# Queues with simultaneous loss on breakdowns

Dave Thornley [*]

### Abstract

We take a queue with breakdowns and repairs of processors in which
the queue length does not change on the corresponding Markov modula-
tion transitions as given by Mitrani and Chakka [1], and introduce simul-
taneous losses on breakdown (and re-sampling on repair) using a range of
methods, most interestingly including enriching the modulation structure
to cause simultaneous loss or re-sampling, which gives distinct behaviour.
We solve for the steady state of the queue using spectral expansion for
convenient access to a range of performance measures.

## 1   Introduction

Multiprocessor servers with breakdowns and repairs can be modelled using a
Markov modulated queue in which the modulation state tracks the number of
active processors. Mitrani and Chakka [1] provide an efficient solution for such
a queue in which the jobs associated with a given processor are re-sampled when
that processor breaks down, such that the queue length is not altered.

We describe a new system derived from this, but which loses the job in service
at the instant a processor breaks down. Modeling this requires spreading the
modulation structure across queue lengths.

To examine the importance of simultaneity of loss, we contrast the key result
with those derived from using independent streams of negative customers to
approximate the loss rates. We find that the queue length distribution due to
simultaneous losses cannot be adequately matched, particularly at the full and
empty queues.

We can extend our generalised modulation concept to model re-sampling of
a job when a processor recovers. This may correspond to a situation where a
deactivated server has local persistent storage which takes jobs from the queue.

The formulation is further extensible to represent geometrically batched
losses and re-samplings, which can then be solved using our latest techniques
[2].

In this paper, we provide the Kolmogorov balance equations governing the
simultaneous loss/resample behaviour, describe an implementation method, and
present some results for a simple example which emphasis the significance of
simultaneity.

---

[*]Department of Computing Imperial College of Science, Technology and Medicine Huxley
Building 180 Queen's Gate London SW7 2BZ England `djt@doc.ic.ac.uk`

# 2 Breakdowns and repairs

In a uniform multi-processor server, we have $N$ processors, each with the same processing rate[1] $\mu$. These break down independently as a Poisson point process with rate $b$. Inactive processors are similarly repaired independently at rate $r$. The distinct combinations of activity/inactivity are represented as distinct phases in a Markov modulation process. The joint state probabilities of queue length and modulation state form a (semi-)finite lattice strip. Figure 1 shows small excerpts of such a strip for a range of queueing behaviours associated with breakdowns and repairs. The horizontal dimension gives the number of active processors in each column, and the vertical dimension is the queue length, increasing upwards.
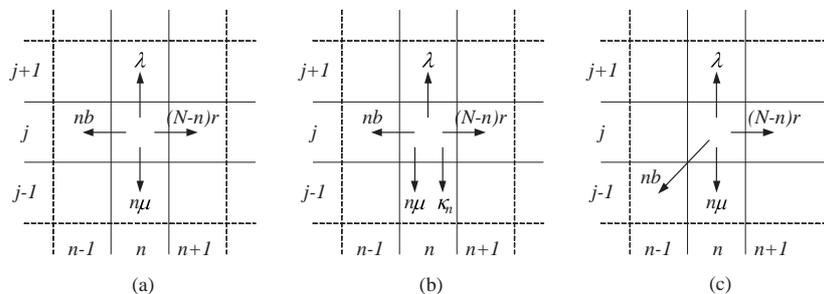


Figure 1: State transitions in a queue with breakdowns and repairs. (a) is as in [1], (b) has losses due to remove-from-head negative customers, and (c) models simultaneous losses.

Figure 1(a) shows the transitions away from state $(n, j)$ in the semi-finite lattice describing the transitions in the Markov modulated queue with single breakdowns and repairs.

$$\mathbf{v}_{j-1}[\Lambda] + \mathbf{v}_j[Q - M - \Lambda] + v_{j+1}[M] = \mathbf{0}$$

Where the modulator's instantaneous transition matrix $Q$ is of the same form as [1] as given below[2]. Figure 1 a) and b) show the loss of the job in process using respectively independent negative customers and simultaneous transition.

