

Efficient algorithms for buffer space allocation

Stanley B. Gershwin^{a,*} and James E. Schor^b

^a *Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA 02139-4307, USA*
E-mail: gershwin@mit.edu

^b *Analytics, Inc., 101 Rogers Street, Cambridge, MA 02142, USA*
E-mail: jschor@analytics-consulting.com

This paper describes efficient algorithms for determining how buffer space should be allocated in a flow line. We analyze two problems: a *primal* problem, which minimizes total buffer space subject to a production rate constraint; and a *dual* problem, which maximizes production rate subject to a total buffer space constraint. The dual problem is solved by means of a gradient method, and the primal problem is solved using the dual solution. Numerical results are presented. Profit optimization problems are natural generalizations of the primal and dual problems, and we show how they can be solved using essentially the same algorithms.

1. Introduction

1.1. Problem description

Production systems are often organized with machines or work centers connected in series and separated by buffers. This arrangement is often called a *flow line*, or *transfer line*, or *production line*. A five-machine line is represented in figure 1, in which the squares represent machines and the circles represent buffers. Material moves in the direction of the arrows, from upstream inventory to the first machine for an operation, to the first buffer where it waits for the second machine, to the second machine, etc.

Material flow may be disrupted by machine failures or variable processing times. Buffers are inserted between machines to limit the propagation of disruptions, and this increases the average production rate of the line. Inclusion of buffers requires additional capital investment and floor space, which may be expensive. Buffering also increases in-process inventory. If the buffers are too large, the work-in-process inventory and capital costs incurred will outweigh the benefit of increased productivity. If the buffers are too small, the machines will be underutilized or demand will not be met.

Often, the limitation on the amount of in-process inventory is not due to physical constraints. Instead this limitation is a control policy, for example, a version of the

* Corresponding author.

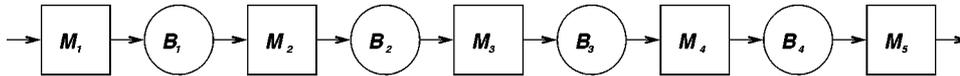


Figure 1. Flow line.

kanban policy [38]. It is essential to determine buffer sizes to achieve the desired performance.

This paper describes efficient algorithms for buffer space allocation. Much literature analyzes the effects of machine unreliability and buffer sizes on system performance. Less effort has been expended to develop methods for optimizing or improving the allocation of buffer space. This paper is a summary of Schor [44], in which can be found all details of algorithms and numerical results.

1.2. Approach and outline

This paper provides a unified view of solutions to several versions of the problem. They are organized according to problem type and model type.

- **Problem type.** We describe algorithms which choose buffers sizes for a flow line. They are based on the qualitative properties of the production rate described in section 2. The number of machines and the reliability and speed parameters of each machine are assumed to be specified.
 - * **Primal.** The goal of the *primal* is to minimize the total buffer space required for the line to meet or exceed a given average production rate. This formulation is appropriate if either floor space or the buffering mechanism is expensive, if work-in-process inventory is inexpensive, and if an average production rate is mandated.
 - * **Dual.** The goal of the *dual* is to maximize the production rate achievable with a given total buffer space. This is appropriate in cases where the total floor space is fixed, where the number of buffer locations is fixed, and where the problem is how to space the machines to get the maximum benefit from the buffers. In addition, the dual is used in this paper as a subproblem in solving the primal.
 - * **Constrained profit maximization.** The goal is to maximize profit subject to a constraint on total buffer space. The solution technique is similar to that of the dual problem.
 - * **Unconstrained profit maximization.** The buffer space constraint is relaxed. The solution technique is similar to that of the primal problem. It uses the constrained profit maximization problem in the same way that the primal algorithm uses the dual.
- **Model type.** We consider systems that may be represented by discrete or continuous material flow models. In both, the processing time is deterministic. The discrete material model has the advantage of better representing the discrete nature of typical factories, but it suffers from two important drawbacks: it is restricted to systems

with equal processing times, and the buffer sizes (which are the decision variables in the allocation problems) are integers. The continuous model is better suited to optimization because its variables are real numbers; it can be used for discrete part factory design by rounding the buffer sizes after the optimization has been completed.

- * **Discrete material.** In the discrete material, deterministic processing time system, the machines have equal processing times and geometrically distributed failure and repair times. The probability of Machine M_i failing during an operation time while it is operating (operational and not starved or blocked) is p_i . The probability of a repair during an operation time while it is under repair is r_i . The size of Buffer B_i is N_i . N_i is a nonnegative integer; p_i and r_i are real numbers in $[0,1]$. See [7,18]. We use an approximate decomposition method, developed by Gershwin [17], to determine production rate and average buffer levels. The DDX algorithm [12], solves the decomposition equations efficiently and produces estimates for production rate and average buffer levels.
- * **Continuous material.** The machines have different processing rates and exponentially distributed failure and repair times. The amount of material produced during $(t, t + \delta t)$ by Machine M_i while it is fully operational (operational, not starved or blocked, and not slowed down by an adjacent machine) is $\mu_i \delta t$. The probability of Machine M_i failing during such an interval is $p_i \delta t$. The probability of a repair during $(t, t + \delta t)$ while it is under repair is $r_i \delta t$. The size of Buffer B_i is N_i . μ_i , p_i , r_i , and N_i are nonnegative real numbers. The decomposition for the continuous model, and the ADDX algorithm, with which we evaluate its performance measures, were developed by Burman [6]. The ADDX algorithm approximately determines the production rate and average in-process inventory for this kind of system, by a decomposition method which is an extension of the Gershwin [17]–Dallery–David–Xie [12] method.

Numerical results are provided for selected cases.

Qualitative properties of flow line models are discussed in section 2. The primal and dual problem formulations are presented in section 3. The primal and dual problem solution methods are summarized in sections 4 and 5. The behavior of the algorithms is discussed in section 6. We reformulate the problems to maximize profit (and therefore to include inventory costs) in section 7. We conclude in section 8.

1.3. *Prior literature*

There is a substantial literature on the analysis of transfer or flow lines [13]. This literature is mainly concerned with the prediction of performance. Much of it is aimed at evaluating the production rate of a system with specified machines and buffers. One ultimate goal of the literature is to aid the designers of such systems to make efficient choices of machines and buffers. To do this, two elements are needed: a fast, accurate method of evaluating systems, and an efficient way of selecting systems to evaluate.

Fast evaluation methods exist, but as Martin [36] says, “. . . the more practical goal of optimizing system design is one less visibly pursued in the literature.”

Here, we do not review the literature on evaluation; we concentrate instead on flow line optimization papers and on papers about relevant qualitative properties of production lines. The main contribution of this paper is a set of algorithms that efficiently select the minimal buffer space in a flow line to achieve a specified production rate. The algorithms are based on analytical approximations of the Buzacott model of a flow line, so they are appropriate for automated systems, and they are fast. By contrast, some papers in the literature are based on simulation; others use different analytic models; and nearly all have other objective functions.

1.3.1. *Qualitative properties*

In order for gradient methods to work well, certain qualitative properties are needed. Although the independent variables often are integers, there should be some kind of *continuity*: a small change in a buffer size should lead to a small change in the system’s performance. (Many line optimization papers assume this implicitly.) In addition, many methods rely on *monotonicity*: an increase in a buffer’s size (while all the other buffer sizes are increased or held constant) increases production rate. Finally, for a gradient method to find a unique optimum, it is helpful for some form of *concavity* to hold. That is, the increase in production rate due to a unit increase in buffer size decreases as the buffer size increases. (It is interesting to observe how much of the work on line optimization was done *before* the publication of the papers described in this section.) Combinatorial approaches, which do not exploit these properties, can be expected to be relatively inefficient.

Several papers have studied qualitative properties of models of manufacturing systems that are relevant to optimization. The properties have been proved for classes of systems other than those we study here, but we believe that the systems are similar enough to justify assuming the properties we need. Numerical experiments support this assumption.

Shanthikumar and Yao [48] studied cyclic queuing networks with finite buffers. Service processes were exponential, with rates that are increasing functions of the number of customers in the queue. They showed that the production rate is an increasing function of the buffer sizes. Adan and Van der Wal [1] proved monotonicity of the production rate with respect to buffer sizes of acyclic systems using a sample path method. That is, they also showed that the production rate increases as each buffer is enlarged. They prove this for a more general class of networks than the one studied here.

