sub>0</sub> = θ<sub>n</sub>, Δθ = Δθ/10, n = 0 go to Step 2. Otherwise, θ
   = (θ<sub>n</sub>+θ<sub>n+1</sub>)/2. Then use λ<sub>f</sub> = ΛF<sub>θ</sub>(θ), λ<sub>c</sub> = Λ(1 F<sub>θ</sub>(θ)) to compute the equilibrium W(λ<sub>c</sub>, μ<sub>c</sub>) using (5) and W(λ<sub>f</sub><sup>#</sup>, μ<sub>f</sub>) using the MMPP/M/1 queue algorithm (see Appendix) and Little's Law. □

















Sa	arvica	a tim		ith a	rand	e of (	$\Omega = 0$			
50					rang		5013			
Dist.	Erlang	k=infty	Erlang	k=9	Erlang	k=8	Erlang	k=7	Erlang	k=6
COV	0		0.333333		0.353553		0.377964		0.408248	
M	E(Lc) Imp	E(Lf) Imp	E(Lc) Imp	E(Lf) Imp	E(Lc) Imp	E(Lf) Imp	E(Lc) Imp	E(Lf) Imp	E(Lc) Imp	E(Lf) Imp
10	-42.85%	-25.70%	-42.73%	-24.80%	-42.38%	-24.33%	-42.73%	-23.95%	-43.37%	-23.41%
9	-42.45%	-23.34%	-42.35%	-22.44%	-42.06%	-21.43%	-43.63%	-23.11%	-43.47%	-21.18%
8	-42.19%	-20.40%	-42.64%	-19.37%	-42.78%	-18.63%	-42.56%	-18.77%	-42.03%	-16.46%
7	-42.21%	-16.77%	-43.13%	-17.00%	-43.86%	-16.33%	-42.30%	-14.50%	-43.43%	-15.17%
6	-42.13%	-13.07%	-41.36%	-10.11%	-41.39%	-9.95%	-42.76%	-13.61%	-42.81%	-13.26%
5	-40.63%	-9.04%	-40.11%	-11.04%	-39.19%	-7.63%	-40.45%	-12.24%	-38.86%	-7.64%
4	-36.05%	-5.75%	-34.42%	-6.99%	-33.78%	-3.40%	-33.76%	-6.50%	-32.87%	-4.93%
Average Imp	-41.21%	-16.30%	-40.96%	-15.97%	-40.78%	-14.53%	-41.17%	-16.10%	-40.98%	-14.58%
Dist.	1									
	Erlang	k=5	Erlang	k=4	Erlang	k=3	Erlang	k=2	Erlang	k=1
COV	Erlang 0.447214	k=5	Erlang 0.5	k=4	Erlang 0.57735	k=3	Erlang 0.707107	k=2	Erlang 1	k=1
COV M	Erlang 0.447214 E(Lc) Imp	k=5 E(Lf) Imp	Erlang 0.5 E(Lc) Imp	k=4 E(Lf) Imp	Erlang 0.57735 E(Lc) Imp	k=3 E(Lf) Imp	Erlang 0.707107 E(Lc) Imp	k=2 E(Lf) Imp	Erlang 1 E(Lc) Imp	k=1 E(Lf) Imp
COV M 10	Erlang 0.447214 E(Lc) Imp -42.74%	k=5 E(Lf) Imp -22.50%	Erlang 0.5 E(Lc) Imp -42.28%	k=4 E(Lf) Imp -21.09%	Erlang 0.57735 E(Lc) Imp -43.89%	k=3 E(Lf) Imp -22.29%	Erlang 0.707107 E(Lc) Imp -44.86%	k=2 E(Lf) Imp -18.93%	Erlang 1 E(Lc) Imp -46.26%	k=1 E(Lf) Imp -19.03%
COV M 10 9	Erlang 0.447214 E(Lc) Imp -42.74% -41.95%	k=5 E(Lf) Imp -22.50% -19.02%	Erlang 0.5 E(Lc) Imp -42.28% -43.36%	k=4 E(Lf) Imp -21.09% -19.96%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21%	k=3 E(Lf) Imp -22.29% -19.86%	Erlang 0.707107 E(Lc) Imp -44.86%	k=2 E(Lf) Imp -18.93% -15.45%	Erlang 1 E(Lc) Imp -46.26% -44.00%	k=1 E(Lf) Imp -19.03% -16.92%
COV M 10 9 8	Erlang 0.447214 E(Lc) Imp -42.74% -41.95% -42.03%	k=5 E(Lf) Imp -22.50% -19.02% -15.89%	Erlang 0.5 E(Lc) Imp -42.28% -43.36% -43.41%	k=4 E(Lf) Imp -21.09% -19.96% -16.57%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21% -44.09%	k=3 E(Lf) Imp -22.29% -19.86% -16.76%	Erlang 0.707107 E(Lc) Imp -44.86% -44.48%	k=2 E(Lf) Imp -18.93% -15.45% -13.70%	Erlang 1 E(Lc) Imp -46.26% -44.00% -41.59%	k=1 E(Lf) Imp -19.03% -16.92% -13.77%
COV M 10 9 8 7	Erlang 0.447214 E(Lc) Imp -42.74% -41.95% -42.03% -42.57%	k=5 E(Lf) Imp -22.50% -19.02% -15.89% -13.25%	Erlang 0.5 E(Lc) Imp -42.28% -43.36% -43.41% -42.87%	k=4 E(Lf) Imp -21.09% -19.96% -16.57% -12.63%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21% -44.09% -44.57%	k=3 E(Lf) Imp -22.29% -19.86% -16.76% -16.33%	Erlang 0.707107 E(Lc) Imp -44.86% -44.48% -44.65% -43.57%	k=2 E(Lf) Imp -18.93% -15.45% -13.70% -13.10%	Erlang 1 E(Lc) Imp -46.26% -44.00% -41.59% -38.45%	k=1 E(Lf) Imp -19.03% -16.92% -13.77% -13.34%
COV M 10 9 8 7 6	Erlang 0.447214 E(Lc) Imp -42.74% -41.95% -42.03% -42.57% -41.64%	k=5 E(Lf) Imp -22.50% -19.02% -15.89% -13.25% -10.08%	Erlang 0.5 E(Lc) Imp -42.28% -43.36% -43.41% -42.87% -41.46%	k=4 E(Lf) Imp -21.09% -19.96% -16.57% -12.63% -9.74%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21% -44.09% -44.57% -41.28%	k=3 E(Lf) Imp -22.29% -19.86% -16.76% -16.33% -11.17%	Erlang 0.707107 E(Lc) Imp -44.86% -44.48% -44.65% -43.57% -39.49%	k=2 E(Lf) Imp -18.93% -15.45% -13.70% -13.10% -10.30%	Erlang 1 E(Lc) Imp -46.26% -44.00% -41.59% -38.45% -33.76%	k=1 E(Lf) Imp -19.03% -16.92% -13.77% -13.34% -9.94%
COV M 10 9 8 7 6 5	Erlang 0.447214 E(Lc) Imp -42.74% -41.95% -42.03% -42.57% -41.64% -38.46%	k=5 E(Lf) Imp -22.50% -19.02% -15.89% -13.25% -10.08% -7.57%	Erlang 0.5 E(Lc) Imp -42.28% -43.36% -43.41% -42.87% -41.46% -38.19%	k=4 E(Lf) Imp -21.09% -19.96% -16.57% -12.63% -9.74% -7.76%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21% -44.09% -44.57% -41.28% -36.89%	k=3 E(Lf) Imp -22.29% -19.86% -16.76% -16.33% -11.17% -6.85%	Erlang 0.707107 E(Lc) Imp -44.86% -44.48% -44.65% -43.57% -39.49% -34.92%	k=2 E(Lf) Imp -18.93% -15.45% -13.70% -13.10% -10.30% -8.70%	Erlang 1 E(Lc) Imp -46.26% -44.00% -41.59% -38.45% -33.76% -27.67%	k=1 E(Lf) Imp -19.03% -16.92% -13.77% -13.34% -9.94% -7.22%
COV M 10 9 8 7 6 5 4	Erlang 0.447214 E(Lc) Imp -42.74% -41.95% -42.03% -42.57% -41.64% -38.46% -32.35%	k=5 E(Lf) Imp -22.50% -19.02% -15.89% -13.25% -10.08% -7.57% -5.71%	Erlang 0.5 E(Lc) Imp -42.28% -43.36% -43.41% -42.87% -41.46% -38.19% -32.23%	k=4 E(Lf) Imp -21.09% -19.96% -16.57% -12.63% -9.74% -7.76% -5.69%	Erlang 0.57735 E(Lc) Imp -43.89% -44.21% -44.09% -44.57% -41.28% -36.89% -30.62%	k=3 E(Lf) Imp -22.29% -19.86% -16.76% -16.33% -11.17% -6.85% -4.20%	Erlang 0.707107 E(Lc) Imp -44.86% -44.48% -44.65% -43.57% -39.49% -34.92% -28.01%	k=2 E(Lf) Imp -18.93% -15.45% -13.70% -13.10% -10.30% -8.70% -4 19%	Erlang 1 E(Lc) Imp -46.26% -44.00% -41.59% -38.45% -33.76% -27.67% -21.35%	k=1 E(Lf) Imp -19.03% -16.92% -13.77% -13.34% -9.94% -7.22% -7.02%

