

1

2

Performance Evaluation and Optimization of Open Zero-Buffer Multi-Server Queueing Networks

3

4

R. Andriansyah^a, T. van Woensel^{b,*}, F. R. B. Cruz^c, L. Duczmal^c

5

^a*Department of Mechanical Engineering, Eindhoven University of Technology, Eindhoven, The Netherlands*

6

^b*Department of Technology Management, Eindhoven University of Technology, Eindhoven, The Netherlands*

7

^c*Department of Statistics, Federal University of Minas Gerais, Belo Horizonte, MG, Brazil*

8

Abstract

9 Open zero-buffer multi-server general queueing networks occur throughout a number of physical systems in the semi-
10 process and process industries. In this paper, we first analyze and quantify the performance of these systems in terms
11 of throughput using the generalized expansion method (GEM). We compare our results with simulation, and show
12 that the GEM gives an effective approximation for the performance measures of the systems under study. Secondly,
13 we embed the performance evaluation via the GEM in a multi-objective optimization setting. This multi-objective
14 optimization approach results in the Pareto efficient curves optimizing both the number of servers used and the
15 throughput. Experiments for a large number of settings and network topologies are presented in detail.

16 *Key words:* Zero Buffer Systems, Generalized Expansion Method, Performance Evaluation and Optimization, Genetic
17 Algorithms

18 1. Introduction and Motivation

19 Bufferless networks occur throughout a number of real-life physical systems in the semi-process and process
20 industries. A zero-buffer production environment might be created due to the processing technology of the
21 product itself, or simply due to the absence of any storage capacity between two consecutive operations of

* Corresponding author.

Email addresses: `r.andriansyah@tue.nl` (R. Andriansyah), `t.v.woensel@tm.tue.nl` (T. van Woensel), `fcruz@est.ufmg.br`
(F. R. B. Cruz), `duczmal@est.ufmg.br` (L. Duczmal).

22 a job. More specifically, the system of interest is a zero-buffer multi-server general queueing network, which
23 operates as follows: jobs arrive to the first node in the system with a certain arrival rate λ . The job is then
24 processed by an available server with service rate μ . After the service is finished, the job proceeds to the
25 next connected node in the system only if a server is available at this node. If this is not the case, then the
26 finished job has to wait at the previous server until there is an available server at the downstream node.
27 There is no queueing space between servers: the only buffer space available in each node are the servers
28 themselves. Due to the possibility of blocking, the throughput rate of the network might be smaller than
29 the arrival rate.

30 Hall and Sriskandarajah (1996) describe a steel production process. Here, molten steel goes through
31 a series of operations ranging from ingots, un-molding, reheating, soaking, and preliminary rolling. In this
32 process, the steel must pass through one operation to another continuously, without any waiting or buffering
33 of work in process, since such waiting would result in cooling down the steel to a temperature that is
34 not acceptable for the next process. Hence, either a job is finished and transferred directly to the next
35 process, or it is buffered in the machine itself until the downstream process is ready to receive another
36 job. In food-processing, no buffer space is allowed between the cooking operation and before the canning
37 operation. This is due to the requirement that the product should still be fresh when it is canned. Similar
38 issues can be found in producing juice and beer (see *e.g.*, Fey (2000)). In this examples, restrictions in
39 the processing technology creates zero-buffer production system. Ramudhin and Ratliff (1995) studied a
40 condiments manufacturer, with mayonnaise and various types of salad dressing as the product. The nature
41 of the product dictates hygienic consideration as one of the critical factor in production. Hence, there is
42 no space for work-in-process inventory and the product must never wait between two operations. As a last
43 example, third generation mobile communication networks are characterized by a multi-server zero-buffer
44 queueing system (Tsybakov, 2002). In such systems, arrivals are represented by requests of audio, data, and
45 video messages, whereas the service time is the message transmission time. Here, zero-buffers are caused
46 due to absence of storage capacity between operations.

47 Despite the high industrial relevance of zero-buffer networks particularly in process and semi-process
48 industries, only scant literature is available quantifying the performance or optimizing these networks. We
49 refer to Fransoo and Rutten (1994) for more information on process industries. This paper suggests a solution
50 to this challenging problem by providing a methodology for the performance evaluation and optimization of
51 zero-buffer systems in an arbitrary network configuration. The major contributions of this paper are:

- 52 (i) In the literature, few papers are found on the specific zero buffer systems analyzed in the present paper.
53 Moreover, these papers are mainly limited to single server tandem line settings. Here, we examine the
54 performance of open zero-buffer multi-servers with arbitrarily configured queueing networks;
- 55 (ii) We compare our results with simulation, and show that the generalized expansion method (GEM)
56 proposed by Kerbache and Smith (1988) gives an excellent approximation for the performance measures

57 of the systems studied.

58 (iii) We present a multi-objective optimization methodology for zero-buffer systems, which provide par-
59 ticularly insightful results. It is the first time to our knowledge that a multi-objective approach is
60 considered for the optimization of bufferless systems.

61 (iv) We demonstrate the usefulness of the developed tools in the system design phase of zero-buffer sys-
62 tems. More specifically, we elaborate on several network topologies, including a complex, arbitrarily
63 configured zero-buffer system, and provide answers to some difficult managerial questions regarding
64 the system's performance and its optimal configuration.

65 The paper is structured as follows. First, a brief literature review on zero-buffer queueing networks is
66 presented. Then, we describe how we obtain the blocking probability, one of the most important performance
67 measures in the analysis. Next, the performance evaluation methodology applied is described. In this paper,
68 we use the GEM to obtain the relevant performance measures. After this, we use a genetic algorithm approach
69 to optimize this kind of queueing networks. In the following section, we elaborate on the experimental results
70 obtained for a large number of situations (*i.e.*, tandem, split, merge cases, and large complex topologies). We
71 analyze the performance of these systems both analytically and by simulation. The last section concludes
72 this paper, with final remarks and topics for future research in the area.

73 2. Literature Review

74 The research in the area of queueing networks is very active, resulting in a vast amount of journal and
75 conference papers, books, reports etc. For a general and complete classification of queueing networks, the
76 reader is referred to, *e.g.*, Walrand (1988). Queueing networks can be divided into two categories (Perros,
77 1994): unrestricted and restricted queueing network. Unrestricted queueing networks include cases where
78 all nodes within the network have an unlimited capacity. Restricted queueing networks have limited (finite)
79 capacity for all nodes in the network. These queueing networks often represent real-life systems, where there
80 normally exist finite spaces for holding entities. In this paper, we focus on investigating restricted queueing
81 networks.

82 2.1. Restricted Queueing Networks

83 Restricted queueing networks have a finite capacity in each node, referred to as the total buffer capacity
84 of with size K_j . That is, a finite node j can only hold entities up to a certain quantity K_j including those
85 entities in service. The buffer capacity at finite node j causes blocking to occur when the arriving quantity to
86 node j exceeds its buffer capacity K_j (Buzacott and Shanthikumar, 1993). As a consequence, each node in
87 the network might be affected by events at other nodes, leading to the phenomena of blocking and starvation

88 (Perros, 1994). In general, three blocking mechanisms can be distinguished, blocking-after-service (BAS),
89 blocking-before-service (BBS), and repetitive-service (RS). These three types of blocking mechanisms are
90 presented in Figure 1. Most production lines operate under the BAS system. Moreover, in the literature it
91 is the most common assumption regarding buffer behavior (Dallery and Gershwin, 1992).

Figure 1 goes around here

92 BAS systems occur when a finished job in node i attempts to move downstream to node j whose buffers
93 are still being occupied by other jobs. In this sense, buffers at downstream nodes are full and hence cannot
94 receive the finished job from the upstream node. It is referred to as BAS due to the fact that the blocked job
95 is already finished with service in the previous node. It is therefore blocked while occupying the upstream
96 server. In the literature, this type of blocking is also known as type-1 blocking (Perros, 1994).

97 BBS systems are characterized by situations in which jobs in upstream node i have already determined
98 its destination downstream node j even before they are serviced. If a job at upstream node i finds full its
99 downstream destination node j , the server at node i becomes blocked and thus cannot process the job. In
100 this sense, the blocking occurs before the job is processed in upstream node i . Jobs in upstream node will
101 have to wait until there is available space in its downstream destination node, either in the queue or server.
102 In the literature, this type of blocking is known as type-2 blocking (Perros, 1994).

103 RS systems occur when a finished job at upstream node i cannot continue to downstream node j because
104 the downstream node is full, and the job immediately receives another service at node i . This repeating
105 process continues until downstream node j is not full. In the literature, this type of blocking is known as
106 type-3 blocking (Perros, 1994).

107 2.2. Zero-Buffer Queueing Networks

108 A particular case where a queueing network has a finite capacity but no buffers before servers is denoted
109 as a zero-buffer queueing network (also denoted in the literature as bufferless, no intermediate buffers etc.).
110 In this specific case, the buffer space at node j is equal to the number of server c in that node, that is,
111 $K_j = c$. Given that there is no space to queue, a job in the upstream node can only enter the downstream
112 node if the servers have finished processing their jobs (see Figure 2).

Figure 2 goes around here

113 For small tandem line zero-buffer networks, exact solutions are available. Small asynchronous single server
114 tandem lines with zero buffer and reliable machines were partly analyzed by several authors (Hildebrand,
115 1968; Hillier and Boling, 1967; Muth and Alkaff, 1987; Rao, 1976a,b). Most of this exact analysis is based
116 on Continuous Time Markov processes. Here, we will rely on simulations and on approximation techniques
117 for analyzing more complex queueing networks (including splits, merges, and multiple servers per node and
118 symmetrical/asymmetrical settings). More specifically, for the system on-hand in this paper, we will obtain

119 the performance measures via the GEM (discussed in detail in section 3). In Kendall notation, we look at
 120 $M/M/c/K$ queueing models and set the buffer size equal to the number of servers, $K = c$, resulting in
 121 $M/M/c/c$ queueing models.

