# An Analytical Formula for Variance of Output from a Series-Parallel Production System with No Interstation Buffers and Time-Dependent Failures

B. TAN
Graduate School of Business, Koç University
Çayır Cad. Istinye, Istanbul, 80860, Turkey
btan@ku.edu.tr

**Abstract**—This paper presents a method to determine the mean and the variance of the amount of materials produced in a fixed time interval by a continuous materials flow production system with $N$ stations in series and $M$ stations in parallel and no interstation buffers. Unreliable stations with exponential failure and repair times, time dependent failures, and deterministic processing times are considered. Closed-form expressions for the asymptotic mean and variance of the amount of materials produced per unit time are given for series, parallel, and series-parallel production systems with identical stations. It is shown that the distribution of the amount of materials produced in a fixed time is asymptotically normal. By using this property, effects of variability on the due-date performance are investigated by considering the probability of meeting a customer order on time. Numerical experiments that explore some relationships among performance measures and production system parameters are also presented.

**Keywords**—Variance of the output, Markov models, Production systems, Throughput, Sojourn time, Performance evaluation.

## 1. INTRODUCTION AND PAST WORK

In this study, our scope is limited to production systems that can be modelled as continuous time Markov chains (CTMC). Specifically, continuous materials flow production systems with no interstation buffers and time-dependent failures are considered. Continuous materials flow production systems are commonly used in process industries. Furthermore, discrete material flow production systems can be approximated with continuous materials flow models especially if the processing rates are much smaller than the failure and repair rates [1]. A special interest has been devoted to production systems with no interstation buffers. Such systems are commonly used in industry. Since there are no buffers in between stations, the behavior of each machine is highly dependent on the others [2,3].

There are many studies on the performance modeling of production systems [1]. Most of these studies focus on the mean performance such as the production rate or the average Work-In-Progress inventory levels. The new emerging management methods in manufacturing require designing production systems that can meet customer requirements on time, every time. In other words, these methods require predictable and dependable delivery schedules. Gershwin [4] reports that his informal numerical and simulation experimentation and factory observations indicate that the standard deviation of weekly production can be over 10% of the mean. This

Typeset by $\mathcal{A}_{\mathcal{M}}\mathcal{S}$-TEX

means that it is very likely that customer requirements cannot be met on time most of the time. Therefore, designing production systems today requires a focused study on the variability of the output as well as on the average performance. Furthermore, analytical models that focus on the variability of the output can also be used to operate production systems more effectively. For example, these analytical models help a production manager to decide on production quota in a pull manufacturing environment [5] or to set the due-date according to a prespecified service level [6]. Furthermore, the optimal production quota is set by studying the variability of the output process. This study is motivated by the need for developing analytical models for the variance of the output that are to be used in design and control of manufacturing systems.

Tan [7] presents a methodology that yields a closed-form expression for the asymptotic variance of the amount of materials produced per unit time of a production line with $N$-stations with exponential failure and repair times, no interstation buffers, and time dependent failures. This methodology is based on determining the limiting mean and variance of total sojourn time in a specific state of a continuous time Markov chain by using the transient probability functions. The methodology presented in [7] does not allow the rate matrix of a CTMC to have multiplicity of eigenvalues. Therefore this result cannot be used to analyze general CTMCs. This study is built on [7] and extends it in two dimensions. First, we extend the result presented in [7] to general CTMCs with multiplicity of eigenvalues. This allows us to analyze any manufacturing system that can be modelled as a CTMC. The covariance of the total sojourn times in two specific states of a CTMC is also determined. This result can also be used in other applications where the limiting expectation, variance, and covariance of the total sojourn time in a specific state of a CTMC are of interest. Next, we extend the closed-form expressions for the mean and variance of the output from a series production line to also parallel, and series-parallel arrangement of workstations.

The number of studies on estimating the variance of the output process of a production system is limited. Miltenburg [8] considers discrete parameter, discrete state space Markov chain models of discrete material flow production lines. Hendrics [9] develops a technique to analytically describe the output process of a serial production line with $N$ reliable machines with exponential processing time distributions and finite buffer capacities. Hendrics and McClain [10] extended the results for exponential processing time distribution to general processing time distribution. Gershwin [4] develops a method for calculating the variance of the number of parts produced by a transfer line during a fixed time interval. The method yields an exact formula for a single machine and uses decomposition for two-machine and longer lines. Carrascosa [11] uses the analytical results presented in [4] and simulation to explore the variability of output produced by a deterministic two-machine finite buffer line with unreliable machines. Duenyas and Hopp [12] present an estimate of the variance of the output of a single closed loop queueing system with exponential servers. Duenyas *et al.* [5] extend this approach to characterize the output process of a CONWIP line with unreliable servers with deterministic processing times and exponentially distributed repairs and failures. Papadopoulos [2] presents an algorithm for calculating the mean sojourn time of a $K$-station production line with no interstation buffers, exponential service times, manufacturing blocking. Grassman [13] presents a method to determine the mean and variance of time averages in Markovian systems. This method can also be used to determine the mean and variance of the output per unit time from a production system that can be modeled as a CTMC.