Meester and Shanthikumar [37] studied tandem queueing systems with reliable exponential servers and finite storage areas. They showed that the production rate is an increasing, concave function of the buffer sizes. Anantharam and Tsoucas [3] studied systems of queues in series with single exponential servers, finite buffers, and communication blocking (blocking before service). They showed that the throughput (production rate) is a concave function of the buffer sizes. Glasserman and Yao [22]

studied similar systems under a wide variety of blocking mechanisms, and demonstrated concavity for systems with exponential processing times, and other qualitative properties for general service time distributions. Dallery et al. [14] generalized the results of Shanthikumar and Yao [48] and Meester and Shanthikumar [37] by establishing the concavity of throughput as a function of buffer sizes (and initial conditions, for a closed network) for a greatly extended class of distributions of service times. Useful properties for a general class of related systems were studied by Tayur [52] and Muckstadt and Tayur [38]. Buzacott and Shanthikumar [9] summarized much of the relevant literature on manufacturing system design.

1.3.2. *Simulation-based approaches*

Many papers have proposed design methodologies based on evaluation by simulation. Because simulation takes longer than analytical methods on the same problem, even the most recent of such methods are slower than analytical methods. The main advantage of simulation is that it has the potential for dealing with a larger class of systems than analytical methods.

Barten [5] simulated a set of cases, and interpolated a curve fit for delay time as a function of storage capacity for several line lengths, assuming normally distributed operation times. He formed a cost function, and calculated the optimal storage capacity. Freeman [16] simulated short lines (two and three stages) with unreliable machines. Even under these limited conditions, he was able to show that the placement and sizes of buffers is important, and he provided some preliminary rules about how to allocate storage space.

Anderson and Moodie [4] simulated two classes of production lines with reliable machines: one class had identical, normally distributed service times; the other had identical, exponentially distributed service times. All the buffers in each line had identical sizes. In each case, average in-process inventory and *delay* (evidently 1.0-utilization) were fitted to simple functions of the length of the lines and the buffer sizes. Cost functions were also developed, and optimal buffer sizes were calculated from them. Martin [35,36] extended this work by refining the delay and inventory functions and generalizing the cost and profit functions. He studied the performance of the line as a function of line length and buffer size.

Powell [42] and Powell and Pyke [43] simulated flow lines to determine the optimal locations and amounts of buffering. In their systems, operation times are random but machines or workstations do not fail. As a consequence, the buffer sizes they found were small. They also suggested some rules of thumb for buffer sizing and placement.

Chow [11] and Liu and Lin [34] used simulation to construct functions to predict the throughput and the CV of the interdeparture time of a reliable two-machine line. Liu and Lin [34] improved Chow's approach by designing a procedure that required fewer runs to construct these functions. They also proposed an aggregation method, which represented the long line as a single two-machine line, to extend their results to longer lines. They used this in a dynamic programming formulation to solve the production

rate maximization problem. They solve the minimum total buffer space problem by using the solution to the production rate maximization problem. Starting with buffers of size zero, they evaluated the maximum production rate for each value of total space. They increased the space by one and repeated until the throughput constraint was satisfied. They performed experiments on lines of less than ten machines with a total buffer space of less than fifty. When evaluating longer lines they only considered bottleneck machines.

Another important class of simulation-based papers optimize lines by means of *perturbation analysis*. Ho et al. [26] was the first of these. In this paper, a Buzacott-like model of a transfer line (one with discrete parts, identical constant processing times, and geometrically distributed repair and failure times) was simulated. By means of a sophisticated analysis, an estimate of the gradient of the production rate with respect to all the buffer sizes was determined from only one run of the simulation. This estimate was then used in a steepest descent method for solving the dual problem (see also Ho et al. [27]). Caramanis [10] applied a similar technique to a line with deterministic but different processing times (and exponential repair and failure times). The cost function depended on buffer capacity, average inventory, and production rate. A more recent paper, which studied the optimization of production rate in a continuous material model of a flow line (subject to linear constraints on the line parameters) by similar techniques is Plambeck et al. [41].

1.3.3. *Exploration of the optimum by exact numerical evaluation*

Hillier and So [24] studied systems of unreliable machines and finite buffers. They represented the systems as finite queues in series with two-stage Coxian distributed operation times. Using an exact numerical method for short lines (5 stations or less) and small buffers (4 or less), they evaluated all plausible optimal distributions of buffer sizes in order to find the optimum. Hillier et al. [25] studied a problem which is similar to the dual described here, the maximization of the production rate subject to a total buffer space constraint. As is typical of the papers in the series by Hillier and his co-authors, they evaluated a large number of cases to find the optimum, and characterized the optimal distribution. This methodology requires small, simple lines: their machines were reliable with exponential processing times, their lines had no more than eight stations, and their buffers were no larger than 14. When all machines were the same, they observed an *inverted bowl phenomenon*, in which the optimal distribution has smaller buffers near the ends of the line, and larger buffers in the interior.

The main goal of these papers has been to develop an understanding of the structure of optimal systems. By contrast, our goal is to develop efficient methods for obtaining an optimum (or a near-optimum). In addition, these papers tend to deal with models of reliable work centers with processing times (which may be appropriate for manual systems). This paper deals with deterministic processing time models which are appropriate for automated systems with fixed processing times and random failures.

1.3.4. Other algorithms and approaches

Early papers. The early papers on buffer allocation in flow lines made various simplifications or restrictions in order to obtain results. Small systems are easier to analyze than large systems, and reliable machines with exponential processing times are easier to analyze than machines with more realistic behavior.

One of the earliest papers in the production line literature was that of Sevast'yanov [47]. This paper included an exact solution of a two-machine line, and an approximate analysis of longer lines. (Some ideas that have been implemented in the last fifteen years were anticipated by this paper.) It also contained formulations of storage size optimization problems that are similar to those here. Numerical results were not included, no doubt due to the unavailability of computers, but theoretical results demonstrated that buffer sizes should be chosen so that the apparent upstream failure rate should equal the apparent downstream failure rate, for an observer in each buffer.

Hatcher [23] studied reliable systems with exponential processing times and finite buffers, and investigated the allocation of capacity in three-stage, two-buffer systems. (However, there is an error in the evaluation of the line. See Knott [31].) Kraemer and Love [32] optimize the profit in a two-stage, one-buffer system with reliable, exponential machines. Profit is revenue minus inventory carrying cost minus storage facility cost; revenue is proportional to production rate; inventory carrying cost is proportional to expected in-process inventory; and storage facility cost is proportional to the buffer size.

Sheskin [49] studied the allocation of buffer space in systems like ours: those with unreliable machines with equal, deterministic processing times, and finite buffers. However, he restricted his analysis to the case where $r_i + p_i = 1$ (in our notation), which is equivalent to assuming that the machine repair state at time $t+1$ is independent of the state at time t . In addition, he assumed time-dependent failures, i.e., that machines could fail even if they were starved or blocked. A decomposition method was used to produce numerical results for small systems with small buffers. These results led to some rules of thumb on the allocation of buffer space to maximize production rate. Soyster et al. [51] used the same model to study the maximization of production rate subject to general linear inequality constraints on buffer sizes. They approximated the production rate for small systems and used an integer programming package to find optimal allocations of buffer space.

Altioek and Stidham [2] studied systems with exponentially distributed service, repair, and failure times. Their goal was to maximize average profit, which increased with production rate and decreased with total average inventory. There were no constraints on total buffer space. They evaluated the steady state distribution by numerically solving the transition equations, so they were limited to small systems.

More recent papers. Chow [11] developed an aggregation procedure for evaluating production lines. He developed a dynamic programming procedure to maximize production rate subject to a total buffer space constraint.

Smith and Daskalaki [50] developed a method for general networks of reliable exponential machines with random routing and finite buffers. They formed an unconstrained optimization problem. In their profit objective, production rate had a positive coefficient, and average buffer level had a negative coefficient. A standard optimization method was used. By contrast, we deal with unreliable machines with constant processing time and fixed routing in a tandem network, and we develop a special purpose optimization method.

Jafari and Shanthikumar [30] studied a problem similar to our dual: maximizing the production rate subject to the total storage space constraint. They included scrapping in their model, and they evaluated their lines in a manner similar to the method used here. They studied two optimization methods: a dynamic programming algorithm, and a greedy heuristic.

Park [40] developed a two-phase tree search branch and bound technique to solve some of the same problems that we investigate here. He restricted his analysis to the same discrete material model that we study. He identified four kinds of problems: problem 1 is the same as our primal; problem 2 is a generalization of our dual (in which the buffer sizes may be weighted with different coefficients, and in which there may be more than one linear constraint); problem 3 seeks to minimize the total buffer space subject to the production rate requirement of problem 1 (and our primal) and to the generalized buffer limits of problem 2. Problem 4 seems to have two objectives: to simultaneously maximize production rate while minimizing total buffer space. Park also included a survey of earlier buffer allocation work.

Jacobs and Meerkov [28,29] studied the *improvability* of production line systems. While their definition of “improvable” appears equivalent to “non-optimal”, their goal was not the design of systems. Rather, they sought methods for efficiently improving the performance of existing systems using data that is available as the systems operate. Jacobs and Meerkov [28,29] employ the concept of *improvability* to determine how buffer space should be allocated. They propose a method for determining whether the individual buffer sizes may be changed, while holding the total buffer space constant, so that the production rate will increase. The flow line model they study is similar to that of Sheskin [49] and Soyster et al. [51]. Kuo et al. [33] described an application.