122 2.3. Blocking probabilities

123 For $M/M/c/c$ models, it is known that the probability of the system being full at any time in steady
 124 state, given that the maximum number of jobs in the system equals to the number of servers in the system,
 125 corresponds to the Erlang's loss formula:

$$126 \quad p_c = \frac{\left(\frac{\lambda}{\mu}\right)^c / c!}{\sum_{i=0}^c \left(\frac{\lambda}{\mu}\right)^i / i!} \quad (1)$$

127 with λ being the arrival rate, μ the service rate and c , the number of servers.

128 It turns out that Eq. (1) also holds for any service-time distribution (other than Markovian). More
 129 specifically, this means that for $M/G/c/c$ queueing systems with no waiting space, the above equation is
 130 also valid regardless of the service time variability or the distribution itself. The expressions for the steady-
 131 state probabilities for this $M/G/c/c$ case are identical to those for the corresponding $M/M/c/c$ system (see
 132 Gross and Harris, 1998). This means that all results presented here, also immediately hold for the general
 133 service-time distribution case.

134 3. Performance Evaluation of Zero Buffer Networks

135 The GEM is a robust and effective approximation technique developed by Kerbache and Smith (1987). It
 136 has been successfully used to estimate performance measures for finite queueing networks. As described in
 137 previous papers, this method is basically a combination of repeated trials and node-by-node decomposition
 138 in which each queue is analyzed separately and then corrections are made in order to take into account the
 139 interrelation between the queues in the network. The GEM uses BAS, which is prevalent in most production
 140 and manufacturing, transportation, and other similar systems.

141 In this section, we will only present a high-level overview of the method. For more detailed information
 142 and applications of the GEM, the reader is referred to the papers by Kerbache and Smith (1988), Spinellis
 143 et al. (2000), Smith and Cruz (2005), Cruz and Smith (2007), and Cruz et al. (2007a). The GEM involves
 144 three stages: network reconfiguration, parameter estimation, and feedback elimination. Before discussing
 145 each stage in detail, the variables used in this section are defined in Table 1.

Table 1 goes around here

146 3.1. *Network reconfiguration*

147 The first step in the GEM involves reconfiguring the network: an artificial node is added for each finite node
 148 in the network. The artificial node is added to register the blocked customers at the finite node (Kerbache
 149 and Smith, 1988).

150 Following Kerbache and Smith (2000), the GEM thus creates for each finite queue, represented by vertex
 151 j , an auxiliary vertex h_j , modeled as an $M/G/\infty$ queue (see Figure 3). When an entity arrives to the system,
 152 vertex j may be blocked with probability p_{c_j} , or unblocked, with probability $(1 - p_{c_j})$. Under blocking, the
 153 entities are rerouted to vertex h_j for a delay while node j is busy. Vertex h_j helps to accumulate the time
 154 an entity has to wait before entering vertex j and to compute the effective arrival rate to vertex j . In other
 155 words, the GEM ultimate goal is to provide an approximation scheme to update the service rates of upstream
 156 nodes that takes into account all blocking after service in there, caused by downstream nodes:

157
$$\tilde{\mu}_i^{-1} = \mu_i^{-1} + p_{c_j}(\mu'_{h_j})^{-1}.$$

Figure 3 goes around here

158 3.2. *Parameter estimation*

159 In the second stage, the parameter estimation of all unknowns for the system described above are deter-
 160 mined. These unknowns are in this case (Kerbache and Smith, 1987):

- 161 (i) The probability of a customer being blocked, p_{c_j} ;
- 162 (ii) The probability that a customer is forced back to the holding node given he was rejected at the previous
 163 trial, p'_{c_j} ;
- 164 (iii) The service rate at the holding node, μ_{h_j} ;

165 The reader is referred to the paper by Jain and Smith (1994) for more details.

166 3.3. *Feedback elimination*

167 The repeated visits to the holding nodes (due to the feedbacks), create strong dependence in the arrival
 168 process. Therefore, the repeated immediate feedback is eliminated. This is done by giving the customer
 169 enough service time during the first passage through the holding node. The adapted service rate for the
 170 holding node μ'_{h_j} then becomes:

171
$$\mu'_{h_j} = (1 - p'_{c_j})\mu_{h_j}. \tag{2}$$

172 And the adapted squared coefficient of the service rate for the holding node $(c'_{sh_j})'$ becomes:

173
$$(c'_{sh_j})' = p'_{c_j} + (1 - p'_{c_j})c_{sh_j}^2. \tag{3}$$

174 The effective service rate in node 2 will be the average of μ_2^{-1} when there is no blocking and $\mu_2^{-1} + (\mu'_{h_j})^{-1}$
 175 when there is blocking or:

$$176 \quad \mu_{2\text{eff}}^{-1} = \mu_2^{-1} + p_{c_j} (\mu'_{h_j})^{-1}. \quad (4)$$

177 4. Multi-Objective Optimization of Zero Buffer Networks

178 In this section, we consider the optimization of zero-buffer systems where the GEM is incorporated in
 179 a genetic algorithm (GA) approach. Since there are no buffers in the network, the optimization problem
 180 on-hand is essentially a server allocation problem for zero-buffer general service queueing network. As servers
 181 are expensive, one would like to minimize the amount of servers as much as possible. On the other hand,
 182 one would also like to maximize the throughput of the network. This throughput is affected by the allocated
 183 servers, *i.e.* more servers leads to a higher throughput.

184 4.1. Model Formulation

185 The problem described can be formulated as a multi-objective optimization problem with two objec-
 186 tives: minimizing the total number of servers and maximizing the throughput. The following mathematical
 187 formulation represents the multi-objective optimization problem.

$$188 \quad \text{optimize } Z = \{f_1(\mathbf{x}), f_2(\mathbf{x})\} \quad (5)$$

189 s.t.:

$$x_i \in \{1, 2, \dots\}, \forall i, \quad (6)$$

190 in which

$$f_1(\mathbf{x}) = \sum_i c_i \times x_i, \text{ where } x_i \text{ is the server allocation to node } i \text{ and } c_i \text{ is the cost of a server at node } i,$$

$$f_2(\mathbf{x}) = \Theta(\mathbf{x}), \text{ the throughput profit.}$$

191 In the above formulation, we trade-off throughput versus total number of servers. Of course, we could
 192 have used Work-In-Progress (WIP) or cycle time. We feel, however, that the first criterion in the network
 193 optimization should be assuring a certain throughput compared to the arrival rate. As such, we consider the
 194 other variables WIP or cycle time as secondary criteria. Of course, the methodology can be extended taking
 195 into account these variables. For the sake of presentation of the results and the insights, we choose not to
 196 do this (see *e.g.* Smith and Cruz, 2005, for a similar approach).

198 To solve the above problem formulation, we use a powerful class of optimization heuristics, called Genetic
 199 Algorithms (GA). One particular quality of a GA is that it is capable of generating an entire set of multi-
 200 objective solutions called the *Pareto-optimal set*. As such, a GA searches for a set of optimal solutions of a
 201 problem with more than one objective function in a single run. A GA finds the set of *non-dominated* solutions,
 202 that is, the solutions that are better off than another solution with regards to at least one objective, while
 203 not being worse off in any other objectives than another solution. In a minimization problem, for example,
 204 solutions $x \in F$ that are not a part of the Pareto-optimal set are said to be dominated because there are
 205 other solutions $z \in F$ such that $f(z) < f(x)$, and $f(z) \neq f(x)$. As such, $f(z)$ is better off than $f(x)$ in at
 206 least one objective function, without being worse off in the other objective functions. The resulting Pareto
 207 curve allows one to see the tradeoff among the objective function values from different solutions. That is,
 208 a GA allows for evaluating the loss in one objective (*e.g.* the throughput) compared to a simultaneous
 209 enhancement in another objective (*e.g.* the server allocation).

210 The optimization procedure in a GA starts with the evaluation of several points along the search space.
 211 This set of individual points is called the population. The initial population converges to the optimum
 212 population through several operators as follows (Carrano et al., 2006):

- 213 (i) **Mutation:** Randomly modify some individuals to reach other points in the search space.
- 214 (ii) **Crossover:** Combine the randomly organized pairwise individuals such that each former pair of indi-
 215 viduals creates a new pair.
- 216 (iii) **Selection:** Evaluate the individuals after mutation and crossover. Using a probabilistic rule, these
 217 individuals are chosen or not chosen to be included in the new population. The rule is organized
 218 so that individuals with better objective value have greater probability to be selected into the new
 219 population.
- 220 (iv) **Elitism:** The “best” individuals from the former population are deterministically inserted into the
 221 new population.

222 Particularly important for a multi-objective optimization problem is defining the selection and the elitism
 223 operators to determine the best individual from the search space. One of the most widely used selection and
 224 elitist operator is the non-dominated sorting scheme, known as the NSGA-II (see *e.g.* Deb et al. (2002)). It
 225 has been shown that the NSGA-II requires less computational effort as compared to other GA algorithms.
 226 For a complete elaboration on this algorithm, see Deb et al. (2002), and for details on the version of the
 227 particular implementation used in this article, see Cruz et al. (2007b). We will use a GA as the search
 228 method in combination with the GEM to obtain the performance measures of the network.

229 5. Experimental Results

230 We first discuss the results of using the GEM as a performance evaluation tool. Secondly, we use a GA
231 heuristic in combination with the GEM to solve the multi-objective optimization problem as described
232 previously. Finally, we provide a complex topology to demonstrate the performance of our approach in the
233 design of large and complex zero buffer networks. All algorithms described are coded in Fortran and are
234 available from the authors upon request.

235 5.1. Performance Evaluation

236 To evaluate the performance of the GEM in our network settings, we conducted experiments using three
237 different topologies: series, split, and merge topologies. We examine these topologies for both symmetrical
238 as asymmetrical settings. Please bear in mind that the range of possible experiments is exponential itself, so
239 we have determined a select sample to present. The configurations tested are shown in Figures 4, 5, and 6.