We are not aware of any study on variance of the output from a continuous-flow series-parallel production system. Furthermore, unlike other studies that present numerical techniques to determine the variance of the output, we present closed-form expressions for the mean and variance of the output from a series-parallel production system for the first time. This is the main contribution of this study.

The outline of the remaining part of this paper is as follows. In Section 2, an overview of our approach is presented. In Section 3, total up time of a system and its relation to the output

from a production system are discussed. The main result that yields an analytical formula for the mean and variance of the total up time in recurrent continuous time Markov chains is given in Section 4. In Section 5, different measures to evaluate the performance of a production system are defined. The closed-form expressions for the asymptotic mean and variance of the amount of materials produced per unit time of series, parallel, and series-parallel production systems are given in Section 6. Numerical experiments that investigate the relationship among the system parameters and the performance measures are presented in Section 7. Finally, conclusions are given in Section 8.

## 2. OVERVIEW OF APPROACH

The most important performance measure of a production system is the amount of materials processed per unit time in the long run. Let $p$ (volume/time) be the rate at which the materials are processed when a station is up and operating. We assume that $p$ is deterministically given. All the variability in the system is due to the unreliability of stations. Depending on whether the stations are up or down, the production system will be producing at a rate of $p$ or it does not produce at all. Let $S_{UP}(t)$ be the total operational time of the system during the interval $[0, t)$. Thus, the amount of materials processed during the interval $[0, t)$ is $S_{UP}(t)p$. Therefore the output per unit time $\pi$ can be written as

$$\pi = \lim_{t \to \infty} \frac{S_{UP}(t)p}{t}. \tag{1}$$

Due to the randomness in the production line, $S_{UP}(t)$, and therefore also the amount of materials processed during the period $[0, t)$, are continuous random variables. The expected output per unit time is called the throughput or the production rate.

Since $p$ is given, determining the mean and variance of $S_{UP}(t)$ yields the mean and variance of the output per unit time. We first define the total operating time during an interval depending on the total sojourn times in the operational states of the underlying Markov chain. Then we present a result that yields closed-form expressions for the asymptotic mean and variance of the total sojourn times in a specific state of a CTMC in terms of its transient state probabilities. This result allows one to determine the closed-form expressions for the asymptotic mean and variance of the amount of materials produced per unit time of a production system from the closed-form expressions of the probability that the system is operational at time $t$ given that it was initially operational. This probability is the instantaneous availability of the system [14]. Thus the methodology presented here can be used to determine the asymptotic mean and variance of the amount of materials produced per unit time of a production system in terms of its instantaneous availability function.

## 3. TOTAL UP TIME OF A SYSTEM

Consider a production system that can be modeled as a finite state homogeneous CTMC. We assume that all states constitute a single communicating class, thus it does not have any absorbing states. Furthermore, the process is ergodic. Therefore all the states are recurrent. Note that if there exists an absorbing state, then the distribution of the total time in the operating states can easily be obtained by using a first passage analysis from the transient states to the absorbing states [15].

Let $\{Y(t), \ t \geq 0\}$ be the CTMC of the production system on the state space $Z = \{0, 1, 2, \ldots, N-1\}$. Let us assume that $Z$ is separable into two sets. Let $\underline{U}$ be the set of states in $Z$ where the system is operational (up), and let $\underline{D}$ be the set of states in $Z$ where the system is not operational (down). Let $X(t)$ be defined to be 1 if $Y(t)$ is within $\underline{U}$, and 0 otherwise. Then, the total up time of the system is

$$S_{UP}(t) = \int_0^t X(\tau) \, d\tau. \tag{2}$$

The rate at which the materials are processed $p$ can be taken as one without loss of generality. In this case, $p$ defines the time unit in the model, and the failure and repair times are scaled accordingly. Thus, we have

$$\pi = \lim_{t\to\infty} \frac{S_{\text{UP}}(t)}{t} = \lim_{t\to\infty} \frac{1}{t} \int_0^t X(\tau)\, d\tau. \tag{3}$$

Note that $\pi$ is the time average of a specific function of $Y(t)$. Under a wide variety of conditions, time averages of functions defined in Markovian environments are proven to be asymptotically normal [16]. For a discussion of asymptotic normality of time averages, see [13]. Following this argument, $S_{\text{UP}}(t)$ is also asymptotically normal. Therefore the mean and variance of $S_{\text{UP}}(t)$ determine its asymptotic distribution. This distribution can be used to derive other performance measures related to probabilities of certain events such as meeting a customer order on time.