Seong et al. [45] studied our dual problem: they maximized the production rate subject to a constraint on the total buffer size. (This is our dual, and Park’s problem 2.) In Seong et al. [46], they studied the same objective function as Altiok and Stidham [2] with general linear constraints on buffer sizes. Material was assumed continuous, machine operation speeds were deterministic and equal, and repairs and failures were exponentially distributed. They solved the problem with a gradient projection algorithm. They used a similar approximation of the gradient of the production rate as we do below. Seong et al. [45] used a gradient method to solve the production rate maximization problem for a flow line with exponentially distributed failure, repair, and processing times. Seong et al. [46] also solved this problem and the profit maximization problem for a specified total buffer space for a continuous flow line. We compare our results with theirs in section 7. Dogan and Altiok [15] use decomposition

to estimate the throughput gradient vector with respect to repair rates. Their ultimate goal is to optimize “repair rate allocation”. Ahtiok and Yamashita [53] observed a similar primal/dual relationship as ours, and solved the dual problem using a dynamic programming approach.

Gershwin and Goldis [20] employed a gradient method to solve the primal problem. Their algorithm was based on the observation that if the production rate is expanded to first order the problem may be formulated as an integer linear program. They guaranteed that their solutions are optimal or near-optimal. The algorithm presented here converges much more quickly in the majority of the cases we study.

2. Definitions and properties

2.1. Models, parameters, and notation

In the following, N_i , the size of Buffer B_i , is the decision variable, and P is the production rate of the line. The line has k machines and $k - 1$ buffers. In the discrete material formulation, all machines are assumed to have the same operation time, which we use as the time unit. The parameters of Machine M_i are p_i , the probability of a failure during a time unit while the machine is operating; and r_i , the probability of a repair during a time unit while the machine is down.

In the continuous material formulation, workpieces are treated like a continuous fluid. This is one way of modeling machines that have different speeds. The rate at which Machine M_i processes material, when it is operational and neither starved nor blocked is μ_i . That is, it processes $\mu_i \delta t$ units of material in a time interval of length δt . The reliability parameters of Machine M_i are p_i and r_i . They are probability rates: $p_i \delta t$ is the probability of a failure during a time period of length δt during which the machine is operating and neither starved nor blocked; and $r_i \delta t$ is the probability of a repair during a δt time period during which the machine is down. Details of both models can be found in Gershwin [18].

Although the production rate P is a function of machine speeds and reliabilities (and, in general, other attributes not treated here), we vary only buffer sizes, so we write $P = P(N_1, \dots, N_{k-1})$.

2.2. Machines

Isolated quantities are quantities of a machine that are independent of the rest of the line. The *isolated efficiency* of Machine M_i is given by

$$e_i = \frac{r_i}{r_i + p_i}.$$

In the discrete material model, the operation time is 1 and the production rate is the average number of parts produced in one time unit. Therefore, the *isolated production rate* is the same as the isolated efficiency.

In the continuous material model, the production rate of Machine M_i is μ_i when it is operating and not influenced by other machines or buffers. Therefore, its isolated production rate is given by $\rho_i = \mu_i e_i$.

2.3. Lines

2.3.1. Quantitative properties

The production rate P when all buffer sizes are infinite is the minimum of the isolated production rates of all the machines in the line. That is, in the discrete material case,

$$P(\infty, \dots, \infty) = \min_{i=1, \dots, k} e_i.$$

In the continuous material case,

$$P(\infty, \dots, \infty) = \min_{i=1, \dots, k} \rho_i.$$

In the deterministic processing time model, the production rate when all buffers have size 0 is [7]

$$P(0, \dots, 0) = \frac{1}{1 + \sum_{i=1}^k p_i/r_i}. \quad (1)$$

A similar equation is available for the continuous material case. It is more complex when the μ_i are not the same.

When there are more than two machines and all N_i are neither infinite nor zero, the production rate and average inventory levels cannot be calculated analytically or exactly even numerically. Decomposition methods [6,17,18] and iterative algorithms [6,12,18] must be used. These algorithms work by calculating $r_u(i)$, $p_u(i)$, $r_d(i)$, and $p_d(i)$ (and also $\mu_u(i)$ and $\mu_d(i)$ for the continuous material case), the parameters of fictitious machines that represent the behavior of the flow of material into and out of Buffer B_i . As indicated in section 5.4, these methods should be adapted for use in the algorithms described here.

2.3.2. Qualitative properties

We assume the following qualitative properties:

- **Continuity.** A small change in any N_i creates a small change in P .
- **Monotonicity.** The production rate increases monotonically in each N_i .
- **Concavity.** The production rate appears to be a concave function of the vector (N_1, \dots, N_{k-1}) .

These properties are evident in figure 2, which shows how production rate varies with buffer sizes for the three-machine line whose parameters are given in table 1. Figure 3 shows the curves of constant P for this case. (The *optimal curve* is described in section 5.2.) The literature provides proofs of these properties for similar systems.

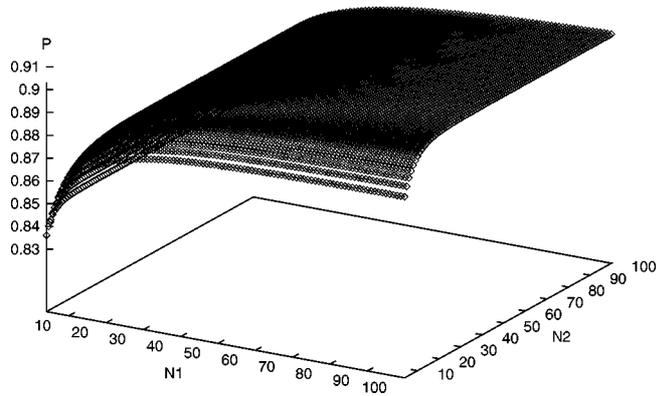


Figure 2. P vs. N_1 and N_2 .

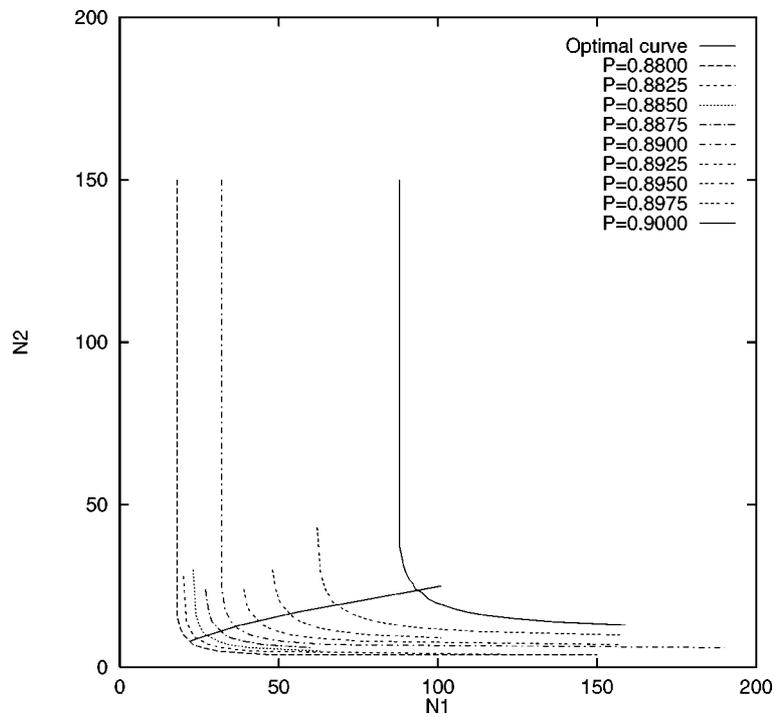


Figure 3. Iso- P lines.

Table 1
Three-machine line parameters.

Machine	r_i	p_i
1	0.35	0.037
2	0.15	0.015
3	0.4	0.02

3. Primal and dual problem

In all of these problems, the buffer sizes have lower bounds. In the discrete material formulation, it is convenient to restrict N_i to be greater than or equal to 4. This is because the functional forms of the production rate and average buffer levels are different for $N_i < 4$ than for $N_i \geq 4$. We therefore write $N_i \geq N^{\text{MIN}}$, where $N^{\text{MIN}} = 4$ in the discrete material case and $N^{\text{MIN}} = 0$ in the continuous material case.