Figure 4 goes around here

Figure 5 goes around here

Figure 6 goes around here

240 5.1.1. Symmetrical Networks

241 In each topology, we analyze several number of nodes, $N \in \{3, 5, 9\}$, and servers, $c \in \{2, 4, 10\}$. We set
242 different values for the arrival rates, $\lambda \in \{2, 4, 8, 16\}$. In all settings a service rate with a mean of $\mu = 10$
243 was used. With these combinations of parameters, we ended up with 36 different systems for each of the
244 topologies (series, merge, and split). The series topology consists of a series of nodes (N), in which each node
245 has c number of servers (see Figure 4). In the split topology (see Figure 5), we divide the arriving jobs to
246 subsequent nodes after departing from the first node in the system. That is, the splitting node is positioned
247 directly after the first node. The routing probability for both routes is set to be equal to each other (0.5).
248 Every case in this topology has two ending nodes in the system ($N - 1$ and N). In the merge topology (see
249 Figure 6), the jobs arrive from two different source nodes. The overall arrival rate is then divided equally
250 for the two source nodes. The last node merges the two streams.

251 Table 2 shows the results for all different sets and topologies. The results reported are the overall through-
252 put at the network. Table 2 is read as follows: the 3-node experiment with an arrival rate $\lambda = 2$ and the
253 number of servers c equal to $(2, 2, 2)$ results in a throughput of 1.996. In general, we observe that the blocking
254 effect becomes more severe for high arrival rates relative to the service rates. For those cases, the throughput
255 rate is reduced more compared to the arrival rate than the less congested settings. This rule holds true for
256 all three types of topologies involved in the experiment.

Table 2 goes around here

257 In order to evaluate the solution quality of the GEM, we set up a simulation experiment of four selected
258 cases in the series, merge, and split topologies. We used an observation period of 200,000 time unit and a
259 warm-up period of 2,000 time units for 20 independent replications. The simulation was carried out using
260 ARENA (Kelton et al., 2001). Table 3 is based on the simulation results. In the column labeled analytical, we
261 give the throughput result from the GEM for each of the cases. We then compare this analytical result with
262 the average result obtained via the simulation. The column δ refers to the half-width of the 95% confidence
263 interval. Also included in the tables is the % deviation for the analytical results on the throughput, $\Delta\% \theta \stackrel{\text{def}}{=} \frac{(\theta - \theta^s)}{\theta^s} \times 100\%$.
264 We see that the analytical results for all cases fall within the 2% error range, which is
265 very reasonable and acceptable (see column $\Delta\% \theta$). Based on the simulation output, we conclude that the
266 analytical results are reliable.

Table 3 goes around here

267 5.1.2. *Asymmetrical Networks*

268 In the previous section, we considered cases with symmetrical settings in the service rates for the different
269 nodes. In this section, we set up some experiments unbalancing the routing probabilities (for the split
270 topologies) and the arrival rates (for the merge topologies), and assuming different service rates, μ 's, and
271 different number of servers, c 's, along the queueing network. The results are shown in Table 4. Simulation
272 results from selected cases are presented in Table 5. Based on the asymmetrical results, we can see that
273 the main conclusions, as observed for the symmetrical results, hold being the highest error below 5%. We
274 see that the blocking effects becomes more important with the increase of the arrival rates. With regards
275 to the simulations, we note that the GEM tends to both overestimate and sometimes underestimate the
276 throughput as compared to the "true" throughput from simulation (as reflected in the values for $\Delta\% \theta$).
277 This is in line with what should be expected when comparing simulation with analytical results.

Table 4 goes around here

Table 5 goes around here

278 5.2. *Network Optimization*

279 In this section, we identify the optimal server configuration given the network structure, arrival rate,
280 and processing rates similar to that of the previous sections. We consider both symmetrical as well as
281 asymmetrical settings. In this section, we will again focus on the earmarked networks from Tables 2 and 4.
282 Using the GA in conjunction with the GEM, we identify the Pareto-optimal set, which is the set of optimal
283 solutions for both objectives in the objective function, *i.e.* the minimal number of servers and the maximal
284 throughput. After much experimentation, the population size was set to 20 and the number of generations

285 to 100, 1,000, and 10,000, for 3, 5, and 9 nodes, respectively, in order to reach a good trade-off between a
286 fair approximation of the Pareto-optimal set and low computational times.

287 The resulting Pareto-optimal sets for the earmarked symmetrical settings are provided in Figure 7. The
288 figures show that the throughput of all the earmarked *symmetrical* settings lie on the Pareto-optimal set
289 (pointed with an arrow). As such, these settings are optimal given the corresponding network structure,
290 arrival rate, and processing rates. A valuable insight from the Pareto-optimal set generated by the GA is
291 that we are able to quantify the additional throughput gained by adding one or more servers into the network
292 and optimizing the allocation of the additional servers. For example, the Pareto graphs in 7 (c), (f), and
293 (i) show that too many servers are utilized and that less servers could have been used without reducing
294 the throughput significantly. Eventually, one may use this information to evaluate whether the cost of
295 adding extra servers into the network is justified by the benefit gained from the additional throughput. This
296 information is critical given the fact that there are cases where additional servers only give marginal increase
297 to the throughput, and other cases where additional servers give significant increase to the throughput.

Figure 7 goes around here

298 Figure 8 shows the GA optimization results for the earmarked *asymmetrical* settings. We observed that
299 the throughput from all the earmarked settings is below the Pareto-optimal set (see in Figure 8, pointed
300 with an arrow), which means that our settings for the network structure, arrival rates, and processing rates
301 were not optimal in the previous section. More specifically, for the same total number of servers in the
302 network, there exists another server allocation that results in a higher throughput. For example, in the
303 case of asymmetrical serial topology with 3 nodes and processing rates of (12, 11, 10) for nodes 1, 2, and 3
304 respectively, the GA found the optimal server allocation as (7, 6, 3) for node 1, 2 and 3. This configuration
305 results in a throughput of 2.000 as opposed to a throughput of 1.998 when the server allocation is 2,4,10,
306 as in the earmarked asymmetrical setting. Note that both results have a total number of 16 servers, but
307 that the actual configuration can be improved such that the throughput can be increased. Depending upon
308 the managerial preference, one can thus decide to increase the throughput by using the same number of
309 servers but in a different configuration (*i.e.* find vertically the nearest configuration on the Pareto-optimal
310 set). Alternatively, one can also accept the current throughput but achieve this throughput with less servers
311 used (*i.e.* find horizontally the nearest configuration on the Pareto curve).

Figure 8 goes around here

312 5.3. A large and complex topology

313 In this section, we analyze a large and complex topology. This topology is inspired by an example fruit
314 juice blending and packaging plant, of which a detailed description can be found in Fey (2000). Using this
315 complex topology, we elaborate in detail on the merits of the methodology, by applying the GEM and the

316 multi-objective approach to this large design problem. More specifically, the queueing network structure
317 given in Figure 9 is analyzed.

Figure 9 goes around here

318 It should be clear that this queueing network combines all three topologies that have been investigated
319 separately in the previous section (namely series, merge, and split). It is interesting to see whether the
320 proposed methodology performs well in such a complicated queueing network setting. In this section, three
321 critical decisions are investigated in designing and optimizing this zero buffer production system:

- 322 (i) The effect of different arrival rates to the throughput of the system.
- 323 (ii) The effect of different service rates of the servers to the throughput of the system.
- 324 (iii) The effect of a different network structure to the throughput of the system

325 We run the GEM using the parameter values given in Table 6. We set the arrival rate to 30.215 units/h;
326 we consider two situations, namely a slow server setting and a fast server setting (see Table 6). The slow
327 server setting is created by using low processing rates of machines as compared to the arrival rate. On the
328 contrary, the fast server setting is characterized by high processing rates of machines as compared to the
329 arrival rate. As such, we would like to see the performance of the GEM in the presence of both low and high
330 blocking situations.

Table 6 goes around here

331 We compare the results from simulation with that of the GEM and obtain the results depicted in Table 7.
332 The GEM performs quite well given such a complex queueing network setting (that is, queueing network
333 with serial, merge and split topologies with a total of 137 servers involved). There is a notable reduction
334 of the GEM's accuracy under a queueing network setting with high blocking probabilities. However, it is
335 clear that the GEM is able to approximate the throughput of a very complex queueing setting in a matter
336 of seconds. This eliminates the necessity to conduct lengthy simulations, as indicated in Table 7 for the
337 required CPU time.

Table 7 goes around here

338 Again using the GA in conjunction with the GEM, we derive the Pareto-optimal set for the minimal
339 number of servers and the maximal throughput. We focus on the results for the slow server case discussed in
340 Table 6. Figure 11 gives the results for the Pareto set for increasing populations sizes. Clearly, the Pareto-set
341 converges for a reasonable size of the population (from 100). In the figure, the white dot corresponds to
342 the throughput for the server vector as used in the higher experiments shown in Tables 7 (without the GA
343 optimization). The Figure clearly shows that it is possible to either increase the throughput with the same
344 number of servers (vertical line from white dot), or to reduce the total number of servers with the same
345 throughput (horizontal line from white dot). Similar results are obtained for the fast server case.

Figure 11 goes around here

346 The fact that the methodology is able to approximate the throughput with high accuracy and within

347 a short amount of time allows us to analyze diverse system design problems, which otherwise would be
348 extremely time consuming to do by using simulation. Following these encouraging results, we will next
349 consider three particular system design problems of queueing networks, namely (1) evaluating the effect
350 of the arrival rate, (2) evaluating the effect of different processing rates, and (3) evaluating the network
351 structure itself.

352 *Different arrival rates*

353 We consider the case in which it is possible to change the arrival rate into the queueing network, while
354 holding the service rate of each server constant. This situation typically denotes the case where a plant
355 manager wants to investigate the threshold of arrival rate beyond which the current server setting is incapable
356 of delivering the desirable throughput. To tackle this issue, we vary the arrival rate into the system and
357 analyze the percentage of blocking compared to the arrivals. Using the same service rate as provided in
358 Table 6, we obtain the throughput for different arrival rates.

