

Reducing work-in-process inventory in certain classes of flexible manufacturing systems

J. George SHANTHIKUMAR

School of Business Administration, University of California, Berkeley, CA, U.S.A.

Kathryn E. STECKE

Graduate School of Business Administration, The University of Michigan, Ann Arbor, MI, U.S.A.

Abstract: In this paper, we establish that maintaining a balanced workload on each machine over time stochastically minimizes the work-in-process inventory in certain types of flexible manufacturing systems (FMSs) with finite or infinite common input buffer storage and an ample buffer at each machine. The results obtained here complement those obtained by Stecke and Morin (1985), in which it is established that balancing workloads maximizes expected production, again for the same, particular types of FMSs. Stecke and Morin (1985) treats a static FMS loading problem, while this paper addresses a dynamic problem which considers three strategies to control the release of parts into the system.

Keywords: Flexible manufacturing systems, work-in-process inventory, workload balancing, stochastic minimization, open and closed queueing networks

1. Introduction

A flexible manufacturing system (FMS) is an automated alternative to the conventional means of batch manufacturing, to date applied mainly in the metalcutting industry. An FMS consists of a number of computer numerically controlled machine tools, which are linked together by an automated material handling system. Individual parts of different types can be machined simultaneously in unit batch sizes.

Efficient use of such FMSs requires careful pre-production planning (see, e.g., Buzacott and Shanthikumar (1980) and Stecke (1981)). One of such planning problems concerns the appropriate loading of the tool magazines of all machines in an FMS. By FMS machine loading, we mean the

allocation of all operations and their associated cutting tools among the FMS machine tools, in order to define which machines will be able to perform each operation of the part types that have been chosen to be machined next. This decision, in part, determines each machine tool's workload.

Earlier analyses of this and similar planning issues in the context of job shops and flow lines have indicated the superiority of either unbalancing or balancing each machine's workload (depending on some particular aspects of the production system, such as buffer size, processing time distribution, ...) with respect to the maximization of production rate or with respect to the minimization of the work-in-process inventory. (See Stecke and Morin (1982) for references to such works.)

The question of balancing workloads in the context of FMS was first raised by Buzacott and Shanthikumar (1980). Based on an asymptotic analysis, they established that balancing the workload on all machines maximizes the expected production rates whenever the number of parts in the system is *very large*, i.e., approaching infinity. Subsequently, Stecke (1981) and Stecke and Morin

Kathryn E. Stecke's work was supported in part by NSF Grant No. ECS 8406407 and a Grant from the Graduate School of Business Administration, The University of Michigan.

(1982), based on empirical results, supported the optimality of balanced workloads with respect to the expected production of FMSs with single machine workcenters (i.e., with no pooling of similar machines) and with a *finite* number of parts.

However, it was also observed in Stecke (1981) that balancing the workload per machine need not maximize the expected production of FMSs having a different number of machines in different machine groups (i.e., with pooling of machines into machine groups). Shanthikumar (1982) for the first time provided an analytical proof for the optimality of the balanced load with respect to the production rate of an FMS with no pooling. An empirical proof can be found in Stecke (1981). An alternative analytical proof can be found in Stecke and Morin (1985). Several extensions and simpler proofs of similar results can be found in Yao (1984a, b), Yao and Kim (1984) and Stecke (1986). All of these FMS loading problems involve static allocations of workload among machines (and eventually, of operations and cutting tools). These problems are solved and the solutions implemented before the FMS begins to operate. The detailed, actual FMS loading problem of allocating operations and tools to machines to maximize production has been addressed in Stecke (1983, 1986a). Again, the problem is static.

The scenario in this paper is different, but complementary to the previous studies, Stecke and Morin (1985), in particular. An open queueing network is used to model the random arrival of individual parts to the FMS. Given that there are N parts in the system, a closed queueing network then provides the expected production, if the workload on each machine is also provided. Three strategies are considered here to release the parts, which have arrived into a production control area, into the FMS. Under these three release policies, if the workload per machine remains balanced over time, not only is the expected production maximized, but we also show that the in-process inventory is stochastically minimized.