Figure 1 depicts a sample realization of the process $\{Y(t),\ t \geq 0\}$ and $X(t)$. In this specific system, there are five states in the state space, $Z = \{0, 1, 2, 3, 4\}$. The system is up in states 2, 3, and 4, i.e., $\underline{U} = \{2, 3, 4\}$, and down in states 0 and 1, i.e., $\underline{D} = \{0, 1\}$. In this figure, $T_i$ denotes the sojourn time in state $i$.
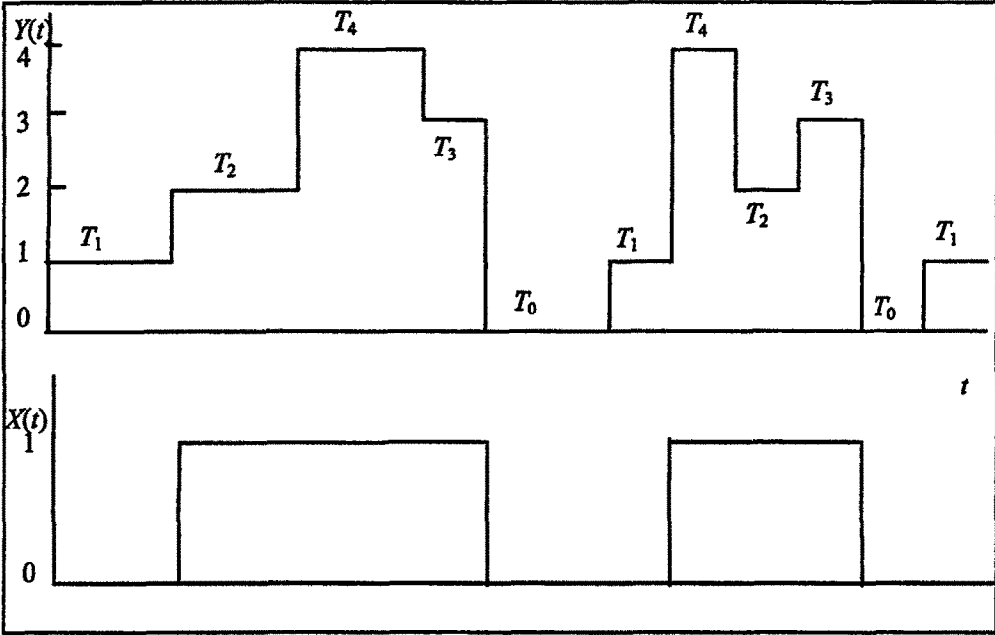


Figure 1. A sample realization of the process $\{Y(t),\ t \geq 0\}$ and $\{X(t),\ t \geq 0\}$ for a specific system with a state space $Z = \{0, 1, 2, 3, 4\}$ and $\underline{U} = \{2, 3, 4\}$, $\underline{D} = \{0, 1\}$.

Let $S_j(t)$ be a random variable that is the total residence time in state $j$ during the period $[0, t)$. Let $E_i[S_j(t)]$ and $\text{Var}_i[S_j(t)]$ be the expectation and the variance of this random variable given the initial state $Y(0) = i$. Similarly let $\text{Covar}_i[S_i(t), S_j(t)]$ be the covariance of the total residence times in states $i$ and $j$ given the initial state $Y(0) = i$. The total up time is the sum of total sojourn times in the states where the system is operational. Considering the specific system depicted in Figure 1, the total up time for this system is the sum of the total sojourn times in states 2, 3, and 4, i.e., $S_{\text{UP}}(t) = S_2(t) + S_3(t) + S_4(t)$. Therefore, the expectation of the total up time $E[S_{\text{UP}}(t)]$ is

$$E[S_{\text{UP}}(t)] = \sum_{i \in \underline{U}} E[S_i(t)].$$

Since the total sojourn times are not independent, the variance of the total up time $\text{Var}[S_{\text{UP}}(t)]$ is

$$\text{Var}[S_{\text{UP}}(t)] = \sum_{i \in \underline{U}} \sum_{j \in \underline{U}} \text{Covar}[S_i(t), S_j(t)]. \tag{4}$$

When we consider the time averages, we obtain

$$m = \lim_{t \to \infty} \frac{E[S_{\text{UP}}(t)]}{t} = \lim_{t \to \infty} \sum_{i \in \underline{U}} \frac{E[S_i(t)]}{t}, \tag{5}$$

$$\sigma^2 = \lim_{t \in \infty} \frac{\text{Var}[S_{\text{UP}}(t)]}{t} = \sum_{i \in \underline{U}} \sum_{j \in \underline{U}} \lim_{t \to \infty} \frac{\text{Covar}[S_i(t), S_j(t)]}{t}, \tag{6}$$

where $m$ and $\sigma^2$ are called the limiting mean and variance of $S_{\text{UP}}(t)$. Next we present a method to determine the limiting expectation and covariance of the total sojourn times from the transient probabilities. That is, we determine the right-hand sides of equations (5) and (6) to determine $m$ and $\sigma^2$.