359 In Figure 10, the results from different arrival rates give us insight about the system's capability of handling
360 different arrival volumes. We can see, for example, that the slow server configuration of the system is able
361 to handle an arrival rate up to 10 units/h with a resulting blocking of less than 5%, while the fast server
362 configuration would handle much more. Higher arrival rates will cause a significant increase in blocking
363 percentage, as it can be seen from Figure 10. As a practical implication, the plant manager can use this
364 insight to determine the appropriate arrival rate such that the blocking probability is under the predefined
365 blocking probability threshold. This example shows how the GEM can be used effectively to predict system
366 performance given different arrival rates, and suggest the threshold arrival rate that will seriously impair
367 the performance of the system in terms of blocking.

Figure 10 goes around here

368 *Different processing rates*

369 Another common case in system design, involves the choice of the facilities used in the network. In
370 particular, we assume that there are several choices of facilities characterized with different processing rates.
371 We treat the arrival rate as an external variable that is difficult or impossible to change. Given this setting,
372 we would like to evaluate the effect on throughput of changing the service rates, while holding the arrival
373 rate and all other parameters constant. We again set the arrival rate fixed at 30.215 units/h, and use
374 three alternative server processing rates as given in Table 8. The number of servers for each facility is that
375 of Table 6. Given the fixed arrival rate, each alternative configuration is evaluated on the throughput as
376 provided in the last line of Table 8.

Table 8 goes around here

377 We assume that configuration #1 is the default configuration of the queueing network. Note that this con-
378 figuration is similar to the previous case where we consider high utilization setting. The resulting throughput
379 from the GEM is 17.422 units/h, as before. Given the processing rates and the number of available servers,
380 one can easily notice that *Phase 2* acts as the bottleneck in the network. As such, one may increase the
381 processing rate of this particular node, as this could be advantageous for the throughput. Furthermore, the
382 last part of the system (*Phase 3*) also has a relatively high utilization, creating heavy imbalance in the
383 system. Table 8 gives the resulting throughput when the processing rates of these phases are increased.
384 Using this configuration, the throughput is increased by 6.1%. As a last example, we increase the process-
385 ing rates of all last three nodes in *Phase 3* in the system and the bottleneck node (see configuration #3).
386 Consequently, the throughput increases further to 16% compared to the default configuration.

387 This example shows how the GEM can be used effectively to predict system performance under varying
388 processing rates of servers in the zero buffer network. The result from this analysis can be used as a basis
389 for *i.e.* benefit/cost analysis in order to justify the potential investment for increasing the servers' capacity.
390 Following this line of analysis, one can argue whether or not the cost of increasing the capacity of the servers
391 is justified by the benefit from the additional throughput gained.

392 *Different network structures*

393 We now consider the case where it is possible to manipulate the structure of the queueing network in
394 order to increase the throughput of the network. One way to reduce congestion (and thus increase the
395 throughput) of a queueing network is by combining multiple sources of variability, which is known as the
396 concept of *variability pooling* (Hopp and Spearman, 2000). This can be achieved, in one way, by sharing the
397 queue through the use of more flexible machines.

398 Let us assume that by using more flexible machines it is possible to reduce the number of lines in *Phase 3*
399 from twelve to only three lines. As such, we assume that the new machines are capable of performing multiple
400 tasks. The new queueing network structure is given in Figure 12. Note that for the three-line structure, we
401 will use the data of configuration #1 in Table 8.

Figure 12 goes around here

402 Using the same total number of servers as in configuration #1, we obtain a throughput of 22.887 units/h,
403 higher than before for this three-line structure, as compared to the throughput of 17.422 units/h for the
404 twelve-line structure (see Table 8, configuration #1). The additional throughput gain can be the used to
405 evaluate whether it is beneficial to use more flexible packaging machines in the network.

406 6. Conclusions and Future Research

407 Throughout this paper, we have demonstrated the use of the generalized expansion method (GEM) to both
408 evaluate and optimize the performance of a zero-buffer queueing network. We considered three topologies of
409 queueing networks, namely series, split, and merge topologies, taking into account both symmetrical as well
410 as asymmetrical composition of the topologies. We obtained the throughput of the networks using the GEM
411 and compared the results with simulation. The insights from the multi-objective optimization of zero-buffer
412 queueing networks are proven to be invaluable particularly in the design phase of such systems.

413 We provided an large scale complex topology for which the GEM was used to approximate the throughput
414 rate. We argued that in addition to its accuracy both in low and high blocking probability settings, the
415 expansion method provides a relatively simple means to see how changes in arrival rate and/or processing
416 rate affects the performance of a zero buffer queueing networks. All in all, we found the GEM as a fast and
417 accurate approximation method to measure the throughput rate of zero-buffer queueing networks as well as
418 for optimization purposes of such networks.

419 *Acknowledgements*

420 The research of prof. Frederico Cruz has been partially funded by the CNPq (*Conselho Nacional de*
421 *Desenvolvimento Científico e Tecnológico*) of the Ministry for Science and Technology of Brazil, grants
422 201046/1994-6, 301809/1996-8, 307702/2004-9, and 472066/2004-8, the FAPEMIG (*Fundação de Amparo*
423 *à Pesquisa do Estado de Minas Gerais*), grants CEX-289/98, CEX-855/98, and TEC-875/07, and PRPq-
424 UFMG, grant 4081-UFMG/RTR/FUNDO/PRPq/99.

425 References

- 426 Buzacott, J., Shanthikumar, J. G., 1993. Stochastic Models of Manufacturing Systems. Prentice-Hall.
- 427 Carrano, E. G., Soares, L. A. E., Takahashi, R. H. C., Saldanha, R. R., Magela Neto, O., April 2006. Electric
428 distribution network multiobjective design using a problem-specific genetic algorithm. IEEE Transactions
429 on Power Delivery 21 (2), 995–1005.
- 430 Cruz, F. R. B., Duarte, A. R., van Woensel, T., 2007a. Buffer allocation in general single-server queueing
431 network. Computers & Operations Research. URL <http://dx.doi.org/10.1016/j.cor.2007.03.004>
- 432 Cruz, F. R. B., Duczmal, L., van Woensel, T., Smith, J. MacGregor., 2007b. A multi-objective approach
433 for buffer allocation in general queueing networks. In: Sixth Conference on the Analysis of Manufacturing
434 Systems. BETA and EURANDOM, TU/e, Lunteren, The Netherlands.

435 Cruz, F. R. B., Smith, J. MacGregor, 2007. Approximate analysis of $M/G/c/c$ state-dependent queueing
436 networks. *Computers & Operations Research* 34 (8), 2332–2344.

437 Dallery, Y., Gershwin, S. B., 1992. Manufacturing flow line systems: A review of models and analytical
438 results. *Queueing Systems* 12, 3–94.

439 Deb, K., Agrawal, S., Pratap, A., Meyarivan, T., 2002. A fast elitist non-dominated sorting genetic algorithm
440 for multi-objective optimization: NSGA-II. *IEEE Transactions on Evolutionary Computation* 6, 182–
441 197.

442 Fey, J., 2000. Design of a fruit juice blending and packaging plant. Ph.D. thesis, Eindhoven University of
443 Technology.

444 Fransoo, J. C., Rutten, W. G. M. M., 1994. A typology of production control situations in process industries.
445 *International Journal of Operations and Production Management* 14 (12), 47–57.

446 Gross, D., Harris, C. M., 1998. *Fundamentals of Queueing Theory*, 3rd Edition. Wiley Series in Probability
447 and Statistics, New York.

448 Hall, N. G., Sriskandarajah, C., 1996. A survey of machine scheduling problems with blocking and no-wait
449 in process. *Operations Research* 44, 510–525.

450 Hildebrand, D., 1968. On the capacity of tandem server, finite queue, service systems. *Operations Research*
451 16, 72–82.

452 Hillier, F., Boling, R., 1967. Finite queues in series with exponential or Erlang service times: A numerical
453 approach. *Operations Research* 16, 286–303.

454 Hopp, W.J. and Spearman, M.L., 2000. *Factory physics*. MacGraw Hill International Editions, New York,
455 NY, USA.

456 Jain, S., Smith, J. MacGregor, 1994. Open finite queueing networks with $M/M/C/K$ parallel servers. *Com-*
457 *puters & Operations Research* 21 (3), 297–317.

458 Kelton, D., Sadowski, R. P., Sadowski, D. A., 2001. *Simulation with Arena*. MacGraw Hill College Div.,
459 New York, NY, USA.

460 Kerbache, L., Smith, J. MacGregor, 1987. The generalized expansion method for open finite queueing net-
461 works. *European Journal of Operational Research* 32, 448–461.

462 Kerbache, L., Smith, J. MacGregor, 1988. Asymptotic behavior of the expansion method for open finite
463 queueing networks. *Computurs & Operations Research* 15 (2), 157–169.

464 Kerbache, L., Smith, J. MacGregor, 2000. Multi-objective routing within large scale facilities using open
465 finite queueing networks. *European Journal of Operational Research* 121, 105–123.

466 Muth, E., Alkaff, A., 1987. The throughput rate of three-station production lines: A unifying solution.
467 *International Journal of Production Research* 25, 1405–1413.

468 Perros, H. G., 1994. *Queueing Networks with Blocking*. Oxford University Press.

469 Ramudhin, A., Ratliff, H. D., 1995. Generating daily production schedules in process industries. *IEE Trans-*

- 470 actions 27, 646–656.
- 471 Rao, N., 1976a. A generalization of the bowl phenomenon in series production systems. *International Journal*
472 *of Production Research* 14 (4), 437–443.
- 473 Rao, N., 1976b. A viable alternative to the method of stages solution of series production systems with
474 Erlang service times. *International Journal of Production Research* 14 (6), 699–702.
- 475 Smith, J. MacGregor, Cruz, F. R. B., 2005. The buffer allocation problem for general finite buffer queueing
476 networks. *IIE Transactions on Design & Manufacturing* 37 (4), 343–365.
- 477 Spinellis, D., Papodopoulos, C., Smith, J. MacGregor, 2000. Large production line optimisation using sim-
478 ulated annealing. *International Journal of Production Research* 38 (3), 509–541.
- 479 Tsybakov, B., 2002. Optimum discarding in a bufferless system. *Queueing Systems* 41, 165–197.
- 480 Walrand, J., 1988. *An Introduction to Queueing Networks*. Prentice-Hall, Englewood Cliffs.