3.1. Primal problem

In the primal problem, we seek the buffer sizes (N_1, \dots, N_{k-1}) that minimize total buffer space N^{TOTAL} such that the production rate P is greater than or equal to a specified value P^* . Or,

$$\text{Minimize } N^{\text{TOTAL}} = \sum_{i=1}^{k-1} N_i \quad (2)$$

$$\text{subject to } P(N_1, \dots, N_{k-1}) \geq P^*; \quad P^* \text{ specified,} \quad (3)$$

$$N_i \geq N^{\text{MIN}}, \quad i = 1, \dots, k-1.$$

This problem is difficult because the constraint $P \geq P^*$ cannot be expressed in closed form. Consequently, even if we know a value of (N_1, \dots, N_{k-1}) satisfying $P(N_1, \dots, N_{k-1}) = P^*$, it is difficult to construct another (N_1, \dots, N_{k-1}) that will also be on that surface.

The qualitative properties ensure that a solution will be found with constraint (3) satisfied with equality. Since the properties appear to hold strictly, there will be a unique solution.

3.2. Dual problem

In the dual problem, we seek the buffer sizes $N = (N_1, \dots, N_{k-1})$ that maximize the production rate P such that the total buffer space is equal to a specified value N^{TOTAL} . That is,

$$\text{Maximize } P(N_1, \dots, N_{k-1}) \quad (4)$$

$$\text{subject to } N^{\text{TOTAL}} = \sum_{i=1}^{k-1} N_i; \quad N^{\text{TOTAL}} \text{ specified,} \quad (5)$$

$$N_i \geq N^{\text{MIN}}, \quad i = 1, \dots, k-1.$$

The input to this problem is N^{TOTAL} and the outputs are the optimal P_{MAX} and (N_1, \dots, N_{k-1}) . The resulting maximum value of P is written $P_{\text{MAX}}(N^{\text{TOTAL}})$.

4. Dual algorithm

The dual algorithm is a gradient algorithm. An initial guess \mathbf{N} is selected and a direction to move in (N_1, \dots, N_{k-1}) space is determined. A linear search is then conducted in that direction until a maximum is encountered. This becomes the next guess. A new direction is chosen and the process continues until no improvement is realized.

The dual problem is relatively easy to solve, and the dual algorithm is relatively easy to implement, because the constraint set (5) is a plane. As a result, large steps may be taken.

To determine the search direction, we find the gradient \mathbf{g} according to a forward or backward difference formula. The forward difference is

$$g_i = \frac{P(N_1, \dots, N_i + \delta N, \dots, N_{k-1}) - P(N_1, \dots, N_i, \dots, N_{k-1})}{\delta N}. \quad (6)$$

Monotonicity guarantees that $g_i \geq 0$ for all i .

We must select $\delta \mathbf{N} = \widehat{\mathbf{N}} - \mathbf{N}$ (where \mathbf{N} is the most recent guess of the dual solution and $\widehat{\mathbf{N}}$ is the next guess) to achieve the greatest increase in P and to continue to satisfy constraint (5). We construct the search direction \mathbf{p} by projecting \mathbf{g} onto the

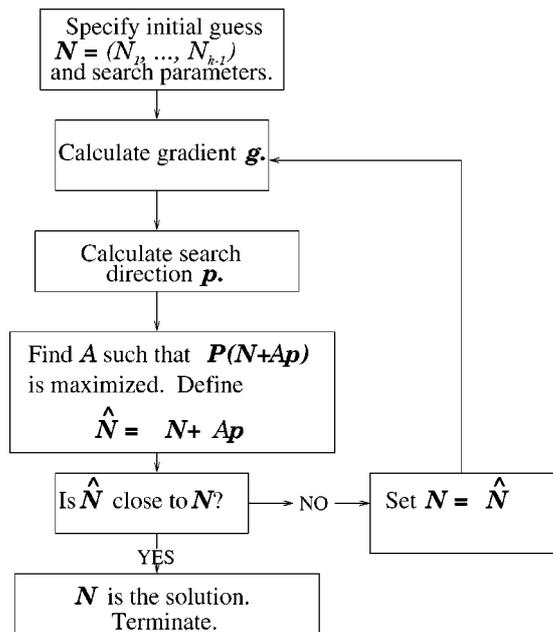


Figure 4. Block diagram of dual algorithm.

hyperplane $\sum_{i=1}^{k-1} \delta N_i = 0$ as follows:

$$\bar{g} = \frac{1}{k-1} \sum_{i=1}^{k-1} g_i; \quad p_i = g_i - \bar{g}. \quad (7)$$

An exception occurs when $N_q = N^{\text{MIN}}$ and $p_q \leq 0$ for some q . We calculate the new direction by deleting the q -th component of \mathbf{g} and setting $p_q = 0$. Then we recalculate the other components of \mathbf{p} using (7), and the new \mathbf{g} . This process is continued until a feasible step may be taken or all components of \mathbf{g} are deleted.

We find the scalar A such that $\widehat{\mathbf{N}} = \mathbf{N} + \delta \mathbf{N} = \mathbf{N} + A\mathbf{p}$ maximizes P . Then we calculate the next gradient \mathbf{g} and repeat. This process ends when all components of \mathbf{g} are sufficiently small.

A block diagram of this algorithm appears in figure 4.

5. Primal algorithm

5.1. Algorithm structure

It is difficult to construct (N_1, \dots, N_{k-1}) to satisfy constraint (3) with equality. By using the dual algorithm as part of the solution of the primal problem, we avoid having to search surfaces like those of figure 2.

We solve the primal problem by dividing it into two subproblems: the dual and the one-dimensional primal problem (ODPP). The primal problem is: find N^{TOTAL} such that $P_{\text{MAX}}(N^{\text{TOTAL}}) = P^*$. The dual algorithm is used to evaluate $P_{\text{MAX}}(N^{\text{TOTAL}})$.

Let $\mathbf{N}(N^{\text{TOTAL}})$ be the distribution of buffer space that solves the dual problem with total buffer space N^{TOTAL} . The *optimal curve* is the set of $\mathbf{N}(N^{\text{TOTAL}})$ as N^{TOTAL} varies over a range of values. An example of an optimal curve appears in figure 3.

The nearly straight optimal curve is helpful in making the algorithm efficient. Linearizing it leads to a very good initial guess of (N_1, \dots, N_{k-1}) for the dual algorithm.

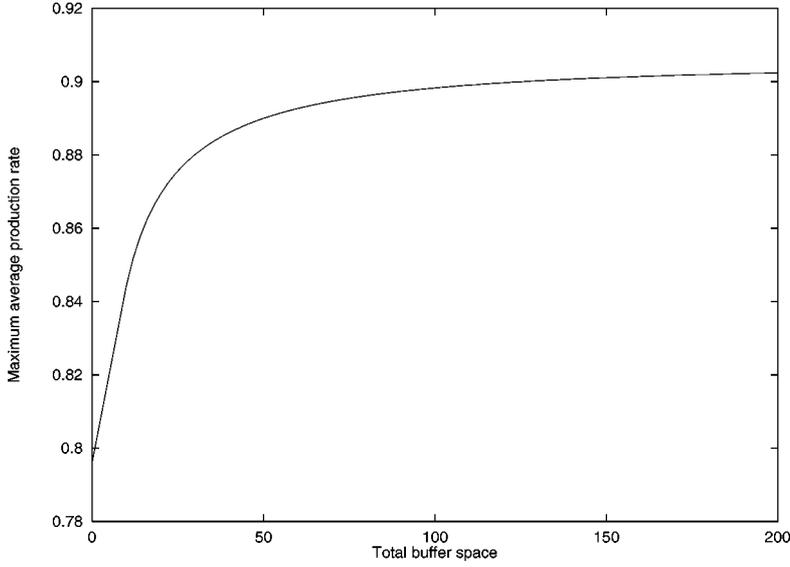
5.2. One-Dimensional Primal Problem (ODPP)

Figure 5 is the graph of $P_{\text{MAX}}(N^{\text{TOTAL}})$ for the three-machine, two-buffer system determined in table 1. (This is the same system that generated figure 3.)

Consider the following one-dimensional problem:

$$\begin{aligned} &\text{Minimize} && N^{\text{TOTAL}} && (8) \\ &\text{subject to} && P_{\text{MAX}}(N^{\text{TOTAL}}) \geq P^*; && P^* \text{ specified,} \\ &&& N^{\text{TOTAL}} \geq (k-1)N^{\text{MIN}}. \end{aligned}$$

Let the solution be called $N^{\text{TOTAL}*}$. This is essentially the same as the primal problem because the same $N^{\text{TOTAL}*}$ satisfies both, and because $\mathbf{N}(N^{\text{TOTAL}*})$ satisfies the primal.

Figure 5. P_{MAX} vs. N^{TOTAL} .