# 4. ASYMPTOTIC MEAN AND VARIANCE OF TOTAL SOJOURN TIMES

## 4.1. Model

Let $Q = \{q_{ij}\}$ be the infinitesimal generator of the rate matrix of the process $\{Y(t), \ t \geq 0\}$ defined on the state space $Z = \{0, 1, 2, \ldots, N-1\}$. Let $\mathbf{P}(t)$ be the row vector of state probabilities where its $i^{\text{th}}$ component is the probability that the state of the system is $i$ at time $t$, i.e., $P_i(t) = P[Y(t) = i]$. Given the initial state probability vector $\mathbf{P}(0), \mathbf{P}(t)$ satisfies the Kolmogorov differential equations:

$$\frac{d}{dt}\mathbf{P}(t) = \mathbf{P}(t)Q. \tag{7}$$

The solution of the Kolmogorov differential equation given in equation (7) is in the form of a matrix exponential given by

$$\mathbf{P}(t) = \mathbf{P}(0)e^{Qt}. \tag{8}$$

We focus on the spectral properties of the rate matrix $Q$. Let us assume that the rate matrix $Q$ has $K$ distinct eigenvalues and the $k^{\text{th}}$ eigenvalue has a multiplicity of $d_k$. Note that the total number of eigenvalues of $Q$, including the multiplicities, is equal to the size of the state space $N$, i.e., $\sum_{k=1}^{K} d_k = N$.

Let $-\xi_1, -\xi_2, \ldots, -\xi_k$ be the distinct eigenvalues of $Q$ arranged in nondecreasing order of magnitude. Since $Q$ is singular and it is the rate matrix, $\xi_1 = 0$, and all the eigenvalues are nonpositive. Then it can be shown that any time-dependent state probability, say $P_{i,j}(t) = P[Y(t) = j \mid Y(0) = i]$ can be written as

$$P_{i,j}(t) = \sum_{k=1}^{K} \sum_{m=1}^{d_k} a_{k,m} t^{m-1} e^{-\xi_k t}, \tag{9}$$

where $a_{k,m}$ $(k = 1, 2, \ldots, K; m = 1, 2, \ldots, d_k)$ are scalars. If the failure and repair mechanisms of components in a multicomponent system are statistically and structurally independent, then it is possible to determine the time-dependent state probabilities without using the solution of the Markov model. That is, the scalars $a_{k,m}$ can be obtained from the closed-form expressions. On the other hand, if the time dependent probabilities cannot be determined directly, semi-symbolic solution of a CTMC can be used first to determine the scalars $a_{k,m}$. For efficient algorithms to obtain semi-symbolic solution of a CTMC, the reader is referred to [17]. If $\{Y(t), \ t \geq 0\}$

is an irreducible CTMC with one communicating class, then the eigenvalue with value 0 has a multiplicity of one, i.e., $d_1 = 1$. Then,

$$P_{i,j}(t) = a_{1,1} + \sum_{k=2}^{K} \sum_{m=1}^{d_k} a_{k,m} t^{m-1} e^{-\xi_k t}, \tag{10}$$

$\xi_k > 0$, $k > 1$. Note that if the eigenvalue with value 0 has a multiplicity of more than one, then the process does not have a single communicating class. In this case, the distribution of the total up time can be determined by analyzing each communicating class separately.

### 4.2. The Limiting Mean, Variance, and Covariance of Total Sojourn Times" for a Recurrent Continuous Time Markov Chain

Let $P_{ii}(t)$ be the probability that the process is in state $i$ at time $t$ given the initial state $i$ is given as

$$P_{ii}(t) = a_{1,1} + \sum_{k=2}^{K} \sum_{m=1}^{d_m} a_{k,m} t^{m-1} e^{-\xi_k t}. \tag{11}$$

Let $P_{ij}(t)$ be the probability that the process is in state $j$ at time $t$ given the initial state $i$ is given as

$$P_{ij}(t) = b_{1,1} + \sum_{k=2}^{K} \sum_{m=1}^{d_k} b_{k,m} t^{m-1} e^{-\xi_k t}, \tag{12}$$

and similarly, let $P_{ij}(t)$ be the probability that the process is in state $i$ at time $t$ given the initial state $j$ is given as

$$P_{ji}(t) = a_{1,1} + \sum_{k=2}^{K} \sum_{m=1}^{d_k} c_{k,m} t^{m-1} e^{-\xi_k t}. \tag{13}$$

Note that since the limiting results are independent of the initial state, $i$ and $j$, $a_{1,1}$ appears as the first term in $P_{ii}(t)$ and $P_{ji}(t)$, and $b_{1,1}$ appears as the first term in $P_{ij}(t)$.

