# The Erlang-R Queue: A Model Supporting Personnel Staffing in Emergency Departments

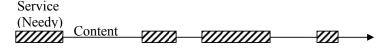
Galit Yom-Toy Avishai Mandelbaum

Industrial Engineering and Management Technion

ORSIS, May 10-11, 2009

**Standard assumption in service models:** service time is continuous.

But we find systems in which: service is dis-continuous and customers re-enter service again and again.



#### **Examples:**

- Machine-Repairman (closed queueing network).
- Medical Unit model (semi-open queueing network).
- Emergency Department (open queueing network).

What is the appropriate staffing procedure? What is the significance of these cycles? Can one still use simple Erlang-C models for staffing?

#### **Related Work**



Massey W.A., Whitt W. Networks of Infinite-Server Queues with Nonstationary Poisson Input. 1993.

Green L., Kolesar P.J., Soares J.

Improving the SIPP Approach for Staffing Service Systems that have Cyclic Demands. 2001.

Jennings O.B., Mandelbaum A., Massey W.A., Whitt W. Server Staffing to Meet Time-Varying Demand. 1996.

Feldman Z., Mandelbaum A., Massey W.A., Whitt W. Staffing of Time-Varying Queues to Achieve Time-Stable Performance. 2007.

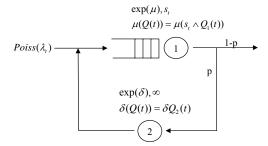
#### Research Outline

- Develop a queueing model that incorporates service and customers' Re-entrance: Erlang-R.
- Develop fluid and diffusion approximations.
- Develop staffing algorithm that attain pre-specified service levels.
- Validate our approach via simulation, based on hospital data.
- Understand when Erlang-R is needed.

#### **Model Definition**

#### The (Time-Varying) Erlang-R Queue:

- λ<sub>t</sub> Arrival rate of a time-varying Poisson arrival process.
- $\mu$  Exponential service rate.
- $\delta$  Delay rate (1/ $\delta$  is the delay time between services).
- p Probability of return to service.
- s<sub>t</sub> Number of servers at time t.
- $Q_i(t)$  Number of customers in node *i* at time t, i = 1, 2.



## Fluid and Diffusion Approximations for Erlang-R

We scale by  $\eta: \lambda_t \to \eta \lambda_t$  and  $s_t \to \eta s_t$ .

#### Theorem: Fluid (FSSLN) and Diffusion (FCLT) Approximations

As  $\eta \to \infty$ ,

$$\frac{Q^{\eta}(t)}{\eta} \rightarrow Q^{(0)}(t), \ \ \text{and} \ \ \sqrt{\eta} \left[ \frac{Q^{\eta}(t)}{\eta} - Q^{(0)}(t) \right] \stackrel{d}{=} Q^{(1)}(t),$$

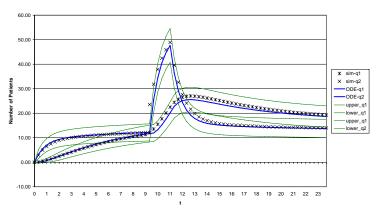
where  $Q^{(0)}(t)$  is the solution of the following ODE:

$$egin{aligned} Q_1^{(0)}(t) &= Q_1^{(0)}(0) + \int_0^t \left(\lambda_ au - \mu_ au \left(Q_1^{(0)}( au) \wedge oldsymbol{s}_ au
ight) + \delta_ au Q_2^{(0)}( au)
ight) d au \ Q_2^{(0)}(t) &= Q_2^{(0)}(0) + \int_0^t \left(p\mu_ au \left(Q_1^{(0)}( au) \wedge oldsymbol{s}_ au
ight) - \delta_ au Q_2^{(0)}( au)
ight) d au, \end{aligned}$$

and  $Q^{(1)}(t)$  is the solution of a SDE (Stochastic Differential Equation).

## Example of Fluid and Diffusion Approximations for Erlang-R

Delta = 0.2; Mu = 1; p = 0.25; s = 50; Lambda=10 (t<9 or t>11), Lambda=50 (9<t<11)



## Staffing: Determine $s_t$ , t > 0

Based on the QED-staffing formula:

$$s = R + \beta \sqrt{R}$$

- Two approaches:
  - PSA / SIPP (lag-SIPP) divide the time-horizon to planning intervals, calculate average arrival rate and steady-state offered-load for each interval, then staff according to steady-state recommendation (i.e.,  $R(t) \approx \bar{\lambda}(t)E[S]$ ).
  - MOL/IS assuming no constraints on number of servers, calculate time-varying offered load. For example in Erlang-C:  $R(t) = E[\int_{t-S}^{t} \lambda(u) du = E[\lambda(t-S)]E[S].$

Staff according to the square-root formula where R(t)replaces R:  $s(t) = R(t) + \beta \sqrt{R(t)}$ .