The random arrivals are of individual parts, and not batches, of a limited number of similar part types. Each operation can be performed on only one machine. The cutting tools for each operation have already been loaded into the appropriate tool magazine. There is no pooling of machines and there is a finite or infinite common buffer area (called here a production control area)

for incoming parts, which are subsequently released into the FMS. This paper is treating this dynamic problem. The results presented in this paper indicate that the work-in-process inventory in certain classes of FMSs is minimized by maintaining a balanced load on each machine over time as the parts are input into the FMS. Methods to determine the minimum inventory requirements are provided in Stecke (1985a).

The open and closed queueing network models are described in Section 2, as well as the three policies that we consider to release parts into the FMS. The optimality of balancing the workload on each machine to maximize expected production using the closed queueing network is presented in Section 3. Section 4 contains results concerning the stochastic minimization of the number of parts in the FMS. A summary and future research needs are provided in Section 5.

2. The queueing network models

A Jackson type queueing network is used to model a flexible manufacturing system. The external arrival of parts forms a Poisson stochastic process with rate λ . An external arrival is received into a production control area from which the parts are dispatched into the flexible manufacturing system. (See Figure 1.) Parts from the production control area are released to the system according to some dispatch policy. In this paper, we consider the following three policies that describe different means to control the input of parts to the FMS:

Policy I. Dispatch parts to the FMS as soon as they are received into the production control area.

Policy II. Any parts that are received into the production control area when there are already Z parts in the FMS are rejected (i.e., lost). Otherwise, parts are released to the FMS as soon as they are received by the production control area.

Policy III. Dispatch a part to the FMS if and only if the total number of parts in the FMS is less than some value Z . This may represent a predetermined input control policy or reflect the fact that the number of pallets available in the shop is limited to Z (e.g., see Buzacott (1982) and Buzacott and

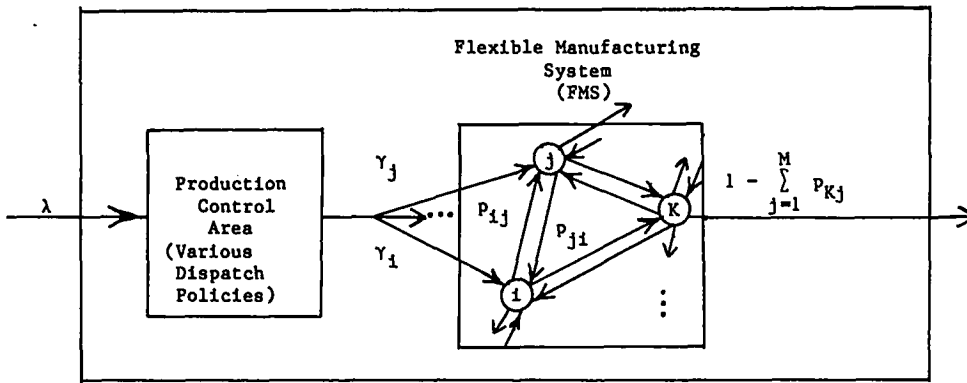


Figure 1. A schematic diagram of a queuing network model of a flexible manufacturing system.

Shanthikumar (1980)). Parts are neither lost nor rejected.

The FMS contains M machines (each is single, there is no pooling). The first operation of a part that is released to the FMS is performed by machine i with probability γ_i . A part that has completed its processing at machine i will proceed next to machine j with probability p_{ij} . This part is an internal arrival to machine j . On the other hand, upon leaving machine i , a part may depart the network of queues (i.e., the FMS) with probability $1 - \sum_{j=1}^M p_{ij}$. The transportation time is assumed to be either negligible compared to the processing or delay times or averaged in with the processing time. The service, or processing, time at machine i is exponentially distributed with mean t_i , $i = 1, 2, \dots, M$. Let q_i be the expected number of visits made to machine i by an arbitrary part. Then, using the Markovian property of the part transitions from machine to machine, it can be shown that:

$$q_i = \gamma_i + \sum_{j=1}^M q_j p_{ji}, \quad i = 1, 2, \dots, M. \quad (1)$$

The mean load imposed on machine i by a single part is then equal to $q_i t_i$ and we denote this by x_i . If λ_e is the effective part arrival rate (i.e., the average number of parts arriving per unit time) to the FMS (different for the different part input policies), then the workload rate utilization of machine i is $\lambda_e x_i$. The total mean workload L imposed on the system by a single part is equal to $\sum_{i=1}^M x_i$. L is the average time required to process a

part through the system. The mean total load remains the constant L , which is independent of how the machines are loaded.