Now, define an indicator variable $I_i(t)$ which is 1 if the process is in state $i$ at time $t$ and 0 otherwise. Then the total time spent in state $i$, i.e., total sojourn time in state $i$, during the interval $[0,T)$ $S_i(T)$ is given by the following stochastic integral:

$$S_i(T) = \int_0^T I_i(\tau) \, d\tau. \tag{14}$$

**Limiting expectation**

Let $E_k[S_i(T)]$ be the expectation of $S_i(T)$ given the initial state $k$. Since $I_i(t)$ is nonnegative, the order of the integral and the expectation in the above equation can be interchanged by using Fubini's theorem which yields

$$E_k[S_i(T)] = E_k \left[ \int_0^T I_i(t) \, dt \right] = \int_0^T E_k[I_i(t)] \, dt. \tag{15}$$

Since $E_k[I_i(t)] = P[Y(t) = i \mid Y(0) = k]$, we obtain

$$E_k[S_i(T)] = \int_0^T P[Y(t) = i \mid Y(0) = k] \, dt = \int_0^T P_{k,i}(t) \, dt. \tag{16}$$

Since the limiting result is independent of the initial state, setting $k = i$ and inserting $P_{i,i}(t)$ given in equation (11) into equation (16) yields

$$E_i[S_i(T)] = \int_0^T \left( a_{1,1} + \sum_{k=2}^{K} \sum_{m=1}^{d_k} a_{k,m} t^{m-1} e^{-\xi_k t} \right) dt$$

$$= a_{1,1}T + \sum_{k=2}^{K} \sum_{m=1}^{d_k} \left[ \frac{a_{k,m}}{\xi_k^m} (m-1)! - \sum_{n=0}^{m-1} \binom{m-1}{n} \frac{T^{m-1-n}}{\xi_k^{n+1}} e^{-\xi_k T} \right], \tag{17}$$

where $\binom{M}{N} = M!/N!(M-N)!$

Dividing equation (17) by $T$ and taking the limit as $T \to \infty$ immediately give the limiting expectation $E_i$ of the total residence time in state $i$ given the initial state $i$:

$$E_i = \lim_{T \to \infty} \frac{E_i[S_i(T)]}{T} = a_{1,1}. \tag{18}$$

Note that the limiting result is independent of the initial state. Since the process is ergodic, the limiting mean of the fraction of time the process is in $i$, $E_i$, is also the steady-state probability that the process is in state $i$, i.e.,

$$E_i = \lim_{t \to \infty} P_{k,i}(t) = a_{1,1}.$$

**Limiting covariance and variance**

The covariance of $S_i(T)$ and $S_j(T)$ given the initial state $k$ is given as

$$\text{Covar}_k[S_i(T), S_j(T)] = E_k[S_i(T)S_j(T)] - E_k[S_i(T)]E_k[S_j(T)]. \tag{19}$$

The first term on the right-hand side of the above equation is evaluated as

$$E_k[S_i(T)S_j(T)] = E_k\left[ \int_0^T I_i(t)\, dt \int_0^T I_j(t)\, dt \right] = \int_0^T \int_0^T E_k[I_i(s)I_j(t)]\, ds\, dt. \tag{20}$$

Since $E_k[I_i(s)I_j(t)] = P[Y(s) = i, Y(t) = j \mid T(0) = k]$, we obtain

$$\begin{aligned}
E_k[S_i(T)S_j(T)] &= \int_0^T \int_0^T P[Y(s) = i, Y(t) = j \mid Y(0) = k]\, ds\, dt \\
&= \iint_{s<t} P_{ij}(t-s)P_{ki}(s)\, ds\, dt + \iint_{t<s} P_{ji}(s-t)P_{kj}(t)\, ds\, dt \\
&= \int_0^T \int_0^{T-s} P_{ki}(s)P_{ij}(u)\, du\, ds + \int_0^T \int_0^{T-t} P_{kj}(t)P_{ji}(u)\, du\, dt.
\end{aligned} \tag{21}$$

Similarly, the second term of the right-hand side of equation (19) is evaluated as

$$E_k[S_i(T)]E_k[S_j(t)] = \int_0^T P_{ki}(t)\, dt \int_0^T P_{kj}(t)\, dt. \tag{22}$$