Offered-Load in Erlang-R queue = The number of busy servers (or the number of customers) in a corresponding  $(M_t/M/\infty)^2$  network. R(t) is denoted by the following two expressions:

$$R_1(t) = E\left[\int_{t-S_1}^t \lambda_u + \delta Q_2^{(0),\infty}(u) du\right]$$

$$R_2(t) = E\left[\int_{t-S_2}^t p\mu Q_1^{(0),\infty}(u) du\right]$$

where  $Q^{(0),\infty}(t)$  is the solution of the following Fluid ODE:

$$\frac{d}{dt}Q_1^{(0),\infty}(t) = \lambda_t + \delta Q_2^{(0),\infty}(t) - \mu Q_1^{(0),\infty}(t), 
\frac{d}{dt}Q_2^{(0),\infty}(t) = p\mu Q_1^{(0),\infty}(t) - \delta Q_2^{(0),\infty}(t).$$

## Staffing: Determine $s_t$ , t > 0

#### MOL Algorithm for Erlang-R:

- Solve differential equations for  $Q^{\infty}(t)$ .
- Calculate time-varying offered load R(t).
- Staff according to square-root formula:  $s(t) = R(t) + \beta \sqrt{R(t)}$ , where  $\beta$  is chosen according to the steady-state Halfin-Whitt formula.

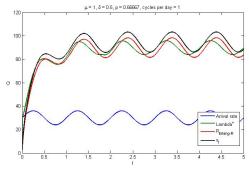
## Case Study: Sinusoidal Arrival Rate

Examine the following periodic arrival rate:

$$\lambda_t = \bar{\lambda} + \bar{\lambda} \alpha \sin(2\pi t/\psi)$$

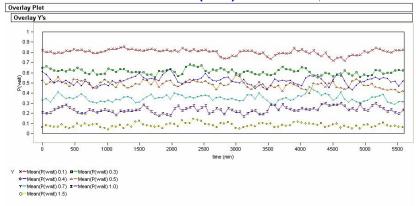
where  $\bar{\lambda}$  is the average arrival rate,  $\alpha$  is the relative amplitude (0 <  $\alpha$  < 1), and  $\psi$  is the period length.

#### Arrival rate, Offered load, and Staffing



## Case Study: Sinusoidal Arrival Rate

#### Simulation results of P(wait) for various $\beta$ values



The performance measure is stable!

## Are Cycles Significant? When?

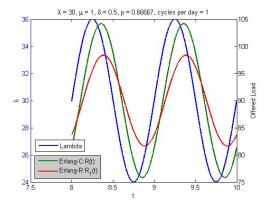
#### Compare to concatenated services:

Define a = average number of re-entrances to service (a = 1/(1 - p)).

• Multi-service Erlang-C: Sum up all service times of each customer (RCCP)

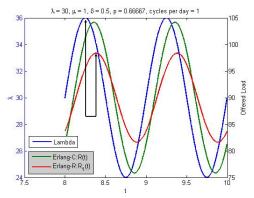
• When  $\lambda(t) \equiv \lambda$ , all have the same steady state! Hence, all measurements (such as P(W > 0)) are identical, though Erlang-R provides information also about the *content* state.

Erlang-C under- or over-estimates the offered load.



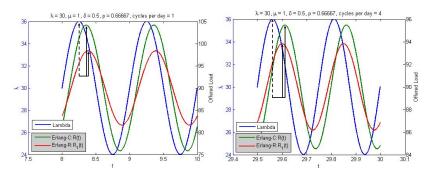
## When Classic Model Fails? Time-Varying Systems

Well known that there is time-lag between peak arrival rate and offered load.



## When Classic Model Fails? Time-Varying Systems

Erlang-C under- or over-estimates this time-lag depending on the period's length.



## Why Erlang-C Does Not Fit Re-entrant Systems?

#### Compare R(t) of Erlang-C and Erlang-R:

*Multi-service* Erlang-C's offered load:  $S_1 \leftrightarrow aS_1$ :

$$R(t) = E[\lambda(t - aS_1)] E[aS_1]$$

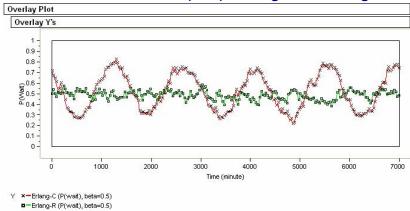
Erlang-R's offered load:

$$R_1(t) = ... = \sum_i p^i E[S_1] E[\lambda(t - (i+1)S_1 - iS_2)]$$

## Nevertheless, Can We Use Erlang-C?

From the previous case study - sinusoidal arrival rate.

#### Simulation results of P(wait): Erlang-C vs. Erlang-R



Using Erlang-C's R(t), does not stabilize systems' performance.

## Thank You