Schor [44] develops a one-dimensional search which is based on the property that $P_{\text{MAX}}(N^{\text{TOTAL}})$ is monotonically increasing in N^{TOTAL} . The algorithm requires an initial guess, N^1 . We evaluate $P_{\text{MAX}}(N^1)$ using the dual algorithm. We also need $P_{\text{MAX}}(0)$, the production rate of the line when all buffers are of size zero. The evaluation of $P_{\text{MAX}}(0)$ is described in section 5.4 under “Necessity”. We construct a linear approximation of $P_{\text{MAX}}(N^{\text{TOTAL}})$, and we use this to estimate N^2 , the second guess.

For $j \geq 2$, $P_{\text{MAX}}(N^j)$ is evaluated using the dual algorithm, j is incremented and then we iterate. Each iteration uses N^j and N^{j-1} to construct the next linear approximation of $P_{\text{MAX}}(N^{\text{TOTAL}})$. The process continues until $P_{\text{MAX}}(N^{\text{TOTAL}})$ is sufficiently close to P^* .

5.3. Algorithm

A block diagram of this algorithm appears in figure 6. The algorithm is initialized with $N^0 = (k-1)N^{\text{MIN}}$, $j = 2$, and with values of N^1 and search parameters specified. $P_{\text{MAX}}(N^1)$ is obtained by solving the dual; $P_{\text{MAX}}(N^0)$ is obtained from Buzacott’s zero-buffer formula in the discrete material case (in which case it is an approximation), or from evaluating the ADDX algorithm for $N_i = 0$ for all i in the continuous case.

N^j is determined from N^{j-2} and N^{j-1} by approximating the inverse function of $P(N^{\text{TOTAL}})$ with a straight line. That is, we approximate

$$N^{\text{TOTAL}} \approx mP_{\text{MAX}} + b,$$

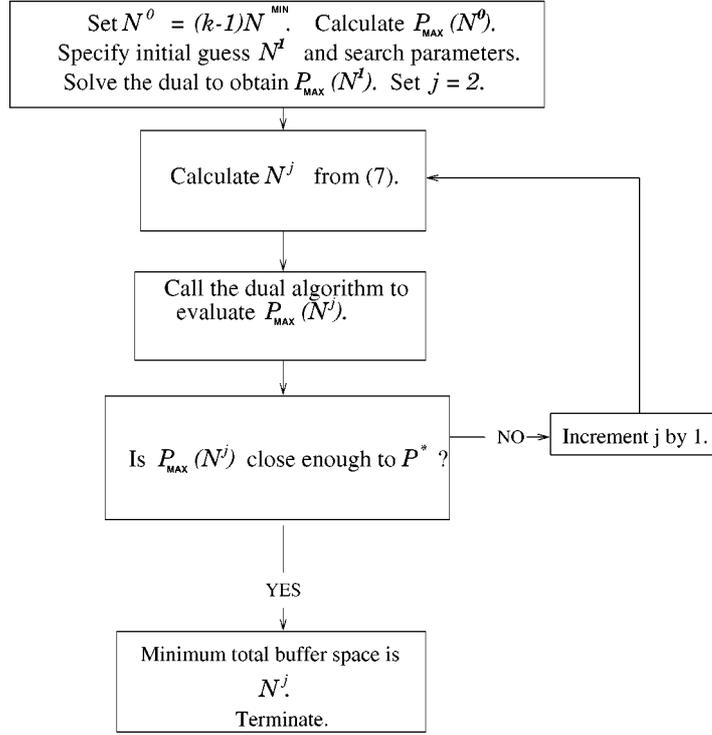


Figure 6. Block diagram of primal algorithm.

where

$$m = \frac{N^{j-1} - N^{j-2}}{P_{\text{MAX}}(N^{j-1}) - P_{\text{MAX}}(N^{j-2})}, \quad b = N^{j-1} - mP_{\text{MAX}}(N^{j-1})$$

and we select

$$N^j = mP^* + b.$$

We can also write this as

$$N^j = AN^{j-1} + BN^{j-2}, \quad (9)$$

where

$$A = \frac{P^* - P_{\text{MAX}}(N^{j-2})}{P_{\text{MAX}}(N^{j-1}) - P_{\text{MAX}}(N^{j-2})}, \quad B = \frac{P_{\text{MAX}}(N^{j-1}) - P^*}{P_{\text{MAX}}(N^{j-1}) - P_{\text{MAX}}(N^{j-2})}.$$

This approach must be modified when N^{TOTAL} is large because the P_{MAX} vs. N^{TOTAL} curve (figure 5) is very flat.

The dual algorithm is used to evaluate $P_{\text{MAX}}(N^j)$. If $P_{\text{MAX}}(N^j)$ is sufficiently close to P^* , the algorithm terminates.

5.4. Implementation issues

There are some important issues to consider when implementing these algorithms.

Feasibility. We check feasibility by requiring that the isolated production rate of each machine is greater than the production rate target. That is, in the discrete material case, a solution exists only if

$$P^* < \min_i \frac{r_i}{r_i + p_i}$$

and in the continuous material case, a solution exists only if

$$P^* < \min_i \frac{r_i \mu_i}{r_i + p_i}.$$

Necessity. We also want to determine if there is a need for buffers. If the production rate of the line with no buffers is greater than P^* there is no reason to run the algorithm. When the deterministic processing time model is used, calculate $P(0, \dots, 0)$ according to (1). If the processing rates of the machines are different, we invoke the ADDX algorithm with $N_i = 0$ for all i to evaluate $P(0, \dots, 0)$. In either case, if $P(0, \dots, 0) \geq P^*$, then all the optimal buffer sizes are 0.

ADDX and DDX algorithm initialization. These are iterative algorithms that must be initialized whenever they are invoked. Here, instead of a standard initialization, we use the final value of (N_1, \dots, N_{k-1}) from the last evaluation. This decreases the number of iterations required for the algorithms to converge. (See Gershwin and Goldis [20].)

Dual initialization. Just as the next guess for N^{TOTAL} is obtained from locally linearizing the $P_{\text{MAX}}(N^{\text{TOTAL}})$ graph, a good guess for (N_1, \dots, N_{k-1}) can be obtained from linearizing the optimal curve.

Rounding. In the final step of the algorithms for continuous material models (if the actual line is discrete), we round N^i so that the solution is integer. This can sometimes reduce P to below the target P^* , so some adjustment is made.

6. Algorithm behavior

6.1. Dual

Schor [44] observes that the solutions are either optimal or very close to optimal in all of the cases studied. The accuracy was established by comparing results with an algorithm of Gershwin and Goldis [20]. Gershwin and Goldis established the accuracy of their solutions by comparing them to the results of an exhaustive search.

When a continuous model is used for a discrete system, Schor observed that the accuracy of the solution is affected by the rounding methodology.

Table 2
Five-machine line (Ho et al. [27]).

Machine	1	2	3	4	5
MTTR = $1/r$	11	19	12	7	7
MTTF = $1/p$	20	167	22	22	26

Table 3
Buffer allocation for a five-machine line.

Case	Buffer i				N	P (DDX)	P (simulation)
	1	2	3	4			
Ho et al.	5	11	8	7	31	0.4914	0.4931
Dual	7	10	10	4	31	0.4943	0.4962

6.1.1. Discrete five-machine system

The optimal allocation algorithm of Ho et al. [26] is an iterative simulation-based algorithm. It has an efficient method that determines the production rate and gradient by simulating a transfer line once rather than k times. In this algorithm, A is determined empirically (and is evidently constant throughout the algorithm), and the next guess \widehat{N} is $N + Ap$, where N is the previous guess. Simulation is then used to estimate $P(N_1, \dots, N_{k-1})$ and g for each guess. When no improvement is realized, the algorithm stops.

In this example, we are determining the optimal buffer allocation for the five-machine discrete material line described in table 2. The final buffer distributions are presented in table 3. The production rates for both buffer distributions are calculated using the DDX algorithm and simulation. The distribution of buffer space seemed to be quite different, but the total buffer space and the production rate were very close.

The simulation was run for 50 runs of 100,000 cycles each and it required more than 5 minutes. To run the dual algorithm required 0.3 seconds. (This experiment was performed on a Sun SparcStation 2.)

6.1.2. Continuous seven-machine system

In this example, the machines have different processing times so we cannot use the discrete material analytical model. Ho et al. have no such restriction because they use simulation. As a result, we compare the optimal distribution that we get from using the continuous material model and then rounding with the distribution obtained by Ho et al.

We are seeking the optimal buffer allocation for the seven-machine line described in table 4. The final buffer distributions are presented in table 5. The production rates reported in the last column of table 5 are calculated using the ADDX algorithm.

Table 4
Seven-machine line (Ho et al. [26]).

Machine	1	2	3	4	5	6	7
MTTR = $1/r$	450	760	460	270	270	650	320
MTTF = $1/p$	820	5700	870	830	970	1900	1100
Cycle time = $1/\mu$	40	34	39	38	37	40	43

Table 5
Buffer allocation for a seven-machine line.