481 **Figure Legends**

482 **Figure 1:** The three different blocking mechanisms

483 **Figure 2:** A zero-buffer queueing network representation

484 **Figure 3:** The generalized expansion method

485 **Figure 4:** Series topology

486 **Figure 5:** Split topology

487 **Figure 6:** Merge topology

488 **Figure 7:** GA optimization for earmarked symmetrical settings

489 **Figure 8:** GA optimization for earmarked asymmetrical settings

490 **Figure 9:** Queueing network structure

491 **Figure 11:** GA optimization for the fruit juice plant

492 **Figure 10:** Plot of throughput for different arrival rates

493 **Figure 12:** New queueing network structure

494 **Table Legends**

495 **Table 1:** Variables used in this paper

496 **Table 2:** Throughput θ for the symmetrical cases

497 **Table 3:** Simulation results for the symmetrical cases

498 **Table 4:** Throughput θ for the asymmetrical cases.

499 **Table 5:** Simulation results for the asymmetrical cases

500 **Table 6:** Processing rates used in the complex topology

501 **Table 7:** Data used in the complex topology

502 **Table 8:** Throughput rates for different configurations

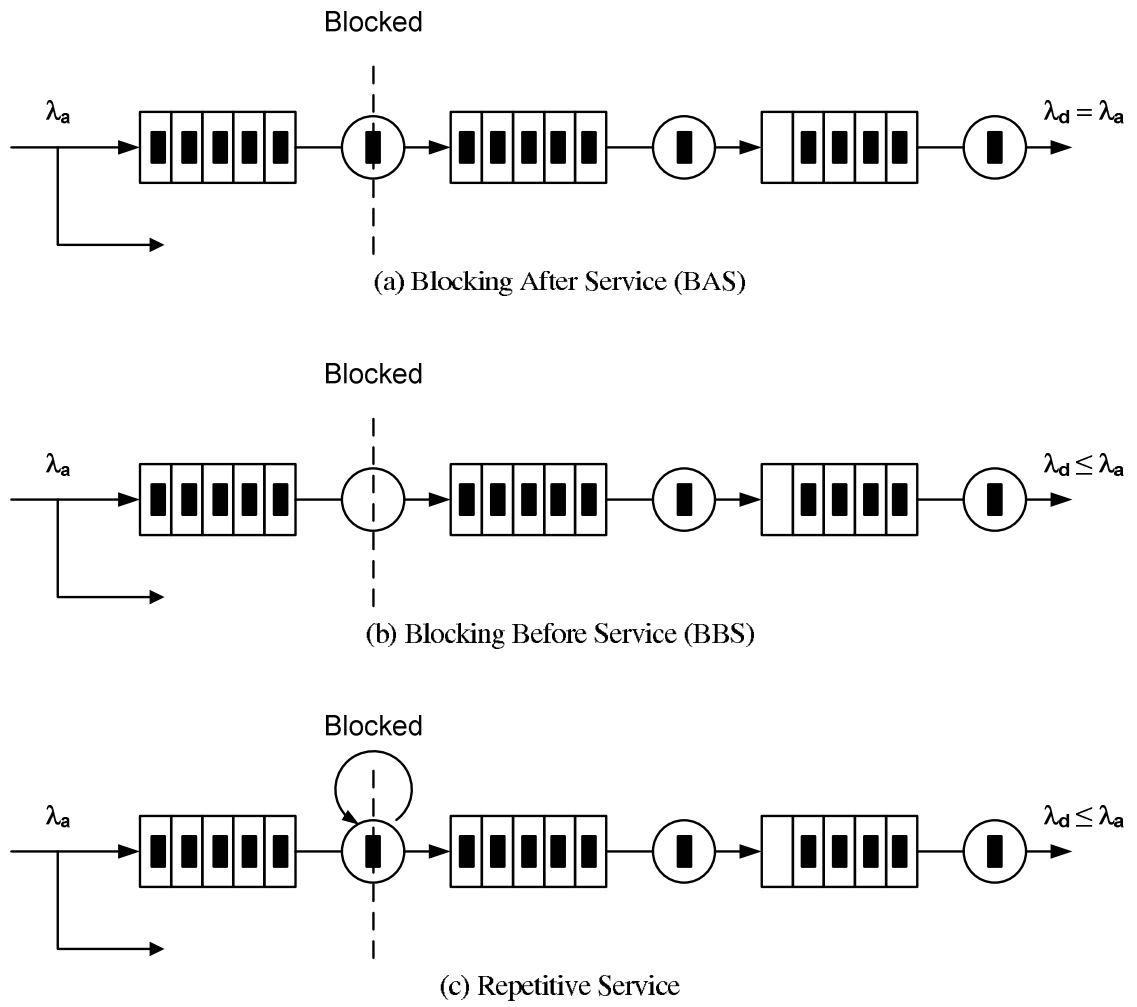


Fig. 1.

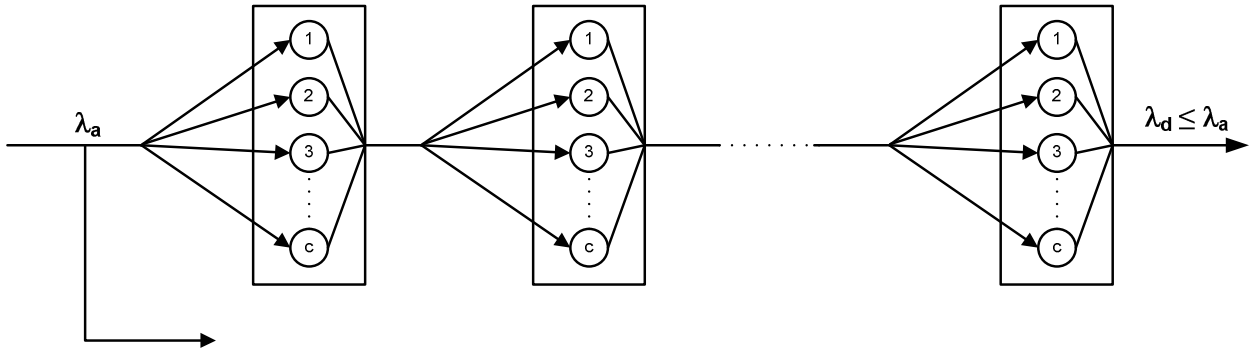


Fig. 2.

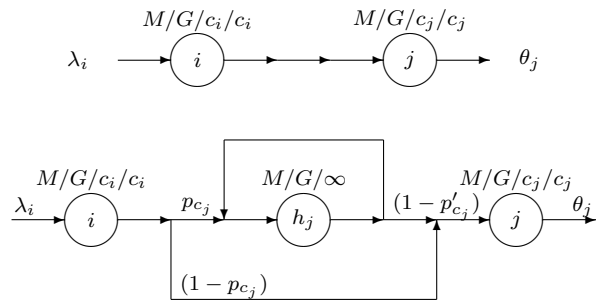


Fig. 3.

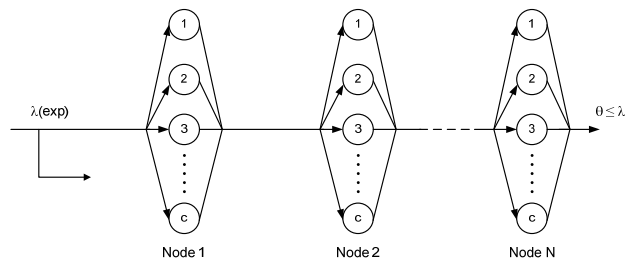


Fig. 4.

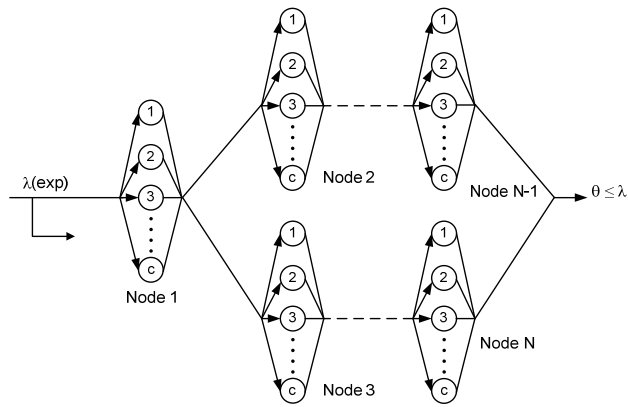


Fig. 5.