Then, under a balanced workload, the mean load imposed by a single part on a machine tool is $L/M (\equiv x^*)$. The balanced workload distribution is then

$$x^* = (x^*, x^*, \dots, x^*).$$

3. Production rate

In this section, the number of parts in the system is kept constant at level N , i.e., the queuing network is closed. Then the expected production rate, $PR_N(x)$, of the FMS with the mean load distribution $x = (x_1, x_2, \dots, x_M)$ is (see, e.g., Solberg (1977) and Buzacott and Shanthikumar (1980)),

$$PR_N(x) = G(x, N-1)/G(x, N), \quad N \geq 1, \quad (2)$$

where

$$G(x, N) = \sum_{n \in S_N} \left(\prod_{i=1}^M x_i^{n_i} \right), \quad N \geq 1, \quad (3)$$

$$G(x, 0) = 1,$$

$$n = (n_1, n_2, \dots, n_M), \quad n_i \geq 0, \quad i = 1, 2, \dots, M,$$

and

$$S_N = \left\{ n: \sum_{i=1}^M n_i = N \right\}, \quad N \geq 1.$$

Here, n_i is the number of parts at machine i ,

including the part in process if $n_i \geq 1$.

The following theorem concerning maximum expected production, using different proofs, has been established in various ways for slightly different systems in Stecke (1981), Shanthikumar (1982), Stecke and Morin (1982), and Yao (1984b).

Theorem 1. *The expected production rate, $PR_N(x)$, is maximized by balancing the workload on all machines. That is,*

$$PR_N(x^*) = \max_{x: \sum_{i=1}^M x_i = L} \{PR_N(x)\},$$

for $x^* = (L/M, L/M, \dots, L/M)$.

Corollary 1. *The expected production rate of a single-server closed queueing network, $PR_N(x)$, is bounded by*

$$PR_N(x) \leq \frac{N}{M+N-1} \frac{M}{L}. \tag{4}$$

We note here that for information concerning the robustness of queueing networks as aggregate models of the steady state performance of FMSs, see Suri (1983), for example.

4. Number of parts in the system

Let $P_n(x)$ be the steady state probability that there are n parts ($n = 0, 1, \dots$) in the system at an arbitrary time epoch under a mean work load distribution $x = (x_1, x_2, \dots, x_M)$. It can be shown (see, e.g., Shanthikumar and Sargent (1981)) that under Policy I:

$$P_n(x) = \frac{\lambda}{PR_n(x)} P_{n-1}(x), \quad n = 1, 2, \dots, \tag{5}$$

$$P_0(x) = \left[1 + \sum_{n=1}^{\infty} \left\{ \prod_{i=1}^n \left(\frac{\lambda}{PR_i(x)} \right) \right\} \right]^{-1}.$$

and under Policy II:

$$P_n(x) = \frac{\lambda}{PR_n(x)} P_{n-1}(x), \quad n = 1, 2, \dots, Z, \tag{6}$$

$$P_n(x) = 0, \quad n = Z + 1, Z + 2, \dots,$$

$$P_0(x) = \left[1 + \sum_{n=1}^Z \left\{ \prod_{i=1}^n \left(\frac{\lambda}{PR_i(x)} \right) \right\} \right]^{-1}.$$

For Policy III, Buzacott and Shanthikumar (1980) has observed that the following approximation is within a 90% confidence interval of the simulation results for the mean number of parts in the system. (See Table 2 there.)

$$P_n(x) = \frac{\lambda}{PR_n(x)} P_{n-1}(x), \quad n = 1, 2, \dots, Z,$$

$$P_n(x) = \frac{\lambda}{PR_Z(x)} P_{n-1}(x), \quad n = Z + 1, Z + 2, \dots,$$

$$P_0(x) = \left(1 + \sum_{n=1}^{Z-1} \left\{ \prod_{i=1}^n \left(\frac{\lambda}{PR_i(x)} \right) \right\} + \prod_{i=1}^Z \frac{\lambda}{PR_i(x)} \cdot \left(\frac{PR_Z(x)}{PR_Z(x) - \lambda} \right) \right)^{-1} \tag{7}$$

