Large scale queueing systems in the Quality/Efficiency (Halfin-Whitt) driven regime, and applications

David Gamarnik MIT

33rd CONFERENCE ON THE MATHEMATICS OF OPERATIONS RESEARCH

January 2008

Joint work with **P. Momcilovic**, U of Michigan,

Talk Outline

- GI/GI/N queueing model
- Applications to call/contact centers
- Challenge: non-Markovian systems
- **Results I**: Markov chain characterization of the queue length process
- **Result II**: Interchange of Heavy Traffic-Steady State limits. Tight decay rate.
- Discussion of methods: Lyapunov functions
- Further challenges

Model: G/G/N queueing system



Model: G/G/N queueing system

Kiefer & Wolfowitz [56] – Steady state regime exists (stability) iff

$$\rho_N \triangleq \frac{\lambda}{N\mu} < 1.$$



... but no explicit formulas for Q_{∞}, W_{∞} except Erlang M/M/N (Erlang-C) formulas



Also

$$\mathbb{P}(\frac{Q_{\infty}^{N}}{\sqrt{N}} > t | Q_{\infty}^{N} > 0) \to e^{-\beta t}$$

$$\mathbb{P}(\frac{I_{\infty}^{N}}{\sqrt{N}} > t | Q_{\infty}^{N} = 0) \to \frac{\Phi(\beta - t)}{\Phi(\beta)}$$

Furthermore:

• Extends to non-Poisson arrival processes: G/M/N

$$\mathbb{P}(\frac{Q_{\infty}^{N}}{\sqrt{N}} > t | Q_{\infty}^{N} > 0) \to e^{-\theta t}, \quad \theta = \frac{2\beta}{1 + c_{a}^{2}}$$

$$\mathbb{P}(\frac{I_{\infty}^{N}}{\sqrt{N}} > t | Q_{\infty}^{N} = 0) \to \frac{\Phi(\beta - t)}{\Phi(\beta)}$$

- Extends to phase-type service times: Puhalskii & Reiman [2000]
- Diffusion approximations:

$$\left(\frac{Q_t^N}{\sqrt{N}}, \frac{I_t^N}{\sqrt{N}}\right) \Rightarrow (\hat{Q}_t, \hat{I}_t)$$

will return to this

Motivation: Call Centers



Tradeoff between utilization $\,
ho\,\,$ (cost of staffing N) and performance $\,\mathbb{P}({\sf Wait}>x)\,$

Motivation: Call Centers



Surveys:

M. Armony: <u>http://www.stern.nyu.edu/om/faculty/armony/research/CallCenterSurvey.pdf</u> A. Mandelbaum: <u>http://iew3.technion.ac.il/serveng/References/references.html</u> W. Whitt: <u>http://www.ieor.columbia.edu/~ww2040/CallCenterF04/call.html</u>

Papers on queueing models of call centers by

Atar, Armony, Baron, Bassamboo, Brown, Dai, Green, Gans, Garnett, Gurvich, Harrison, Jelenkovic, Jennings, Halfin, Kolesar, Kumar, Maglaras, Mandelbaum, Massey, Momcilovic, Randhawa, Reed, Reiman, Sakov, Shen, Shimkin, Tezcan, Whitt, Zeevi, Zeltin, Zhao, Zohar

Brown, Gans, Mandelbaum, Sakov, Shen, Zeltyn, Zhao [2002]

Statistical Analysis of a Telephone Call Center: A Queueing-Science Perspective.

One of the conclusions: service times have log-normal distribution.

General service time distributions

- Deterministic service times: Jelenkovic, Mandelbaum & Momcilovic [2004]
- G/G/N virtual waiting times process limits, discrete service times: Mandelbaum & Momcilovic [2005]

$$\frac{W_t}{\sqrt{N}} \Rightarrow \hat{W}_t$$

G/G/N process limits:
 Reed [2006]

Principal questions for this work:

- What about G/G/N in steady state?
- Does $Q_{\infty}^N \approx \sqrt{N}$ hold in general?

Assumption: service times have a *discrete* probability distribution.

$$\mathbb{P}(S=sk)=p_k, \quad k=1,2,\ldots,K$$

Summary of results:

- Markov chain characterization of the limiting process $\lim_{N\to\infty} \frac{Q^N(t)}{\sqrt{N}} = Q(t)$
- Tight decay rate for the limiting queue length in steady-state

Consider the following Markov chain on $(Q_t, L_t) \in \mathbb{R}^{K+1}$

$$L_{t+1} = \mathcal{T}[L_t] + X + \left((Q_t + \sqrt{\mu}\sigma X) \wedge (\beta - \sum_{j \le 2} L_{j,t}) p, \right.$$
$$Q_{t+1} = \left(\sqrt{\mu}\sigma X + Q_t + \sum_{j \le 2} L_{j,t} - \beta \right)^+,$$

where

$$X \stackrel{d}{=} N(0,1)$$

$$X \stackrel{d}{=} N(0,\Sigma), \qquad \Sigma_{ij} = \begin{cases} (1-p_i)p_i, & 1 \le i = j \le K, \\ -p_i p_j, & 1 \le i \ne j \le K; \end{cases}$$

$$\mathcal{T}(l_1, \dots, l_K) = (l_2, \dots, l_{K-1}, 0)$$

 $\boldsymbol{L}_{t}^{N} = (L_{1,t}^{N}, \dots, L_{K,t}^{N}), \$ - centered work in progress