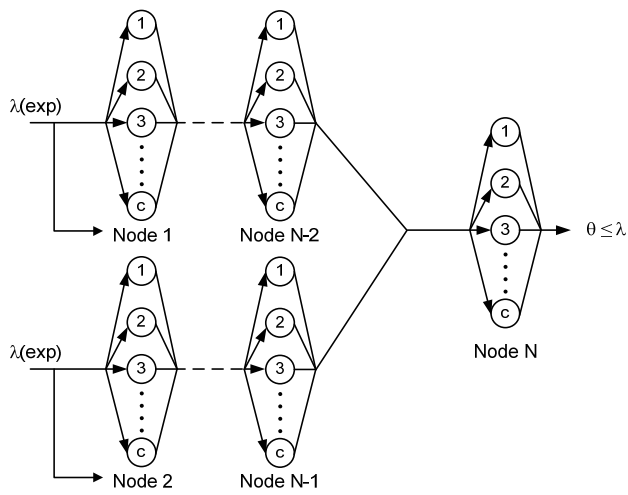
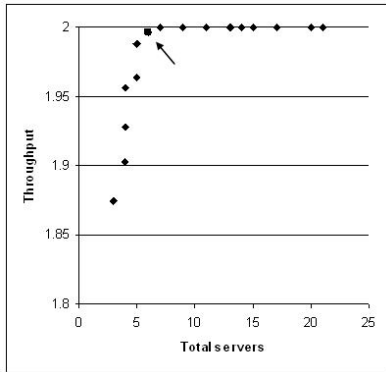
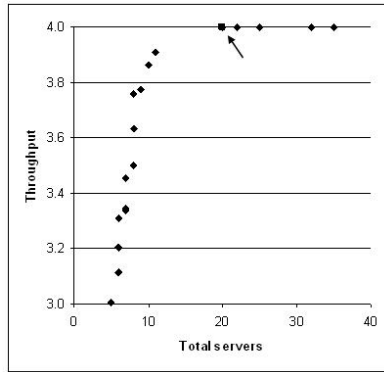


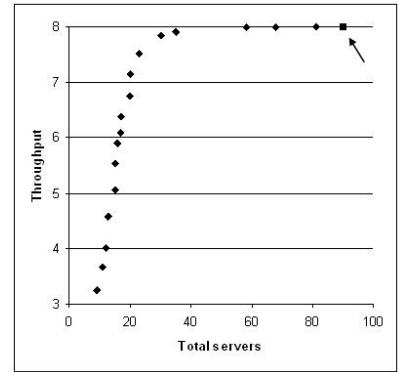
Fig. 6.



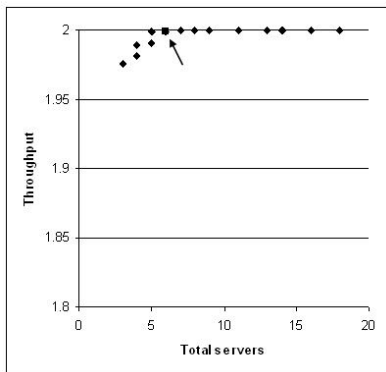
(a) Serial 3 nodes



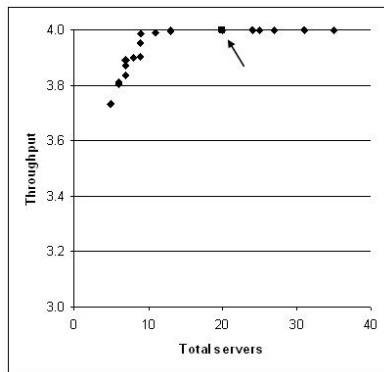
(b) Serial 5 nodes



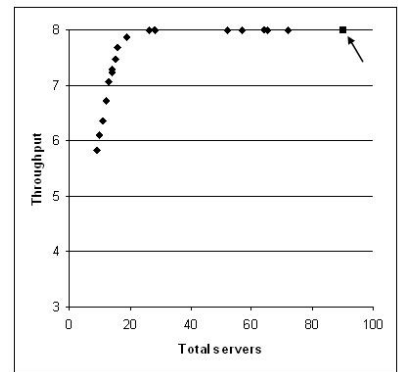
(c) Serial 9 nodes



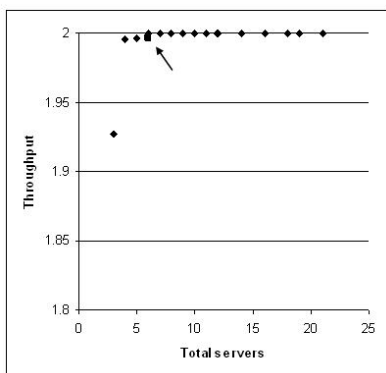
(d) Merge 3 nodes



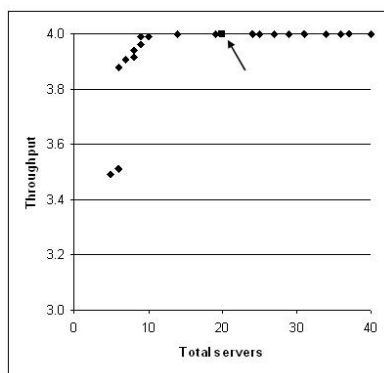
(e) Merge 5 nodes



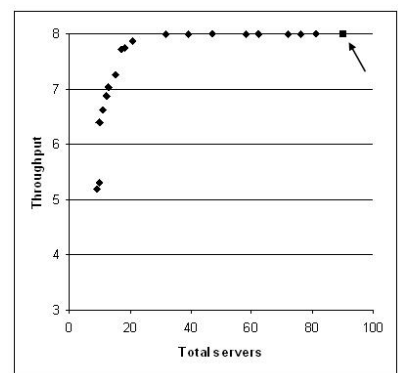
(f) Merge 9 nodes



(g) Split 3 nodes

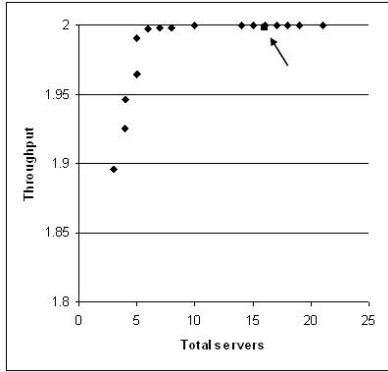


(h) Split 5 nodes

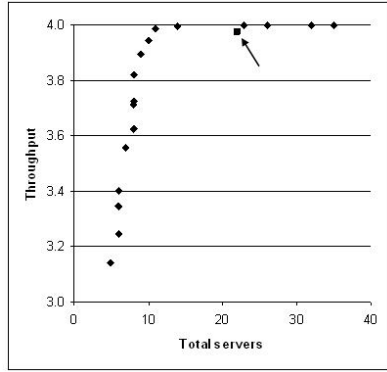


(i) Split 9 nodes

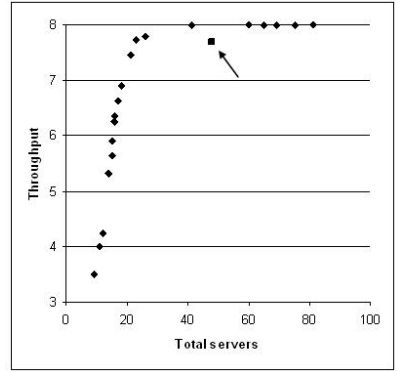
Fig. 7.



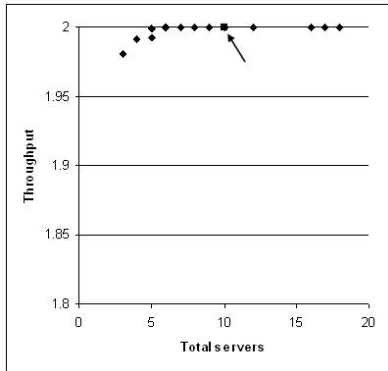
(a) Serial 3 nodes



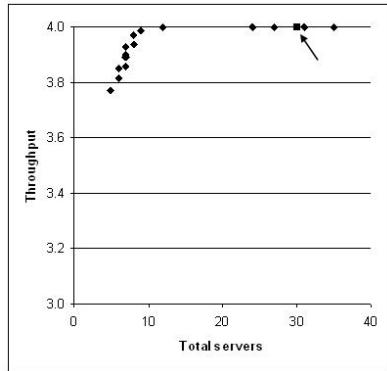
(b) Serial 5 nodes



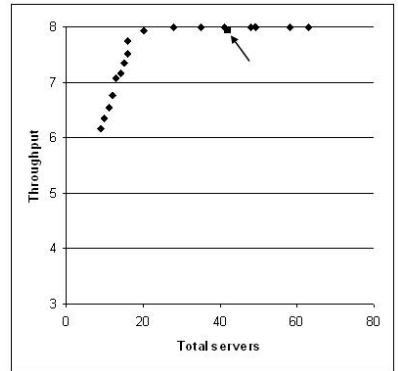
(c) Serial 9 nodes



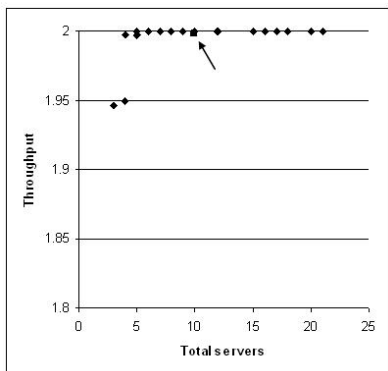
(d) Merge 3 nodes



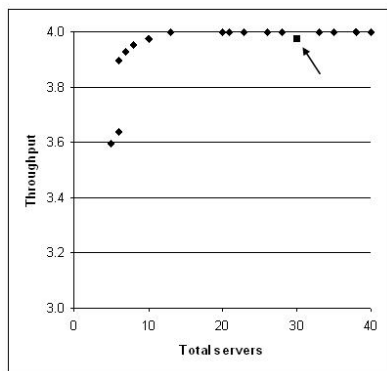
(e) Merge 5 nodes



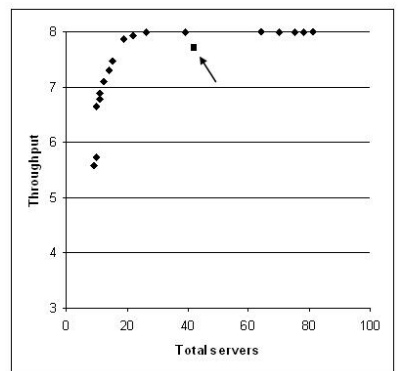
(f) Merge 9 nodes



(g) Split 3 nodes



(h) Split 5 nodes



(i) Split 9 nodes

Fig. 8.