Let $M(n)/M(n)/1$ be a birth-death queueing system with state-dependent arrival rates $\lambda = (\lambda_n)_0^\infty$ and state-dependent service rates $\mu = (\mu_n)_1^\infty$. Also, let $F_n(\lambda, \mu)$ be the steady state probability that the number of parts in the system is less than or equal to n . Then

$$F_n(\lambda, \mu) = \sum_{r=0}^n \rho_r(\lambda, \mu) / \sum_{r=0}^{\infty} \rho_r(\lambda, \mu), \tag{8}$$

where

$$\rho_0(\lambda, \mu) = 1, \tag{9}$$

$$\rho_r(\lambda, \mu) = \prod_{i=1}^r \frac{\lambda_{i-1}}{\mu_i}, \quad r = 1, 2, \dots,$$

The results (5), (6), and (7) are special cases of (8) with:

- (a) $\lambda_n = \lambda$ ($n = 0, 1, \dots$) and $\mu_n = PR_n(x)$ ($n = 1, 2, \dots$) for Policy I;
- (b) $\lambda_n = \lambda$ ($n = 0, 1, \dots, Z - 1$), $\lambda_n = 0$ ($n = Z, Z + 1, \dots$), and $\mu_n = PR_n(x)$ ($n = 1, 2, \dots$) for Policy II; and
- (c) $\lambda_n = \lambda$ ($n = 0, 1, \dots$), $\mu_n = PR_n(x)$ ($n = 1, 2, \dots, Z$), and $\mu_n = PR_Z(x)$ ($n = Z + 1, Z + 2, \dots$) for Policy III.

Now, taking the partial derivatives of (8) with respect to λ_r and μ_r , one can see that

$$\frac{\partial}{\partial \lambda_r} F_n(\lambda, \mu) \leq 0, \quad r = 0, 1, \dots, \tag{10}$$

$$\frac{\partial}{\partial \mu_r} F_n(\lambda, \mu) \geq 0, \quad r = 1, 2, \dots,$$

for $n = 0, 1, 2, \dots$. The next lemma then immediately follows from the equations of (10).

Lemma 1. Let $(M(n)/M(n)/1)_i$ be a birth-death queueing model with arrival rates $\lambda_i = (\lambda_{in})_{n=0}^{\infty}$ and service rates $\mu_i = (\mu_{in})_{n=1}^{\infty}$, for $i = 1$ and 2 . If $\lambda_1 \geq \lambda_2$ and $\mu_1 \leq \mu_2$, then

$$F_n(\lambda_1, \mu_1) \leq F_n(\lambda_2, \mu_2), \quad n = 0, 1, \dots$$

Combining Theorem 1 and Lemma 1, one obtains the following.

Theorem 2. The number of parts in the classes of FMS under consideration, following any of the above three input control policies, is stochastically minimized by balancing the workload on all machines.

Note that the statement of Theorem 2 is heuristic for Policy III, since in this case, it is the approximation, which was given by equation (7), that is used to obtain the result.

5. Conclusions

In this paper, we have established the superiority of balancing workloads in flexible manufacturing systems utilizing no pooling, both to stochastically minimize the number of parts in the system as well as to maximize the expected production. Three policies to control the input of parts into the system are considered. There is a Poisson arrival process to a finite or infinite common input buffer area, exponentially distributed service times, Markovian part transfers from machine to machine, and an ample buffer at each machine.

The optimality of balanced workloads is established with the constraints of single machine work centers and $\sum_{i=1}^M x_i$ is a constant. When these conditions are relaxed, the balanced load need not be optimal. On the other hand, for similar systems, but with pools of identically tooled machines that are unbalanced, expected production is maximized by an unbalanced load, again with the constraint that $\sum_{i=1}^M x_i$ is a constant. (See Stecke (1981) and Stecke and Solberg (1985).)

Yao (1984a,b) and Yao and Kim (1984) have obtained similar balancing results for cases where the above two constraints are relaxed. For example, they prove that balancing maximizes expected production when there is pooling but with the same number of machines in each group. This is also observed in Stecke (1981) and Stecke and Solberg (1985).

Empirical studies using both simulation and approximations are currently underway to extend these results to more general cases.

Additional work is required to implement a balancing workload objective over time. The FMS balancing problem is different than that of the flow shop or job shop in several ways. A flow shop is balanced once, during its design. At the other extreme, it is very difficult to balance a job shop. Usually the work is given, to result in one particular bottleneck machine type, which changes over time.