The work in [1] allows for all active processors to break down simultaneously at a given rate, which we will call $b_0$, and simultaneous repair of all inactive processors, which we will call $r_0$. When queue length does not change on breakdown or repair, this is entirely encapsulated within the transition structure of

---

[1] Heterogeneous servers can easily be provided by taking a Cartesian product of the ensemble of modulation processes for each server's activity
[2] In our examples, $b_0$ and $r_0$ are set to zero for clarity.

the modulator.

$$Q = \begin{pmatrix} -\Sigma\ldots & Nr & & & r_0 \\ b_0+b & -\Sigma\ldots & (N-1)r & & r_0 \\ b_0 & 2b & -\Sigma\ldots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & r+r_0 \\ b_0 & & & Nb & -\Sigma\ldots \end{pmatrix}$$

Associated with this modulation structure, the processing rates $\mu_i = (i-1)\hat{\mu}$ in $M$ give the modulated processing rates. Matrix $\Lambda = \lambda I$ the arrival rate in all modulation states.

Losses on breakdown and/or resampling on repair require the addition of a vertical component to the originally horizontal breakdown and repair transitions. The resulting transitions are diagonal, except for breakdowns at the empty queue which are necessarily constrained to being horizontal, and repair transitions at a full queue, which cannot increase the queue length[3].

It is also interesting to note that, when diagonal transitions are used, a particular quality of traditional Markov modulated queues is also lost: horizontal transition rates in the lattice are no longer independent of queue length, so the sum of probabilities in a given queue length are not equal to the equilibrium state probability of the modulator. Indeed, the modulator is not complete for the majority of the queue.

# 3 Breakdown losses modelled by negative customers

To illustrate the effect of the loss being simultaneous with the server breakdown, we briefly examine the results of attempting to approximate the behaviour of the queue using *independent* negative customers [3] in a Poisson arrival stream. This is contrasted with the behaviour of the queue with simultaneous loss of the job in process, which may be considered to be caused by a negative customer removing from the head of the queue, which is *triggered* by the processor breakdown.

We show results of using the same arrival rate as the processor breakdown rate, and a rate chosen to result in a mean queue length equal to the simultaneous loss distribution. These are presented in graph form for visual comparison.

# 4 Simultaneous loss on breakdowns

To correctly model the loss of a job in progress when the processor dies, we use more appropriate transitions in the lattice. Let function $D(A)$ return a diagonal matrix of the row sums of $A$.

The Kolmogorov balance equations for the queue with simultaneous losses

---

[3] Which job is lost depends on the queueing formalism. In a FCFS queue, it may be more natural to lose the job at the end of the queue, since the re-sample was submitted to the system earlier.

is as follows:

$$\mathbf{v}_j[R - D(R) + B - D(B) - M_j - \Lambda] + v_{j+1}[M_{j+1} + B] = \mathbf{0}, \text{ for } j = 0$$
$$\mathbf{v}_{j-1}[\Lambda] + \mathbf{v}_j[R - D(R) - D(B) - M_j - \Lambda]$$
$$+ v_{j+1}[M_{j+1} + B] = \mathbf{0}, \text{ for } 0 < j < L$$
$$\mathbf{v}_{j-1}[\Lambda] + \mathbf{v}_j[R - D(R) - D(B) - M_j] = \mathbf{0}, \text{ for } j = L < \infty$$

Matrix $M_j$ gives the processing rates $\mu_{m,j} = \min(j\mu, \mu_m)$ at queue length $j$, where $\mu_m = (m-1)\mu$ for a queue with N homogeneous when we enumerate modulation states starting at 1. processors. The matrix $B$ provides the transitions in the lattice due to breakdowns (unspecified elements are zero):

$$B = \begin{pmatrix} 0 & & & & & 0 \\ b & \ddots & & & & \\ & \ddots & & 0 & & \\ & & & (N-1)b & 0 & \\ 0 & & & & Nb & 0 \end{pmatrix}$$

The $R$ matrix describes the repair behaviour, and is detailed below.