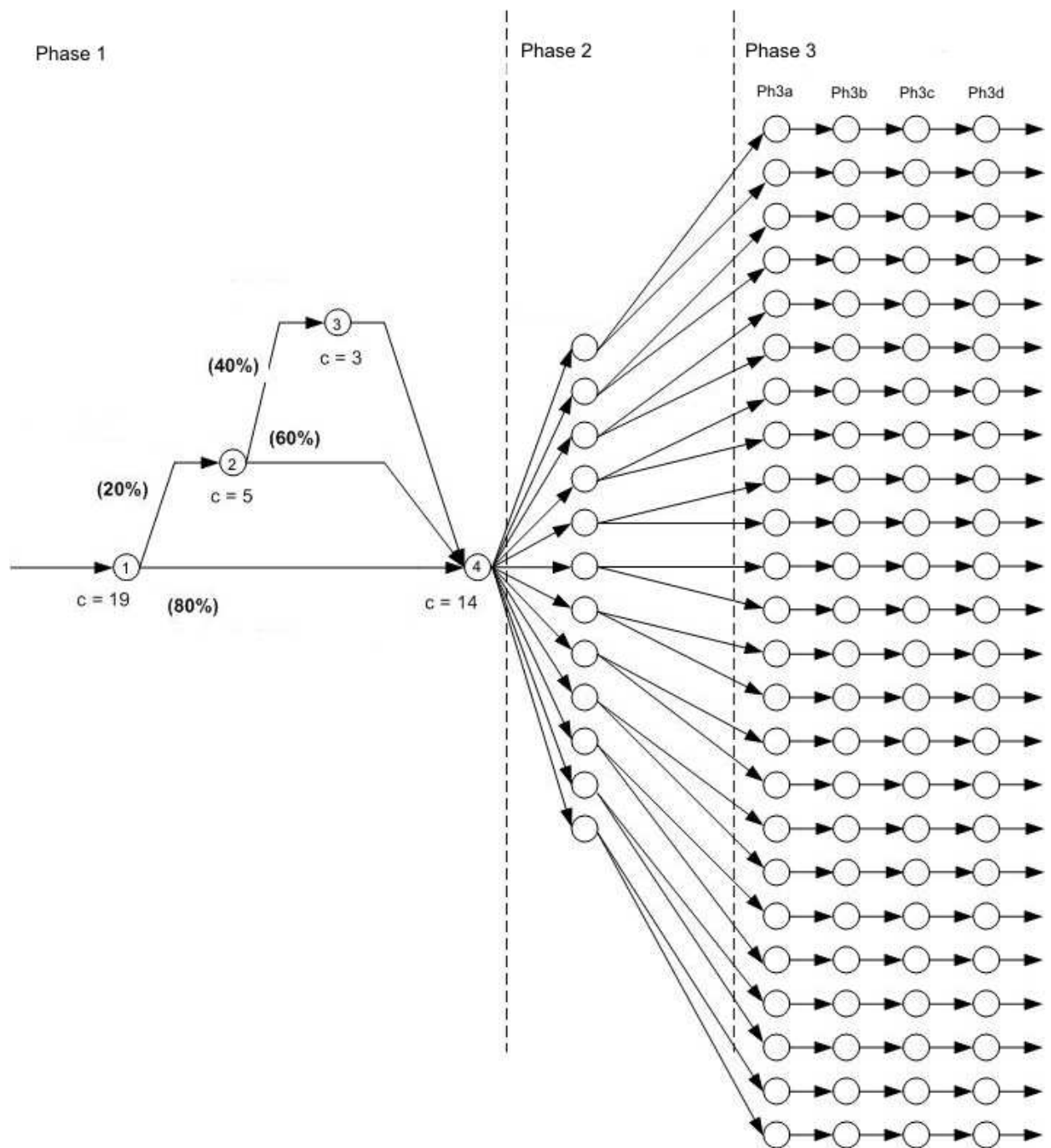


Fig. 9.

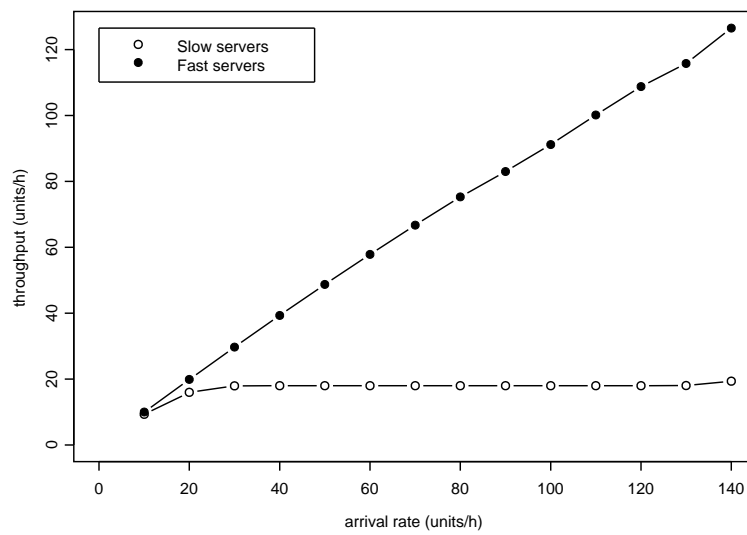


Fig. 10.

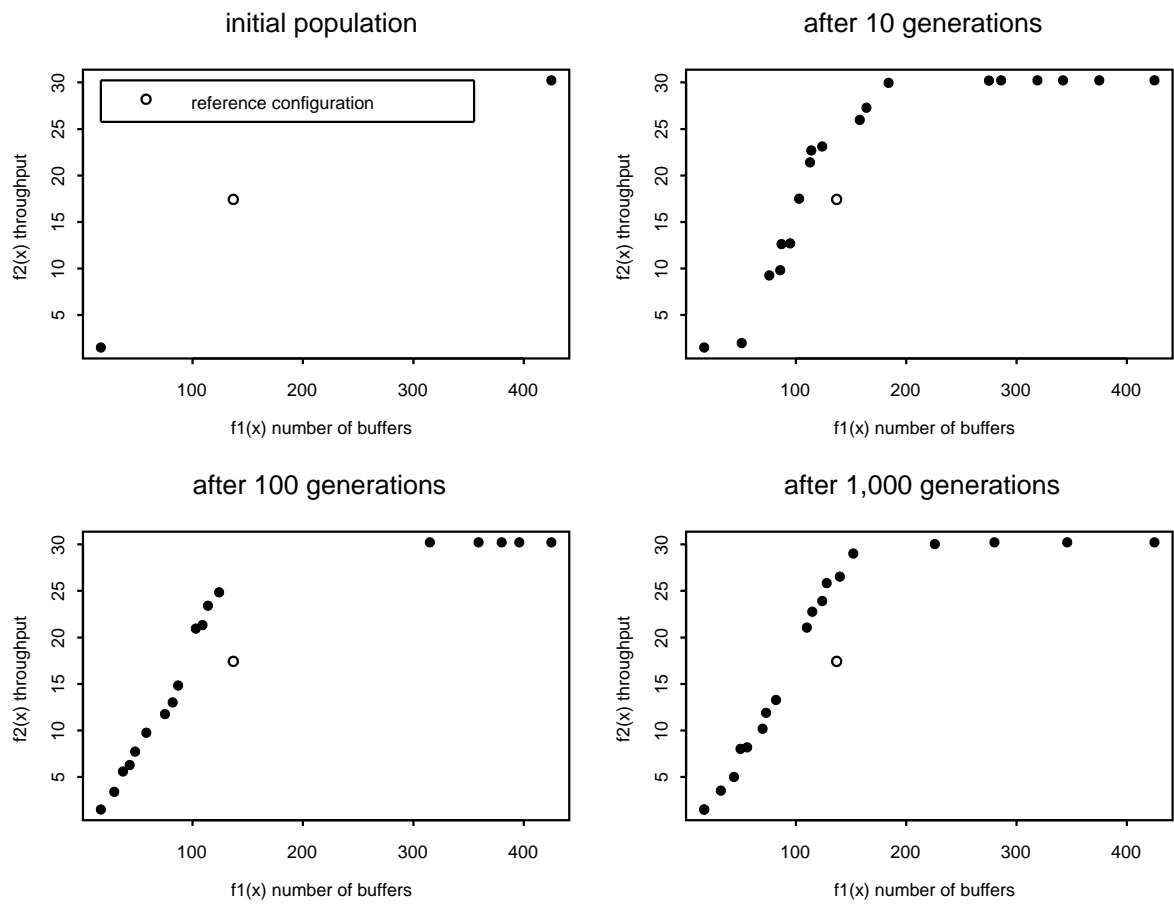


Fig. 11.

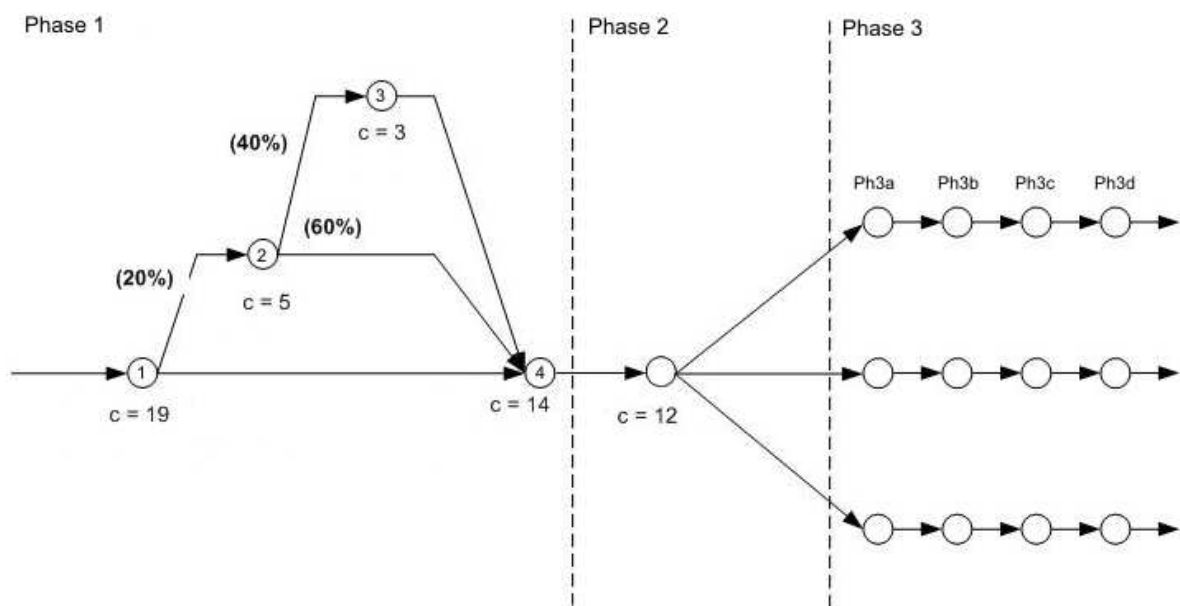


Fig. 12.