Case	Buffer i						N	ADDX production rate
	1	2	3	4	5	6		
Ho et al.	5	11	8	7	19	4	54	0.0126
Dual	8	10	13	10	9	4	54	0.0128

Table 6
Twelve-machine line (Park [40]).

i	1	2	3	4	5	6	7	8	9	10	11	12
r_i	0.35	0.15	0.4	0.4	0.3	0.2	0.3	0.3	0.4	0.4	0.3	0.25
p_i	0.037	0.015	0.02	0.03	0.03	0.01	0.02	0.02	0.02	0.03	0.03	0.01

6.2. Primal

The deterministic processing time model is used for all the examples in this chapter. We have initialized $N_i = 5 \forall i$ for all examples unless we explicitly state otherwise.

Again, Schor [44] compared his results with those of Gershwin and Goldis [20], who established the accuracy of their solutions by comparison with the results of an exhaustive search. The search was performed for several five-machine lines and buffer sizes of up to 20. In all cases, N^{TOTAL} is either the same or it differs by 1.

6.2.1. Comparison with prior work

In this section, we compare the performance of the primal algorithm to algorithms developed by Park [40] and Gershwin and Goldis [20]. Park's algorithm uses a two-phase tree search method and does not guarantee an optimal solution. Gershwin and Goldis' "combined" algorithm (called GG here) is a gradient algorithm which takes advantage of the decomposition to obtain a crude and cheap approximation of the gradient in the early iterations, and uses a more accurate approximation at the end. (This could be done with any iterative algorithm.)

Both the Park and GG algorithms solve the primal problem for the discrete material, deterministic processing time model. Both algorithms use the DDX algorithm to evaluate $P(N_1, \dots, N_{k-1})$. Both papers present results for the twelve-machine line described in table 6.

Table 7
Comparison of algorithms.

Case	Target rate	Park			GG			Primal		
		Actual rate	Long line evals.	N	Actual rate	Evals. long line (two-mach.)	N	Actual rate	Evals. long line (two-mach.)	N
1	0.85	0.8396	739	84						
2	0.85	0.8420	1,107	84						
3	0.85	0.8505	1,659	93	0.8507	182	87	0.8507	114	87
4	0.85	0.8505	1,838	93		(125,923)			(79,140)	
5	0.895	0.8950	510	390	0.8950	1,173	242	0.8950	342	243
						(1,481,928)			(534,820)	

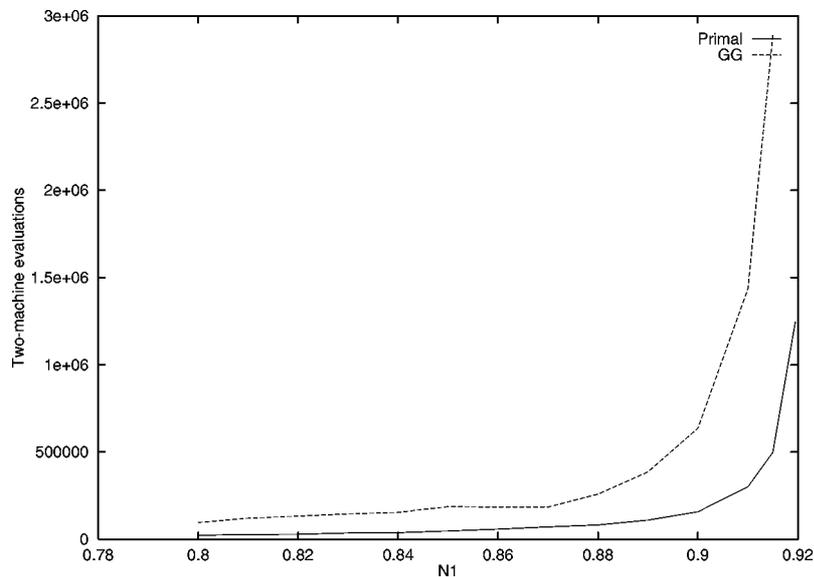


Figure 7. Two-machine evaluations as a function of P .

Table 7 compares the total buffer size, production rates, and number of long line and two-machine evaluations for experiments with two different values of P^* . The initial guess for the primal and GG algorithms, in both cases, is $N = (5, 5, 5, 5, 5, 5, 5, 5, 5, 5)$. (Park's algorithm does not require an initial guess.) Cases 1–4 are all the same problem: Park runs his algorithm with different search parameters.

We feel that the number of two-machine evaluations is the best measure of the run time of the algorithms because this measure is computer-independent. Park reports the number of long line evaluations. This is also independent of the computer, but it is a less absolute measure because there is more flexibility in the long line evaluation. (The initial conditions and the solution accuracy can vary during the course of the algorithm.)

Table 8
Ten-machine block.

Machine	<i>i</i>									
	1	2	3	4	5	6	7	8	9	10
r_i	0.094	0.095	0.045	0.078	0.069	0.094	0.095	0.045	0.078	0.069
p_i	0.007	0.008	0.003	0.004	0.006	0.007	0.008	0.003	0.004	0.006
P_i	0.93	0.92	0.94	0.95	0.92	0.93	0.92	0.94	0.95	0.92

Table 9
Long line computation times.

Number of machines	Two-machine evals.		N	
	Primal	GG	Primal	GG
10	82,176	258,661	433	443
20	1,814,868	6,867,878	995	995
30	15,955,968	47,157,058	1,556	1,557

For cases 1–4 the GG algorithm and the primal algorithm achieve identical solutions but the primal algorithm requires about one third fewer two-machine evaluations. Park's algorithm does not satisfy the production rate constraint in cases 1 or 2. In the cases where Park does satisfy the constraint, a larger N is calculated. When Park's algorithm satisfies the production rate constraint, it requires substantially more long line evaluations than the others.

For $P^* = 0.895$, the solution of the primal algorithm has one more buffer space than the GG algorithm. The primal algorithm requires the least computational effort in all cases.

Increasing P^ .* Figure 7 illustrates how execution time is impacted by P^* for the primal and GG algorithms. The number of two-machine evaluations, for the line described in table 8, is reported. When P^* approaches the infinite buffer production rate (0.92) the number of two-machine evaluations required increases significantly. The primal and GG algorithms calculate the same value for N in all cases.

Long lines. Table 9 demonstrates how increasing the length of the line impacts computational effort. Each line is composed of ten-machine blocks. Each block is described in table 8. $P^* = 0.88$ for all the lines. The significant increase in computational effort is a consequence of three factors. The first is the increased number of two-machine evaluations required for the DDX algorithm to converge. The second is that each gradient calculation requires an additional long line evaluation for each added machine. The third is that as the line increases the total buffer space increases and the algorithm must search a larger space.

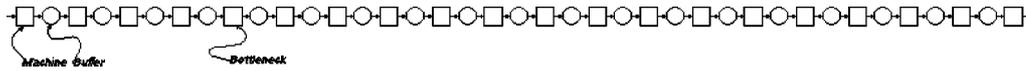


Figure 8. Twenty-machine line.

Table 10
Summary of twenty-machine line results.

Line	Production rate with no buffers	Minimal buffer space for production rate target (0.88 parts/min)
Case 1	0.487	430
Case 2	0.475	485
Case 3	0.475	523

6.2.2. Effect of a bottleneck

This example comes from Gershwin and Goldis [20]. Consider the design of a 20-machine production line. The machines have already been selected, and the only decision remaining is the amount of space to allocate for in-process inventory. The common operation time is one operation per minute. The target production rate is 0.88 parts per minute. The goal is to determine the smallest amount of in-process inventory space so that the line meets that target.

Three versions of this line are analyzed.

- In case 1, each machine is unreliable with mean time to fail (MTTF) 200 minutes and mean time to repair (MTTR) 10.5 minutes. Every machine could therefore operate in isolation at rate 0.95 parts per minute on the average, and easily meet the target. However, if there were no in-process inventory allowed, each machine's failure would cause every other machine to be forced idle, and the actual production rate would only be 0.487. To achieve the target production rate, the minimal required inventory space is room for 430 parts, distributed according to the figure.
- Cases 2 and 3 are similar to case 1, except that Machine 5 (and only Machine 5) is replaced by a less reliable model. There is thus a well-defined bottleneck. See figure 8. In both cases, the reliability of Machine 5 is reduced to 0.905.
 - * In case 2, MTTF = 100 and MTTR = 10.5 minutes.
 - * In case 3, MTTF = 200 and MTTR = 21 minutes.

The results are summarized in table 10 and figure 9. With no buffering between the machines, the production rate of both cases 2 and 3 would be 0.475. Again, each machine (including the bottleneck) could meet the required production rate in isolation, but the complete system cannot without inventory buffering between the machines. The minimal required total inventory space is room for 485 parts for case 2 and 523 for case 3. Figure 9 displays how the space is distributed in each case. The zero-buffer production rates were calculated from Buzacott [8].