Thus, the covariance of $S_j(t)$ and $S_j(t)$ given the initial state $k$ is

$$\begin{aligned}
\text{Covar}_k[S_i(T), S_j(T)] = &\int_0^T \int_0^{T-s} P_{ki}(s)P_{ij}(u)\, du\, ds + \int_0^T \int_0^{T-t} P_{kj}(t)P_{ji}(u)\, du\, dt \\
&- \int_0^T P_{ki}(t)\, dt \int_0^T P_{kj}(t)\, dt.
\end{aligned} \tag{23}$$

Since the asymptotic results are independent of the initial state, $k$ can be set to $i$ in equation (23). Now, inserting equations (11), (12), and (13) into equation (23), evaluating the integral, dividing by $T$, and taking the limit as $T \to \infty$ yield the limiting covariance $C_{ij}$ of the total residence in states $i$ and $j$

$$C_{ij} = \lim_{T \to \infty} \frac{\text{Covar}_i[S_i(T), S_j(T)]}{T} = a_{1,1} \sum_{k=2}^{K} \sum_{m=1}^{d_k} \frac{b_{k,m}}{\xi_k^m}(m-1)! + b_{1,1} \sum_{k=2}^{K} \sum_{m=1}^{d_k} \frac{c_{k,m}}{\xi_k^m}(m-1)! \tag{24}$$

Since $\text{Var}_i[S_i(T)] = \text{Covar}_i[S_i(T), S_i(T)]$, equation (24) yields

$$V_i = \lim_{T \to \infty} \frac{\text{Var}_i[S_i(T)]}{T} = 2a_{1,1} \sum_{k=2}^{K} \sum_{m=1}^{d_k} \frac{a_{k,m}}{\xi_k^m}(m-1)! \tag{25}$$

# 5. PERFORMANCE MEASURES FOR CONTINUOUS MATERIAL FLOW PRODUCTION SYSTEMS

Once the limiting means and covariances of total sojourn times are obtained, equations (5) and (6) yield the limiting mean and variance of the total up time, $m$ and $\sigma^2$. Therefore, the mean and variance of the output from a production system per unit time, $E[\pi]$ and $\text{Var}[\pi]$, are simply given as

$$E[\pi] = \lim_{t \to \infty} \frac{E_i[S_{\text{UP}}(t)]}{t} p = mp, \tag{26}$$

$$\text{Var}[\pi] = \lim_{t \to \infty} \frac{\text{Var}_i[S_{\text{UP}}(t)]}{t} p^2 = \sigma^2 p^2. \tag{27}$$

Similarly, the amount of materials processed in a production system during a period of length $t$ is normally distributed with mean $E[\pi]t$ and variance $\text{Var}[\pi]t$ as $t$ approaches infinity. The asymptotic distribution of the amount of materials processed can also be used to derive other performance measures that depend on the probabilities of certain events such as the probability of meeting a customer order on time. Now consider the probability of meeting a customer order of $V^*$ volume of materials by time $T$. If the amount of materials produced in $T$ time unit $V$ is greater than or equal to $V^*$, then the customer order is said to be met on time. Then the probability of meeting the customer order on time is

$$P[V \geq V^*] = P\left[z \geq \frac{V^* - E[\pi]T}{\sqrt{\text{Var}[\pi]T}}\right] = 1 - \Phi\left(\frac{V^* - E[\pi]T}{\sqrt{\text{Var}[\pi]T}}\right), \tag{28}$$

where $\Phi(.)$ is the cumulative normal distribution function given by

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{z} e^{-(1/2)t^2} dt, \tag{29}$$

and $E[\pi]$ and $\text{Var}[\pi]$ are given in equations (26) and (27). We use the mean and variance of the output and also the probability of meeting a customer order on time as the main measures to evaluate the performance of a production system.

# 6. PERFORMANCE OF SERIES, PARALLEL, SERIES-PARALLEL PRODUCTION SYSTEMS WITH NO INTERSTATION BUFFERS

In this section, the general result presented in the preceding section is used to derive the closed-from expressions for the mean and variance of the output of series, parallel, series-parallel production systems per unit time.

## 6.1. General Assumptions

Consider a multistation production system with no interstation buffers. The workstations are unreliable. The failure time and the repair time of station $i$ are assumed to be exponentially distributed with rates $\lambda_i$ and $\mu_i$, respectively. Failure of a station is assumed to be time dependent. Thus, the production system operates in hot standby. The failure and repair times of a station are independent of the failure and repair times of other stations. Materials flow continuously in the production system. The rate at which the materials are processed in the line is taken as 1 volume/time without loss of generality. All the stations in the system have the same processing rate. The time unit used in the model can be arranged according to the processing rate of the stations such that this rate is one and the failure and repair rates are scaled up or down accordingly. There is always enough material waiting to be processed in the input to the system. That is, the first station is never starved. The output of the system is never blocked. Thus, materials that finish processing in the production system immediately leave the system without any interruption. In addition to the assumptions given here, the series, parallel, and series-parallel models are described separately below.