# 5 Simultaneous re-sampling on processor recovery

We now consider allowing the job which was taken out of the system when the processor broke down to be re-sampled when the processor recovers. This corresponds to the processor having local storage. When it is reactivated, the job becomes another potential departure. This has implications for the specific queueing paradigm in terms of significantly disrupting the implicit service order, especially if repairs are slow.

Matrix $R$ provides the transitions due to repairs in the Kolmogorov balance
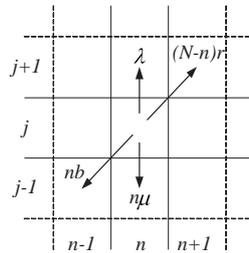


Figure 2: State transitions in a queue with breakdowns and repairs with simultaneous loss of the job in service and resampling on repair

equations for the queue.

$$R = \begin{pmatrix} 0 & Nr & & & 0 \\ & 0 & (N-1)r & & \\ & & 0 & \ddots & \\ & & & \ddots & r \\ 0 & & & & 0 \end{pmatrix}$$

The Kolmogorov balance equations for this queue are as follows:

$$\mathbf{v}_j[R - D(R) + B - D(B) - M_j - \Lambda] + v_{j+1}[M_{j+1} + B] = \mathbf{0}, \text{ for } j = 0$$
$$\mathbf{v}_{j-1}[\Lambda + R] + \mathbf{v}_j[-D(R) - D(B) - M_j - \Lambda]$$
$$+ v_{j+1}[M_{j+1} + B] = \mathbf{0}, \text{ for } 0 < j < L$$
$$\mathbf{v}_{j-1}[\Lambda + R] + \mathbf{v}_j[-D(R) - D(B) - M_j] = \mathbf{0}, \text{ for } j = L < \infty$$

In the examples we present later with one processor, we have the following $B$ and $R$ matrices:

$$B = \begin{pmatrix} 0 & 0 \\ b & 0 \end{pmatrix}, R = \begin{pmatrix} 0 & r \\ 0 & 0 \end{pmatrix}$$

# 6   Geometric batches

The Kolmogorov balance equations for geometrically batched losses as soluble using techniques detailed in [2] on breakdown and re-sampling on repair are as follows:

$$\mathbf{v}_j[-D(R) + B - D(B) - M_j - \Lambda]$$
$$+ v_{j+1}[M_{j+1}] + \sum_{i=j+1}^{L} \mathbf{v}_i B \beta^{i-j-1} = \mathbf{0}, \text{ for } j = 0$$

$$\mathbf{v}_{j-1}[\Lambda] + \sum_{i=0}^{j-1} \mathbf{v}_i R(1-\rho)\rho^{j-i-1} + \mathbf{v}_j[-D(R) - D(B) - M_j - \Lambda]$$

$$+ v_{j+1}[M_{j+1}] + \sum_{i=j+1}^{L} \mathbf{v}_i B(1-\beta)\beta^{i-j-1} = \mathbf{0}, \text{ for } 0 < j < L$$

$$\mathbf{v}_{j-1}[\Lambda] + \sum_{i=0}^{j-1} \mathbf{v}_i R\rho^{j-i-1} + \mathbf{v}_j[R - D(R) - D(B) - M_j] = \mathbf{0}, \text{ for } j = L < \infty$$

Where $\beta$ is the batch size distribution parameter of the breakdown loss batches, and $\rho$ is the equivalent parameter of re-sample batches. These parameters pertain to the probability function of batch size $s$ with parameter $\theta$ of the form $P(s) = (1-\theta)\theta^{s-1}$. Batch transitions which reach the empty queue or the full queue are truncated, so we see the sum of all batches equal to or larger than the transition size.

# 7  Solution by spectral expansion

The elegant[4], efficient[5] and accurate[6] spectral expansion solution method espoused in [1] is applicable to finite and infinite queues. One drawback is that calculation of very large numbers of eigenvalues is highly compute intensive and potentially unstable[7]. The most recent efficient methods [4] do not calculate the eigensystem directly, and we find the direct access to the eigensystem made available by spectral expansion convenient for calculation of performance measures.