Dist.	Hyperexponential		
COV	1.732050808	1.457737974	
Μ	E(Lc) Imp	E(Lf) Imp	
10	-35.56%	-18.53%	
9	-32.76%	-19.39%	
8	-28.98%	-10.11%	
7	-25.61%	-11.59%	
6	-21.74%	-11.95%	
5	-16.28%	-7.57%	
4	-10.52%	-0.81%	
Average Imp	-24.49%	-11.42%	

#### **V** Conclusion

- For a real-time information scenario, we develop a QBD process model where the transition matrix structure is determined by the customer's choice behavior.
- Under a stability condition, we compute the stationary performance measures by utilizing the special structure of QBD process.
- For a non real-time information scenario, we treat the toll and free systems as an M/M/1/K and MMPP/M/1 models, respectively.
- Under the equilibrium condition, we develop a search algorithm for computing the performance measures.

#### Conclusion...

- With the two models developed, we discover:
  - Real-time queue length information reduces expected waiting times of both free and toll systems about 40-45% reduction for the toll system and 15-20% reduction for the free system.
  - The performance characteristics remain similar for both information scenarios.
- We also examine:
  - The effect of decision variables of the toll system toll price and buffer size => inconsistency between the social welfare goal and the private firm' profit goal.
  - The effect of the government subsidy on the performance of the two-tier service systems.

The End
Q & A