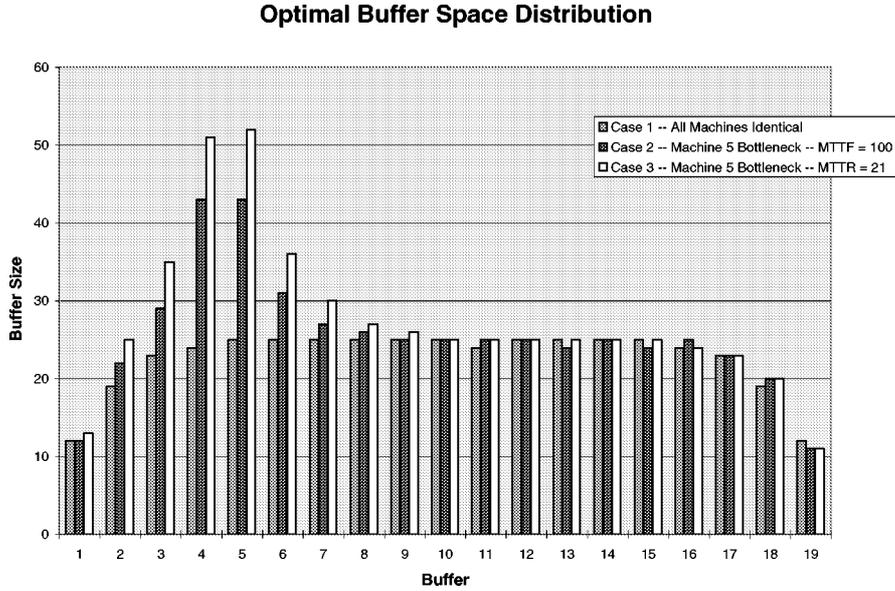


Figure 9. Optimal buffer distributions.

Observations:

- The additional inventory space is allocated approximately symmetrically around the bottleneck, but it is not all placed precisely at the bottleneck.
- Case 3 needs more space than case 2. This suggests that increased repair time is worse than increased failure frequency.
- Case 1 exhibits the inverted bowl phenomenon, observed by Hillier et al. [25]. A line with identical machines requires less buffer space at the ends because of the usual assumption that the first machine is never starved and the last machine is never blocked.
- The graphs are somewhat uneven. This is due mainly to the fact that we have restricted the solution to integers: we are assuming that the line is producing discrete items. In addition, the results shown may be slightly non-optimal.

7. Profit maximization

Here, we generalize the objective function to include inventory holding costs. Let \bar{n}_i be the average amount of inventory in Buffer B_i . This quantity is an output of the DDX and ADDX algorithms. Let Φ , a_i , and b_i be nonnegative cost coefficients. The profit is

$$\Pi = \Phi P(N_1, \dots, N_{k-1}) - \sum_{i=1}^{k-1} a_i N_i - \sum_{i=1}^{k-1} b_i \bar{n}_i,$$

in which Φ is the profit per item sold, when all costs other than inventory costs at this system are considered. The cost of inventory space is linear in the size of the space, and the cost per unit space in Buffer B_i is a_i . The cost of inventory held is linear in the average amount of inventory, and the cost per item stored in Buffer B_i is b_i .

We know of no published papers that discuss the qualitative properties of profit as a function of buffer sizes (comparable to section 2.3.2). Gershwin [18] shows that average buffer level may be either a concave or convex function of buffer size for a two-machine line. The shape of this function depends upon the reliability parameters and processing rates of the machines. Consequently, if profit is a function of average buffer level, it may not be concave.

7.1. Profit maximization with total buffer space constraint

$$\begin{aligned} & \text{Maximize } \Pi & (10) \\ & \text{subject to } N^{\text{TOTAL}} = \sum_{i=1}^{k-1} N_i; \quad N_i \geq N^{\text{MIN}}. \end{aligned}$$

7.2. Profit maximization without space constraints

$$\begin{aligned} & \text{Maximize } \Pi & (11) \\ & \text{subject to } N_i \geq N^{\text{MIN}}. \end{aligned}$$

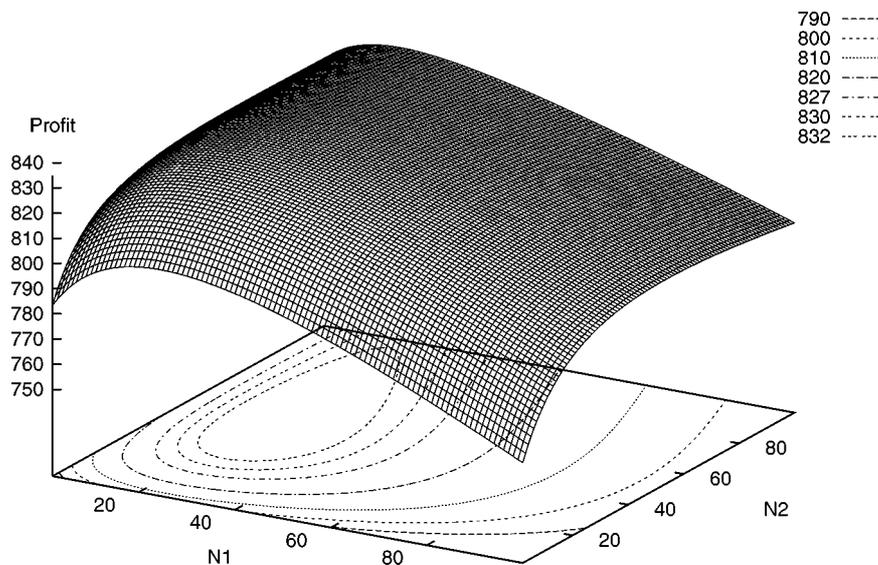
7.3. Properties of Π

We can no longer rely on the objective function being monotonic or concave. However, numerical experience suggests that it has a single maximum, and this appears to be enough for the algorithms to work effectively.

Figure 10 is a graph of $\Pi(N_1, N_2)$ vs. N_1 and N_2 for the balanced three-machine continuous material line with parameters $r_i = 0.1$, $p_i = 0.01$, and $\mu_i = 1$ for $i = 1, 2, 3$. The cost coefficients are $\Phi = 1000$, $a_i = 0$, and $b_i = 1$ for $i = 1, 2$. Therefore $\Pi(N_1, N_2) = 1000P(N_1, N_2) - \bar{n}_1 - \bar{n}_2$. The contour lines indicate that the function $\Pi(N_1, N_2)$ has a global maximum and no other local maxima. The maximum profit for this line is realized when $(N_1, N_2) = (21, 49)$. (These have been rounded to integers.) Note that figure 10 is not symmetric and the buffer sizes are not equal, even though the line is balanced.

7.4. Profit maximization algorithms

The algorithm for the constrained profit maximization problem (10) is essentially the same as the algorithm for the dual problem (section 4). It is the same as figure 4 with P replaced by Π .

Figure 10. Profit vs. N_1 and N_2 .

The algorithm for the unconstrained profit maximization problem (11) follows the primal algorithm (section 5) very closely. It also consists of the iterated solution of two problems: the constrained profit maximization problem (in place of the dual) and a one-dimensional profit maximization problem (in place of the one-dimensional primal).

Details on the algorithms, as well as numerical examples, appear in Schor [44].

8. Conclusions

The algorithms described here are fast and accurate. They can be used in the design of production systems. These methods can easily be extended to assembly/disassembly systems [18]. It would be valuable to extend them to systems that have loops, such as systems with captive pallets or fixtures. Other useful extensions include systems with multiple part types, with set-up changes, etc. The main challenges in such extensions are more in the modeling and analysis than in the optimization algorithm.

These extensions are valuable for designing more general systems. As shown in Gershwin [19] they are also valuable in designing real-time scheduling policies for manufacturing systems.