The spectral expansion component of our solution mechanism follows a scheme derived from one given in [1]. The significant additional functionality relating to localization of geometrically batched processes (introduced in [5], and later satisfactorily implemented using our latest techniques [2]) is applicable to the case of loss or re-sampling of batches.

The region(s) of a queue described by a homogeneous matrix geometric series is(are) most succinctly identified by the probability flux *out* of the component levels (rows of equal queue length in the lattice) being identical. For example, in the case of the single processor queue with simultaneous losses on breakdown, the region of the joint state modulation queue length probability lattice lies between queue lengths 1 and $L-1$, where $L$ is the maximum queue length. When we do not consider resampling, we have:

$$\mathbf{v}_{j-1} \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} + \mathbf{v}_j \begin{pmatrix} -r-\lambda & r \\ 0 & -b-\lambda-\mu \end{pmatrix} + v_{j+1} \begin{pmatrix} 0 & 0 \\ b & \mu \end{pmatrix} = \mathbf{0}$$

The outward probability flux at the empty queue and at a full queue cannot be described as part of an eigensystem, as they have different flux components. At the empty queue, we have no downward transitions due to processing, and the loss component of a breakdown transition is absent. At a full queue, the arrival transitions are absent, and the resampling component of a repair transition is absent.

The solution to the linear homogeneous matrix equation given above is the sum of geometric series provided by the eigenvalues $\xi_i$ of the characteristic equation $Q(\xi) = \sum_{k=0}^{n^u+n^d} \mathbf{v}_{j-n^u} Q_i$ found by setting $\det |Q(\xi)| = 0$, projected onto the corresponding left eigenvectors $\boldsymbol{\psi}_i$ from $\boldsymbol{\psi}_i Q(\xi_i) = 0$.

Each eigenvalue/eigenvector pair defines a basis function component, and by summing the ensemble with each component scaled appropriately, we can satisfy the boundary conditions defined by the balance equations of the levels neighbouring the repeating region, and the normalization constraint.

At queue lengths $j$ falling under the repeating region balance equation, we have:

$$\mathbf{v}_j = \sum_{i=1}^{\epsilon^n} \alpha_i \xi_i^j \boldsymbol{\psi}_i \tag{1}$$

---

[4] The direct use of an eigensystem to represent the system is notionally satisfying, and meshes well with established techniques for formulating distributions of performance measures such as sojourn times [7].

[5] Spectral expansion is often compared to matrix geometric methods, and it is not clear which is the most efficient - each has advantages in different specific situations

[6] The matrix geometric method becomes less accurate at high loads

[7] We have only dealt with systems with approximately 150 eigenvalues in our Mathematica®[6] test-rig [2] so far.

The coefficients $\alpha_i$ are free variables to be constrained by the boundary conditions imposed by the rest of the queue (the processor filling region, and in the case of a finite queue, the full queue region.)

The boundary conditions are imposed by the inclusion of balance equations which include both explicit **v** vectors outside the eigensystem region, and vectors defined by the eigensystem.

When the queue is infinite, any eigenvalues of magnitude greater than or equal to 1 (*i.e.* lying on or outside the unit disk in the Argand plane) must take zero coefficients, as their infinite sum does not converge, and hence cannot be normalized.

# 8 Implementation details

For our purposes, to treat the system within our existing framework [2], the probability flux in the lattice for simultaneous breakdowns can be viewed as negative customers which also disable a processor, removing a job from the head of the queue. For normal negative customers, the transition is solely in the queue length dimension (vertical). We introduce a component in what we can safely call the modulation dimension (horizontal) to represent the change in number of active processors.

Our solver takes the parameters of the queue's behaviour in matrix form, in order that we may associate the rates and batch size distribution parameters with the correct state occupation probabilities within vectors of the probabilities at a given queue length.

To achieve a diagonal transition, we use $B$ (given earlier) as a negative customer rate matrix. This contrasts with the negative customer rate matrix used in the approximation, which is diagonal, thus representing purely vertical transitions in the queue lattice strip, with the horizontal transitions due to breakdown being encapsulated in the modulation strucure.