## 6.2. Series System

In the series system, materials are processed sequentially starting at the first station and leaving the system after being processed at the last station. A production line is basically a series system. Figure 2 depicts a production line with $N$ workstations.
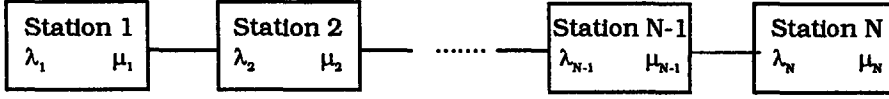


Figure 2. A production line with $N$ unreliable stations and no interstation buffers.

Since there are no interstation buffers in the line, whenever a workstation is down, it halts the production of the line. In other words, a series system is productive when all the stations in the line are up and operating. Since the failures are time dependent and there are no interstation buffers, the failure mechanisms of the stations are statistically as well as structurally independent. Therefore, each station can be analyzed in isolation of other stations in the system. Thus, the instantaneous availability of the series system with $N$ unreliable stations, $P^S_{\mathrm{UP}}(t)$ is

$$P^S_{\mathrm{UP}}(t) = \prod_{n=1}^{N} P[X_n(t) = 1 \mid X_n(0) = 1] = \prod_{n=1}^{N} \left[ \frac{\mu_n}{\lambda_n + \mu_n} + \frac{\lambda_n}{\lambda_n + \mu_n} e^{-(\lambda_n + \mu_n)t} \right]. \tag{30}$$

The equation given above can be expanded for a given $N$ to write in the form given in equation (15). If the failure and repair rates are different, the final expression is extremely lengthy to be of any practical use, and thus it is not given here. When the failure and repair rates are the same for all the stations, i.e., $\lambda_i = \lambda$, $\mu_i = \mu$, $i = 1, \ldots, N$,

$$P^S_{\mathrm{UP}}(t) = \left( \frac{\mu}{\lambda + \mu} + \frac{\lambda}{\lambda + \mu} e^{-(\lambda+\mu)t} \right)^N = \frac{\mu^N}{(\lambda+\mu)^N} + \sum_{k=1}^{N} \binom{N}{k} \frac{\mu^{N-k}\lambda^k e^{-(\lambda+\mu)kt}}{(\lambda+\mu)^N}. \tag{31}$$

Since the resulting expression is in the form of equation (12), equations (26) and (27) yield the expectation and variance of the output per unit time of a series system with $N$ identical stations, $E[\pi^S_N]$ and $\mathrm{Var}[\pi^S_N]$, respectively, given as

$$E\left[\pi^S_N\right] = \frac{\mu^N}{(\lambda + \mu)^N}, \tag{32}$$

$$\mathrm{Var}\left[\pi^S_N\right] = \frac{2\mu^{2N}}{(\lambda + \mu)^{2N+1}} \sum_{k=1}^{N} \binom{N}{k} \frac{(\lambda/\mu)^k}{k}. \tag{33}$$

## 6.3. Parallel System

In the parallel arrangement of workstations, stations perform the same operation. Thus materials leave the system after being processed in any of the stations. There are two different cases. First, we can assume that materials are processed in only one station at a rate of 1 volume/time, and all the other stations are in the hot standby mode waiting to replace the station in the case a station breakdown. Alternatively, we can assume that all the operating stations process the materials together in such a way that the combined processing rate of the system is still 1 volume/time. That is, if there are $k$ operating stations out of $M$ parallel stations at a given time, then the processing rates of these stations are slowed down from 1 volume/time to $1/k$ volume/time. Thus the processing rate of the parallel system does not change as the number of operating stations changes in time. For a discussion of this assumption, the reader is referred to [18, p. 239]. Since the failures are time dependent and the failure and repair times
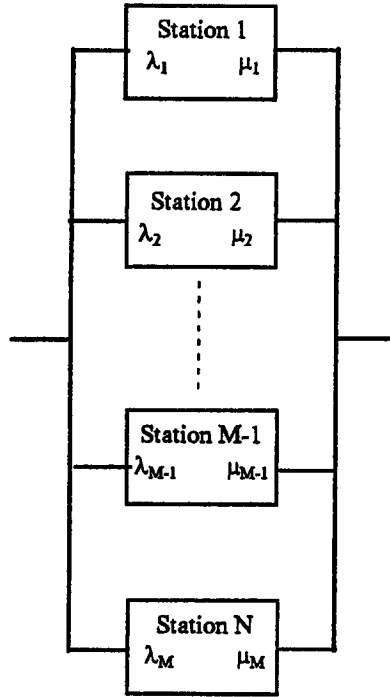
Figure 3. A production system with $M$ unreliable stations and no buffers.

are exponential, these two different cases are identical. Figure 3 depicts a parallel system with $M$ unreliable stations and no buffers.