An FMS can produce in unit batch sizes. There are planning decisions to be made that impact balancing (or unbalancing). These include: selecting the part types to be produced next; determining the ratios at which these part types are to be produced; allocating pallets and fixtures among the part types; determining the minimum number of pallets required in the system; loading tools and assigning operations to machines; determining the appropriate input sequence into the FMS; and finally, the actual scheduling of parts through the system. More information concerning what these problems are for an FMS can be found in Stecke (1983, 1986a).

Some models that can be used to address these problems are overviewed in Suri (1984). An efficient algorithm to determine the most balanced allocations for the types of FMSs considered in this paper is provided in Berrada and Stecke (1985). The algorithm also applies to systems of pooled machines, when each pool contains the same number of machines. Extensions to unbalanced groupings (hence unbalanced loadings) are underway. In any case, further research is required to discover methods to address all of the above-mentioned problems in order to implement either a balancing (or an unbalancing) operational objective.

References

- Berrada, M., and Stecke, K.E. (1986), "A branch and bound approach for machine load in flexible manufacturing systems", *Management Science*, forthcoming.
- Buzacott, J.A. (1982), "'Optimal' operating rules for automated manufacturing systems", *IEEE Transactions on Automatic Control*, 80-86.
- Buzacott, J.A., and Shanthikumar, J.G. (1980), "Models for understanding flexible manufacturing systems", *AIIE Transactions* 12 (4), 339-350.

- Shanthikumar, J.G. (1982), "On the superiority of balanced load in a flexible manufacturing system", technical report, Department of IE & OR, Syracuse University, New York.
- Shanthikumar, J.G., and Sargent, R.G. (1981), "A hybrid simulation/ analytic model of a computerized manufacturing system", *Proceedings of the 9th IFORS Conference*, Hamburg, Fed. Rep. Germany, pp. 901-915.
- Solberg, J.J. (1977), "A mathematical model of computerized manufacturing systems", *Proceedings of the International Conference on Production Research*, Tokyo, Japan.
- Stecke, K.E. (1981), *Production Planning Problems for Flexible Manufacturing Systems*, Ph.D. dissertation, School of Industrial Engineering, Purdue University, West Lafayette, IN.
- Stecke, K.E. (1983), "Formulation and solution of nonlinear integer production planning problems for flexible manufacturing systems", *Management Science* 29 (3), 273-288.
- Stecke, K.E. (1985a), "Procedures to determine both appropriate production ratios and minimum inventory requirements to maintain these ratios in flexible manufacturing systems", Working paper no. 448, Graduate School of Business Administration, The University of Michigan, Ann Arbor, MI.
- Stecke, K.E., (1986a), "A hierarchical approach to solving machine grouping and loading problems of flexible manufacturing systems", *European Journal of Operational Research*, 24 (3), 369-378.
- Stecke, K.E. (1986b), "On the nonconcavity of throughput in certain closed queueing networks", *Performance Evaluation* 6(3), to appear.
- Stecke, K.E., and Morin, T.L. (1982), "Optimality of balanced workloads in flexible manufacturing systems", Working paper no. 289, Graduate School of Business Administration, The University of Michigan, Ann Arbor, MI.
- Stecke, K.E., and Morin, T.L. (1985), "The optimality of balancing workloads in certain types of flexible manufacturing systems", *European Journal of Operational Research* 20 (1), 68-82.
- Stecke, K.E., and Solberg, J.J. (1985), "The optimality of unbalancing both workloads and machine group sizes in closed queueing networks of multiserver queues", *Operations Research* 33 (4), 882-910.
- Suri, R. (1983), "Robustness of queueing network formulas", *Journal of the Association for Computing Machinery* 30 (3), 564-594.
- Suri, R. (1984), "An overview of evaluative models for flexible manufacturing systems", *Proceedings of the First ORSA/TIMS Special Interest Conference on Flexible Manufacturing Systems*, Ann Arbor, MI, pp. 8-15.
- Yao, D.D. (1984a), "Majorization and arrangement orderings in open networks of queues", technical report, Department of Industrial Engineering, Columbia University, New York.
- Yao, D.D. (1984b), "Some properties of the throughput function of closed networks of queues", technical report, Department of Industrial Engineering, Columbia University, New York.
- Yao, D.D., and Kim, S.C. (1984), "Some order relations in closed networks of queues with multi-server stations", technical report, Department of Industrial Engineering, Columbia University, New York.