There remains an element of modulation, and we take care to provide the system with an appropriate $Q$ matrix. Note that it is no longer a generator matrix. For example, consider a system with one unreliable processor. We consider the balance equations given in section 4, and note that the equivalent of $Q$ in more traditional Markov modulated queues is as follows:

$$Q = R - D(R) + B - D(B), \text{ for } j = 0$$
$$R - D(R) - D(B), \text{ for } 0 < j \leq L$$

In a two state example, we have:

$$B = \left( \begin{smallmatrix} 0 & 0 \\ b & 0 \end{smallmatrix} \right), R = \left( \begin{smallmatrix} 0 & r \\ 0 & 0 \end{smallmatrix} \right)$$
$$Q = \left( \begin{smallmatrix} -r & r \\ b & -b \end{smallmatrix} \right), \text{ for } j = 0$$
$$= \left( \begin{smallmatrix} -r & r \\ 0 & -b \end{smallmatrix} \right), \text{ for } 0 < j \leq L$$

This $Q$ matrix at levels $0 < j \leq L$ does not provide a transition back to modulation state 1 from 2. This is instead provided by the breakdown rates given in $B$. At the empty queue, no jobs are lost on breakdown of a server, so this transition becomes horizontal, and is incorporated into the $Q$ matrix.

# 9 Results

We present a number of simple examples of a single unreliable processor queue. We use the terms *active phase* and *inactive phase* to refer to the modulation states representing the corresponding state of the processor. We use a short finite queue as it highlights the effect on the sharp peak in queue length probability at the full queue due to the zero processing rate in the inactive phase. We see that this peak is de-emphasised when the independent negative customer approximation is used.

First, we calculate the queue length distribution of the queue when simultaneous losses are enabled. This is the reference behaviour, and the behaviour resulting from use of independent negative customers is plotted in comparison with this.

We then use independent negative customers to approximate this behaviour. We first use a negative customer rate equal to the breakdown rate in the active phase, and show the distribution for this. We then multiply the negative customer rate by an adjustment factor to match the mean queue length of the reference behaviour. We show these distributions superposed for comparison.

Negative customers at the breakdown rate are then used in the inactive phase, and similar comparisons made, followed by an examination of the results of negative customers present in both phases at equal rates.

## 9.1 Observations

Figure 3 shows the joint queue length activity state characteristics of the single unreliable processor queue with losses due to a) independent Poisson arrivals of negative customers (PAN) in the active state, b) PAN in the inactive state, c) PAN at equal rates in both activity states and d) with simultaneous loss on breakdown (SLB). Note that the PAN characteristics all show a dip in the joint probability of very short queue lengths in the inactive state, which is absent in the SLB characteristic.

Figure 4 shows the queue length probabilities resulting from the independent negative customer approximation methods plotted against the reference characteristic with simultaneous losses. The solid line is the SLB characteristic, the dotted line shows the result of using a negative customer rate equal to the breakdown rate, and the dashed line shows the characteristic after the negative customer rate has been adjusted in order to give the same mean queue length.

The factor by which the negative customer arrival rate was multiplied to achieve the mean queue length matches shown (which are to within 1%) are a) 17 for the active phase, b) 7.7 for the inactive phase and c) 5 for both phases (each being $b$ before adjustment).