References

- [1] I. Adan and J. van der Wal, Monotonicity of the throughput in single server production and assembly networks with respect to the buffer sizes, in: *Queueing Networks with Blocking*, eds. H.G. Perros

- and T. Altiok (Elsevier Science, 1989) pp. 345–356.
- [2] T. Altiok and S. Stidham Jr., The allocation of interstage buffer capacities in production lines, *IIE Transactions* 15 (1983) 292–299.
 - [3] V. Anantharam and P. Tsoucas, Stochastic concavity of throughput in series of queues with finite buffers, *Advances in Applied Probability* 22 (1990) 761–763.
 - [4] D.R. Anderson and C.L. Moodie, Optimal buffer storage capacity in production line systems, *International Journal of Production Research* 7(3) (1969) 233–240.
 - [5] K. Barten, A queueing simulator for determining optimum inventory levels in a sequential process, *Journal of Industrial Engineering* 13(4) (1962) 245–252.
 - [6] M. Burman (1995), New results in flow line analysis, Ph.D. thesis, MIT, Cambridge MA (1995).
 - [7] J.A. Buzacott, Automatic transfer lines with buffer stocks, *International Journal of Production Research* 5(3) (1967) 182–200.
 - [8] J.A. Buzacott, Prediction of the efficiency of production systems without internal storage, *International Journal of Production Research* 6(3) (1968) 173–188.
 - [9] J.A. Buzacott and J.G. Shanthikumar, *Stochastic Models of Manufacturing Systems* (Prentice-Hall, 1993).
 - [10] M. Caramanis, Production system design: A discrete event dynamic system and generalized Benders' decomposition approach, *International Journal of Production Research* 25(8) (1987) 1223–1234.
 - [11] We-Min Chow, Buffer capacity analysis for sequential production lines with variable processing times, *International Journal of Production Research* 25(8) (1987) 1183–1196.
 - [12] Y. Dallery, R. David and X.L. Xie, An efficient algorithm for analysis of transfer lines with unreliable machines and finite buffers, *IIE Transactions* 20(3) (1988) 280–283.
 - [13] Y. Dallery and S.B. Gershwin, Manufacturing flow line systems: A review of models and analytical results, *Queueing Systems* 12(1–2) (1992) 3–94.
 - [14] Y. Dallery, Z. Liu and D. Towsley, Equivalence, reversibility, symmetry and concavity properties in fork/join queueing networks with blocking, *Journal of the ACM* 41(5) (1994) 903–942.
 - [15] E. Dogan and T. Altiok, Gradient estimation of the output rate in transfer lines, Preprint, Rutgers University, Dept. of Industrial Engineering.
 - [16] M.C. Freeman, The effects of breakdowns and interstage storage on production line capacity, *Journal of Industrial Engineering* 15(4) (1964) 194–200.
 - [17] S.B. Gershwin, An efficient decomposition algorithm for the approximate evaluation of tandem queues with finite storage space and blocking, *Operations Research* 35 (1987) 291–305.
 - [18] S.B. Gershwin, *Manufacturing Systems Engineering* (Prentice-Hall, 1994).
 - [19] S.B. Gershwin, Design and operation of manufacturing systems – control- and system-theoretical models and issue, in: *1997 American Control Conference*, Albuquerque, New Mexico (June 4–6, 1997).
 - [20] S.B. Gershwin and Y. Goldis, Efficient algorithms for transfer line design, MIT Laboratory for Manufacturing and Productivity Report LMP-95-005 (1995).
 - [21] S.B. Gershwin and I.C. Schick, Continuous model of an unreliable material flow systems with a finite interstage buffer, MIT LIDS Report LIDS-R-1039 (1980).
 - [22] P. Glasserman and D.D. Yao, Structured buffer-allocation problems, *Journal of Discrete Event Dynamic Systems* 6 (1996) 9–42.
 - [23] J.M. Hatcher, The effect of internal storage on the production rate of a series of stages having exponential service times, *AIIE Transactions* 1(2) (1969) 150–156.
 - [24] F.S. Hillier and K.C. So, The effect of machine breakdowns and interstage storage on the performance of production line systems, *International Journal of Production Research* 29(10) (1991) 2043–2055.
 - [25] F.S. Hillier, K.C. So and R.W. Boling, Toward characterizing the optimal allocation of storage space in production line systems with variable processing times, *Management Science* 39(1) (1993)

- 126–133.
- [26] Y.C. Ho, M.A. Eyler and T.T. Chien, A gradient technique for general buffer storage design in a production line, *International Journal of Production Research* 17(6) (1979) 557–580.
 - [27] Y.C. Ho, M.A. Eyler and T.T. Chien, A new approach to determine parameter sensitivities of transfer lines, *Management Science* 29(6) (1983) 700–714.
 - [28] D. Jacobs and S.M. Meerkov, A system-theoretic property of serial production lines: Improvability, The University of Michigan College of Engineering Control Group Reports, Report No. CGR-93-1 (1993).
 - [29] D. Jacobs and S.M. Meerkov, System-theoretic properties of the process of continuous improvement in production systems, in: *Proceedings of the American Control Conference*, Baltimore, Maryland (1994) pp. 3323–3327.
 - [30] M.A. Jafari and J.G. Shanthikumar, Determination of optimal buffer storage capacities and optimal allocation in multistage automatic transfer lines, *IIE Transactions* 21(2) (1989) 130–134.
 - [31] A.D. Knott, Letter, *AIIE Transactions* 2(3) (1970) p. 273.
 - [32] S.A. Kraemer and R.F. Love, A model for optimizing the buffer inventory storage size in a sequential production system, *AIIE Transactions* II(1) (1970) 64–69.
 - [33] C.-T. Kuo, J.-T. Lim and S.M. Meerkov, Improvability analysis of a machining transfer line: An application, The University of Michigan College of Engineering Control Group Reports, Report No. CGR-94-08 (1994).
 - [34] C.-M. Liu and C.-L. Lin, Performance evaluation of unbalanced serial production lines, *International Journal of Production Research* 32(12) (1994) 2897–2914.
 - [35] G.E. Martin, Predictive formulae for unpaced line efficiency, *International Journal of Production Research* 31(8) (1993) 1981–1990.
 - [36] G.E. Martin, Optimal design of production lines, *International Journal of Production Research* 32(5) (1994) 989–1000.
 - [37] L.E. Meester and J.G. Shanthikumar, Concavity of the throughput of tandem queueing systems with finite buffer storage space, *Advances in Applied Probability* 22 (1990) 764–767.
 - [38] J.A. Muckstadt and S.R. Tayur, A comparison of alternative kanban control mechanisms. I. Background and structural results, *IIE Transactions* 27 (1995) 140–150.
 - [39] J.A. Muckstadt and S.R. Tayur, A comparison of alternative kanban control mechanisms. II. Experimental results, *IIE Transactions* 27 (1995) 151–161.
 - [40] T. Park, A two-phase heuristic algorithm for determining buffer sizes of production lines, *International Journal of Production Research* 31(3) (1993) 613–631.
 - [41] E.L. Plambeck, B.-R. Fu, S.M. Robinson and R. Suri, Throughput optimization in tandem production lines via nonsmooth programming, *Proceedings of the 1993 Summer Computer Simulation Conference*, Boston (1993).
 - [42] S.G. Powell, Buffer allocation in unbalanced three-station serial lines, *International Journal of Production Research* 32(9) (1994) 2201–2218.
 - [43] S.G. Powell and D.F. Pyke, Optimal allocation of buffers in serial production lines with a single bottleneck, The Amos Tuck School of Business Administration, Dartmouth College, Working Paper No. 301 (1994).
 - [44] J.E. Schor, Efficient algorithms for buffer allocation, M.S. thesis, MIT, Cambridge, MA (1995).
 - [45] D. Seong, S.Y. Chang and Y. Hong, Heuristic algorithms for buffer allocation in a production line with unreliable machines, Department of Industrial Engineering, POSTECH, Korea; *International Journal of Production Research* (1994) (accepted for publication).
 - [46] D. Seong, S.Y. Chang and Y. Hong, An algorithm for buffer allocation with linear resource constraints in a continuous flow production line, Technical Report IE-TR-94-05, Department of Industrial Engineering, POSTECH, Korea (1994).
 - [47] B.A. Sevast'yanov, Influence of storage bin capacity on the average standstill time of a production line, *Teoriya Veroyatnostey i ee Primeneniya* 7(4) 438–447 (in Russian). (English translation: *Theory*

- of Probability and its Applications 7(4) (1962) 429–438.)
- [48] J.G. Shanthikumar and D.D. Yao, Monotonicity and concavity properties in cyclic queueing networks with finite buffers, in: *Queueing Networks with Blocking*, eds. H.G. Perros and T. Altiok (Elsevier Science, 1989) pp. 325–344.
 - [49] T.J. Sheskin, Allocation of interstage storage along an automatic production line, *AIIE Transactions* 8(1) (1976) 146–152.
 - [50] J. MacGregor Smith and S. Daskalaki, Buffer space allocation in automated assembly lines, *Operations Research* 36(2) (1988) 343–357.
 - [51] A.L. Soyster, J.W. Schmidt and M.W. Rohrer, Allocation of buffer capacities for a class of fixed cycle production lines, *AIIE Transactions* 11(2) (1979) 140–146.
 - [52] S.R. Tayur, Properties of serial kanban systems, *Queueing Systems* 12 (1992) 297–318.
 - [53] H. Yamashita and T. Altiok, Buffer capacity allocation for a desired throughput in production lines, Presented at the Workshop on *Performance Evaluation and Optimization of Production Lines*, International Workshop held on May 19–21, 1997, on Samos Island, in the Department of Mathematics, University of the Aegean, GR-832 00 Karlovassi, Samos Island, Greece.