Table 1

Variable	Description
h_j	The holding node added to the network, preceding each finite node j
λ_j	The arrival rate to node j
$\tilde{\lambda}_j$	The effective arrival rate to node j
μ_j	The service rate at node j
$\tilde{\mu}_j$	The effective service rate at node j
p_{c_j}	Blocking probability of finite queue of size c_j
p'_{c_j}	Feedback blocking probability in the GEM
μ_{h_j}	The service rate at the holding node h
$c_{D_j}^2$	The squared coefficient of variation of the departure process from node j
$c_{ai,j}^2$	The squared coefficient of variation of the arrival process to the blocked node j
$c_{ah_j}^2$	The squared coefficient of variation of the arrival process at the holding node h_j
$c_{sh_j}^2$	The squared coefficient of variation of the service process at the holding node h_j

Table 2

Series	$N = 3$			$N = 5$			$N = 9$		
\mathbf{c}	(2,2,2)	(4,4,4)	(10,10,10)	(2,...,2)	(4,...,4)	(10,...,10)	(2,...,2)	(4,...,4)	(10,...,10)
$\lambda = 2$	1.996 ^(*)	2.000	2.000	1.996	2.000	2.000	1.995	2.000	2.000
$\lambda = 4$	3.949	4.000	4.000	3.940	4.000 ^(*)	4.000	3.924	4.000	4.000
$\lambda = 8$	7.394	7.988	8.000	7.255	7.988	8.000	7.009	7.987	8.000 ^(*)
$\lambda = 16$	11.49 ^(*)	15.63	16.00	10.676	15.60	16.00	9.572	15.56	16.00
Merge	$N = 3$			$N = 5$			$N = 9$		
\mathbf{c}	(2,2,2)	(4,4,4)	(10,10,10)	(2,...,2)	(4,...,4)	(10,...,10)	(2,...,2)	(4,...,4)	(10,...,10)
$\lambda = 2$	2.000 ^(*)	2.000	2.000	2.000	2.000	2.000	2.000	2.000	2.000
$\lambda = 4$	3.993	4.000	4.000	3.993	4.000 ^(*)	4.000	3.992	4.000	4.000
$\lambda = 8$	7.900	7.999	8.000	7.892	7.999	8.000	7.875	7.999	8.000 ^(*)
$\lambda = 16$	14.52	15.98	16.00	14.40	15.98 ^(*)	16.00	14.16	15.98	16.00
Split	$N = 3$			$N = 5$			$N = 9$		
\mathbf{c}	(2,2,2)	(4,4,4)	(10,10,10)	(2,...,2)	(4,...,4)	(10,...,10)	(2,...,2)	(4,...,4)	(10,...,10)
$\lambda = 2$	1.997 ^(*)	2.000	2.000	1.997	2.000	2.000	1.997	2.000	2.000
$\lambda = 4$	3.957	4.000	4.000	3.957	4.000 ^(*)	4.000	3.956	4.000	4.000
$\lambda = 8$	7.538	7.988	8.000	7.531	7.988	8.000	7.517	7.988	8.000 ^(*)
$\lambda = 16$	12.59	15.65	16.00	12.52	15.65	16.00	12.38	15.65	16.00 ^(*)

(*)Earmarked experiments confirmed by simulation (see Table 3).

Table 3

Topology	λ	\mathbf{c}	Analytical		Simulation		
			θ	θ^s	δ	CPU(min)	$\Delta\% \theta$
Series	2.0	(2,2,2)	1.996	1.966	0.001	2.0	1.51%
	4.0	(4,4,4,4,4)	4.000	3.999	0.002	5.2	0.03%
	8.0	(10,...,10)	8.000	8.002	0.003	36	-0.02%
	16.0	(2,2,2)	11.49	11.32	0.002	22	1.50%
Merge	2.0	(2,2,2)	2.000	1.991	0.002	1.9	0.44%
	4.0	(4,4,4,4,4)	4.000	3.999	0.002	5.8	0.03%
	8.0	(10,...,10)	8.000	7.998	0.003	20	0.03%
	16.0	(4,4,4,4,4)	15.98	15.88	0.003	24	0.63%
Split	2.0	(2,2,2)	1.997	1.967	0.002	2.1	1.55%
	4.0	(4,4,4,4,4)	4.000	3.996	0.002	5.7	0.09%
	8.0	(10,...,10)	8.000	7.999	0.003	25	0.02%
	16.0	(10,...,10)	16.00	16.00	0.004	150	0.00%

Table 4

Series	$N = 3$	$N = 5$	$N = 9$
\mathbf{c}	(2,4,10)	(2,4,10,2,4)	(2,4,10,2,4,10,2,4,10)
$\boldsymbol{\mu}$	(12,11,10)	(12,11,10,12,11)	(12,11,10,12,11,10,12,11,10)
$\lambda = 2$	1.998 ^(*)	1.998	1.998
$\lambda = 4$	3.974	3.974 ^(*)	3.974
$\lambda = 8$	7.698	7.697	7.696 ^(*)
$\lambda = 16$	13.50 ^(*)	13.47	13.45
Merge	$N = 3$	$N = 5$	$N = 9$
\mathbf{c}	(4,4,2)	(10,10,4,4,2)	(2,2,4,4,10,10,4,4,2)
$\boldsymbol{\mu}$	(11,11,12)	(10,10,11,11,12)	(12,12,11,11,10,10,11,11,12)
$\lambda = 2$	2.000 ^(*)	2.000	2.000
$\lambda = 4$	4.000	4.000 ^(*)	3.996
$\lambda = 8$	8.000	8.000	7.947 ^(*)
$\lambda = 16$	15.92	15.93 ^(*)	15.35
Split	$N = 3$	$N = 5$	$N = 9$
\mathbf{c}	(2,4,4)	(2,4,4,10,10)	(2,4,4,10,10,4,4,2,2)
$\boldsymbol{\mu}$	(12,11,11)	(12,11,11,10,10)	(12,11,11,10,10,11,11,12,12)
$\lambda = 2$	1.998 ^(*)	1.998	1.998
$\lambda = 4$	3.974	3.974 ^(*)	3.974
$\lambda = 8$	7.698	7.698	7.698 ^(*)
$\lambda = 16$	13.51	13.51	13.51 ^(*)

(*)Earmarked experiments confirmed by simulation (see Table 5).

Table 5

Topology	λ	\mathbf{c}	$\boldsymbol{\mu}$	Analytical	Simulation			$\Delta\% \theta$
				θ	θ^s	δ	CPU(min)	
Series	2.0	(2,4,10)	(12,11,10)	1.998	1.977	0.001	3.0	1.07%
	4.0	(2,4,10,2,4)	(12,11,10,12,11)	3.974	3.839	0.002	9.9	3.52%
	8.0	(2,4,10,2,4,10,2,4,10)	(12,11,10,12,11,10,12,11,10)	7.969	7.759	0.002	35	2.70%
	16.0	(2,4,10)	(12,11,10)	13.50	13.50	0.002	43	0.00%
Merge	2.0	(4,4,2)	(11,11,12)	2.000	2.001	0.002	1.9	-0.04%
	4.0	(10,10,4,4,2)	(10,10,11,11,12)	4.000	3.999	0.002	5.8	0.02%
	8.0	(2,2,4,4,10,10,4,4,2,2)	(12,12,11,11,10,10,11,11,12)	7.947	7.652	0.003	20	3.86%
	16.0	(10,10,4,4,2)	(10,10,11,11,12)	15.93	16.00	0.005	29	-0.44%
Split	2.0	(2,4,4)	(12,11,11)	1.998	1.973	0.001	2.1	1.27%
	4.0	(2,4,4,10,10)	(12,11,11,10,10)	3.974	3.980	0.002	6.1	-0.15%
	8.0	(2,4,4,10,10,4,4,2,2)	(12,11,11,10,10,11,11,12,12)	7.698	7.654	0.002	22	0.57%
	16.0	(2,4,4,10,10,4,4,2,2)	(12,11,11,10,10,11,11,12,12)	13.51	12.88	0.030	38	4.89%

Table 6

Server	Number of servers	Slow servers (units/h)	Fast servers (units/h)
Phase 1 server 1	19	5.833	15.42
Phase 1 server 2	5	3.400	11.05
Phase 1 server 3	3	3.400	11.05
Phase 1 server 4	14	3.400	20.00
Phase 2 servers	12	2.000	26.25
Ph3a servers	21	7.500	23.63
Ph3b servers	21	2.213	22.13
Ph3c servers	21	2.213	22.13
Ph3d servers	21	2.213	22.13

Table 7

	λ	Analytical	Simulation			$\Delta\% \theta$
		θ	θ^s	δ	CPU (min)	
Slow servers	30.215	17.422	18.464	0.02	120	-5.6%
Fast servers	30.215	29.907	30.213	0.006	70.1	-1.0%

Table 8

Server	Configuration #1 (units/h)	Configuration #2 (units/h)	Configuration #3 (units/h)
Phase 1 server 1	5.833	5.833	5.833
Phase 1 server 2	3.400	3.400	3.400
Phase 1 server 3	3.400	3.400	3.400
Phase 1 server 4	3.400	3.400	3.400
Phase 2 servers	2.000	6.000	6.000
Ph3a servers	7.500	7.500	7.500
Ph3b servers	2.213	2.213	3.320
Ph3c servers	2.213	2.213	3.320
Ph3d servers	2.213	3.320	3.320
Throughput	17.422	18.480	20.124
Improvement		6.1%	16%