Figure 4a) shows that the utilisation is significantly under-estimated when using negative arrivals in the active phase, and that the full queue probability is over-estimated. Conversely, if negative customers are used in the inactive state, the match is improved. Figure 4b) shows that when negative customers are used only in the inactive phase, the utilisation can be closely matched, but the full queue probability is underestimated. In figure 4c) when negative customers are used in both phases, a balance appears to be struck: the utilisation is still well matched, and the full queue probability match is improved.
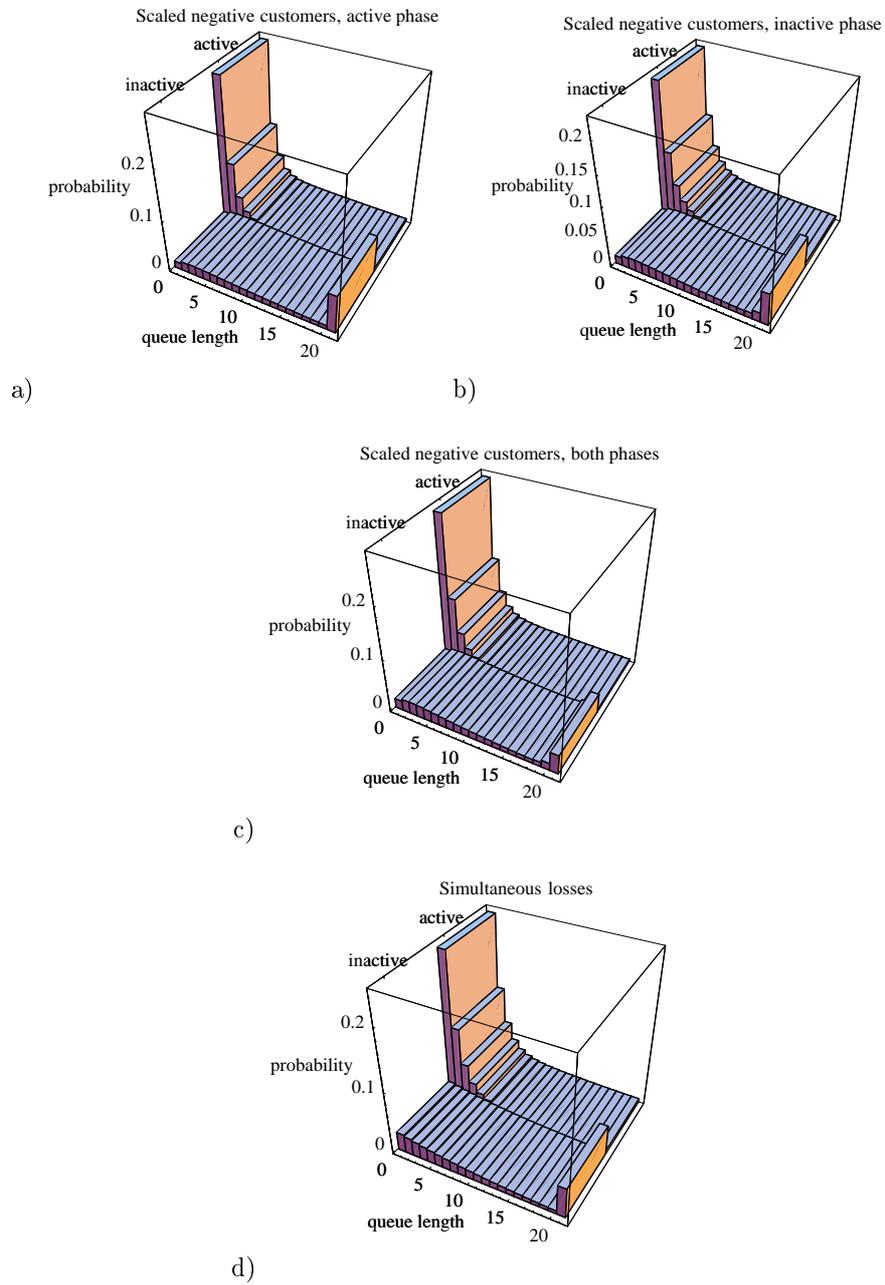
Figure 3: Joint modulation state queue length probability plots for the simultaneous loss case, and the three example approximations.
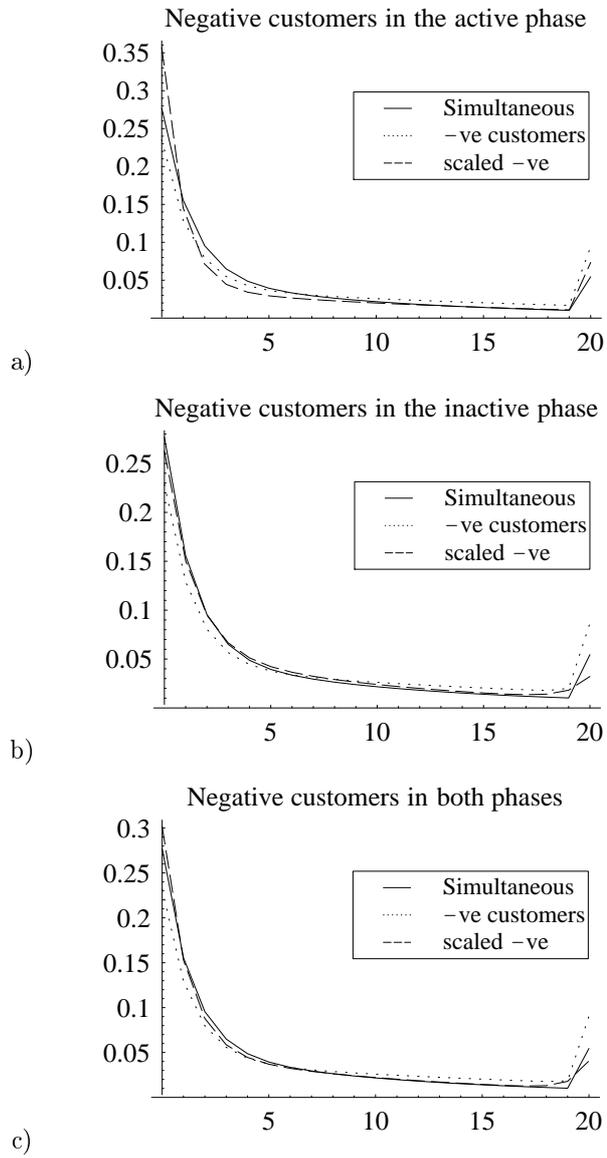
Figure 4: Comparing queue length distributions from the best-estimate using negative customers in the *inactive*, *active* and both phases with the simultaneous loss characteristic.