Theorem I. [Transient] Suppose

$$\mathbb{P}(S = k) = p_k, \quad k = 1, \dots, K, \quad \textbf{Halfin-Whitt regime}$$

$$\lambda_N = \mu N - \beta \mu \sqrt{N} \quad \textbf{Halfin-Whitt regime}$$

$$\left(\frac{Q_0^N}{\sqrt{N}}, \frac{L_0^N}{\sqrt{N}}\right) \Rightarrow (Q_0, L_0). \quad \textbf{Initialization}$$
Then $\left(Q_t^N, L_t^N\right) \Rightarrow (Q_0, L_0) = 1.2$

$$\left(\frac{Q_t}{\sqrt{N}}, \frac{L_t}{\sqrt{N}}\right) \Rightarrow (Q_t, L_t), \quad t = 1, 2, \dots$$

Theorem II. [Steady-state]

• (Q_t, L_t) has a *unique* stationary distribution (Q_∞, L_∞) .





Theorem II. [Steady-state] Moreover (exact decay rate)

$$\mathbb{P}(Q_{\infty} > t) \approx \exp(-\frac{2\beta t}{c_8^2 + c_s^2}),$$
 for large t

Same decay rate as the conventional heavy-traffic model of G/G/N!

• Service time
$$S = \begin{cases} 1, w.p. & 0.5; \\ 2, w.p. & 0.5. \end{cases}$$

• Arrival rate
$$\lambda_N = (2/3)N - \beta \sqrt{N}$$

 $Z_{1,t}$ – number of calls with remaining service time $< \mathbf{1}$

 $Z_{2,t}$ – number of calls with remaining service time in [1,2)







Scaling N does not reveal anything. Look at Gaussian scaling:

$$\hat{A}_{t-1,t} = \frac{A_{t-1,t} - (2/3)N}{\sqrt{N}}$$

$$L_{1,t} = \frac{Z_{1,t} - (2/3)N}{\sqrt{N}}$$

$$L_{2,t} = \frac{Z_{1,t} - (1/3)N}{\sqrt{N}}$$

• Case I $Z_{1,t} > Q_t + A_{t,t+1}$

Then

$$Z_{2,t+1} \approx \frac{1}{2}(Q_t + A_{t,t+1})$$
$$Z_{1,t+1} \approx Z_{2,t} + \frac{1}{2}(Q_t + A_{t,t+1})$$
$$Q_{1,t+1} = 0.$$

Subtract the means

$$L_{2,t+1} \approx \frac{1}{2}(Q_t + A_{t,t+1}) - \frac{1}{3}N$$
$$L_{1,t+1} \approx L_{2,t} + \frac{1}{2}(Q_t + A_{t,t+1}) - \frac{2}{3}N$$
$$Q_{1,t+1} = 0.$$

• Case II similar ...

• Case II similar ...

Putting things together, the Markov chain dynamics is obtained

$$\boldsymbol{L}_{t+1} = \mathcal{T}[\boldsymbol{L}_t] + \boldsymbol{X} + \left((Q_t + \sqrt{\mu}\sigma X) \wedge (\beta + L_{1,t} - \|\boldsymbol{L}_t\|) \right) \boldsymbol{p},$$
$$Q_{t+1} = \left(\sqrt{\mu}\sigma X + Q_t + \|\boldsymbol{L}_t\| - L_{1,t} - \beta \right)^+,$$

Dynamical system is hard to analyze even without stochastic fluctuations.

$$L_{t+1} = \mathcal{T}[L_t] + \left(Q_t \wedge (\beta + L_{1,t} - \|L_t\|)\right)p,$$

$$Q_{t+1} = \left(Q_t + \|L_t\| - L_{1,t} - \beta\right)^+,$$

Steady-State

Proposition. The process $Y_t = Q_t + \sum_j L_{t,j}, Q_t = (Y_t - \beta)^+$ satisfies

$$Y_t = \sum_k p_k (Y_{t-k} - \beta)^+ + \sum_k \alpha_k Z_{t-k},$$

 Z_t - stationary Gaussian process

When Y_t is "large", the drift is pprox -eta

Steady-State

Theorem.
$$\Phi(y,z) = \exp(\theta \sum_k p_k^* y_k + \theta \sum_k \alpha_k^* z_k), \qquad (v_k^* = \sum_{j \ge k} v_j).$$

is a (geometric) Lyapunov function.

$$\mathbb{E}\Big[\Phi(Y_{t+1}, Z_{t+1})\Big|Y_t = y, Z_t = z\Big] < (1 - \gamma)\Phi(y, z)$$

Lyapunov function argument is used to show that

• (L_t, Q_t) has a stationary distribution (Q_∞, L_∞) .

•
$$\left(\frac{Q_{\infty}^{N}}{\sqrt{N}}, \frac{\boldsymbol{L}_{\infty}^{N}}{\sqrt{N}}\right) \Rightarrow (Q_{\infty}, \boldsymbol{L}_{\infty}).$$

$$\bullet \quad \limsup_{N \to \infty} \mathbb{E} \Big[\exp \Big(\frac{\theta Q_\infty^N}{\sqrt{N}} \Big) \Big] < \infty, \qquad \text{ for some } \ \theta > 0.$$

Summary, Challenges and ongoing work:

- G/G/N in Halfin-Whitt regime with discrete service time distribution. Process level and steady-state results.
- Tight decay rate for the queue length in steady-state.
- **Challenge:** non-discrete service times.
- In progress: relaxation in M/M/N in Halfin-Whitt (joint work with D. Goldberg)