There appears to be scope for improving the match achieved by approximation using negative customers by careful choice of different rates in the two phases, but it will not be possible to avoid the dip in the queue length probability near the empty queue in the inactive phase, which is "topped up" from the active phase when simultaneous losses are enabled. In figure 3, compare the steady rise toward the empty queue in the inactive phase with simultaneous losses (top of the figure) with the dip in the characteristics using negative customers.

# 10 Conclusions

This queue with simultaneous loss on breakdown could be described as a queue which incorporates triggers: negative customer arrivals are triggered on processor breakdowns. The enriched modulation structure – which is no longer independent of queue length – is interesting.

While we can match the mean queue length arising from simultaneous loss on breakdown by using independent negative customer arrivals, the relationship between the negative customer rate and the breakdown/repair characteristics is not trivial, and the distribution not matched.

# References

[1] I. Mitrani and R. Chakka. Spectral expansion solution for a class of Markov models: Application and comparison with the matrix-geometric method, *Performance Evaluation* **23** pp. 241-260, 1995.

[2] D.J. Thornley and Harf Zatschler. An automated formulation of queues with multiple modulated batches *submitted to* IFIP WG 7.3 International Symposium on Computer Performance Modeling, Measurement and Evaluation. (Performance 2002)

[3] E. Gelenbe. Product form queueing networks with negative and positive customers, *Journal of Applied Probability* **28**, pp. 656-663, 1991.

[4] Bini DA, Latouche G, Meini B. Solving matrix polynomial equations arising in queueing problems *Linear Algebra and its Applications* 340: 225-244 Jan 1 2002

[5] R. Chakka and P.G. Harrison. A Markov modulated multi-server queue with negative customers - The MM CPP/GE/c/L G-queue. Acta Informatica **37**(11-12), pp. 881-919, 2001

[6] Stephen Wolfram, *The Mathematica Book*, 4th ed., (Wolfram Media/Cambridge University Press, 1999)

[7] P.G. Harrison. The MM CPP/GE/c/L G-queue: sojourn time distribution, *Queueing Systems: Theory and Applications*, accepted to appear 2002.