A parallel system is productive when at least one of the stations is up and operating. Similar to the series case, time dependent failure mechanism assumption and other assumptions allow one to analyze a given station independent of the other stations. Thus, the instantaneous availability of a parallel system with $M$ stations $\underline{P}_{\mathrm{UP}}^{P}(t)$ can be written as

$$
\begin{aligned}
\underline{P}_{\mathrm{UP}}^{P}(t) &= 1 - \prod_{m=1}^{M} \left( 1 - P[X_m(t) = 1 \mid X_m(0) = 1] \right) \\
&= 1 - \prod_{m=1}^{M} \left[ \frac{\lambda_m}{\lambda_m + \mu_m} - \frac{\lambda_m}{\lambda_m + \mu_m} e^{-(\lambda_m + \mu_m)t} \right].
\end{aligned}
\tag{34}
$$

Similar to the series case, equation (34) can be expanded for a given $M$ to write in the form in equation (15). When the failure and repair rates are the same for all the stations, i.e., $\lambda_i = \lambda$, $\mu_i = \mu$, $i = 1, \ldots, M$,

$$
\begin{aligned}
\underline{P}_{\mathrm{UP}}^{P}(t) &= 1 - \left[ \frac{\lambda}{\lambda + \mu} - \frac{\lambda}{\lambda + \mu} e^{-(\lambda+\mu)t} \right]^{M} \\
&= \frac{(\lambda + \mu)^{M} - \lambda^{M}}{(\lambda + \mu)^{M}} + \sum_{m=1}^{M} \binom{M}{m} (-1)^{m+1} \frac{\lambda^{M}}{(\lambda + \mu)^{M}} e^{-m(\lambda+\mu)t}.
\end{aligned}
\tag{35}
$$

Thus, the expectation $E[\pi_M^P]$ and variance $\mathrm{Var}[\pi_M^P]$ of the output per unit time of a parallel system with $M$ identical stations follow equations (32) and (33):

$$
E[\pi_M^P] = \frac{(\lambda + \mu)^{M} - \lambda^{M}}{(\lambda + \mu)^{M}},
\tag{36}
$$

$$
\mathrm{Var}[\pi_M^P] = 2 \frac{\lambda^{M}[(\lambda + \mu)^{M} - \lambda^{M}]}{(\lambda + \mu)^{2M+1}} \sum_{m=1}^{M} \binom{M}{m} \frac{(-1)^{m+1}}{m}.
\tag{37}
$$

## 6.4. Series-Parallel System

In the series-parallel arrangement, the series and the parallel systems given above are placed in series with each other in the same production system. Figure 4 shows a series-parallel system with $M$ parallel stations in series with $N$ stations.
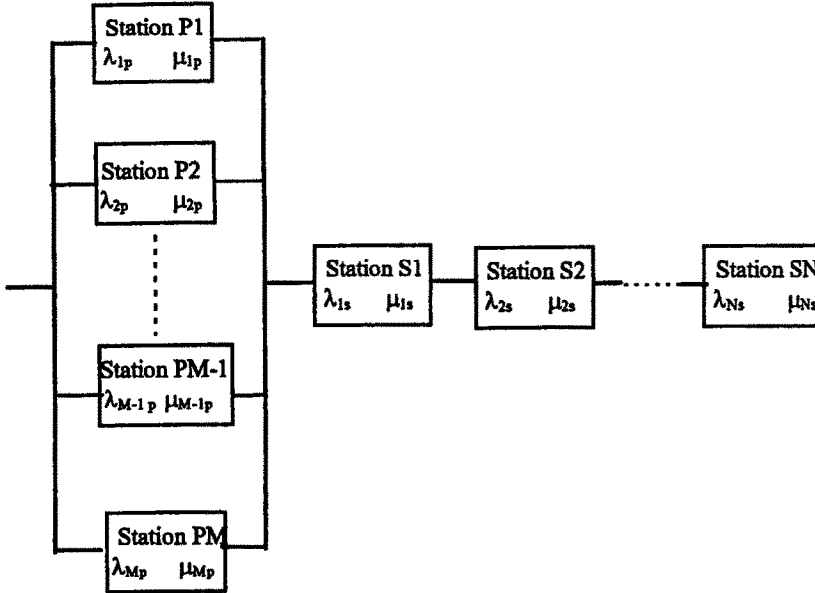


Figure 4. A production system with $M$ unreliable parallel stations in series with $N$ unreliable stations and no interstation buffers.

There are $N + 1$ different processes. Materials are processed first in the parallel system. Then they are processed sequentially starting at station S1 and finishing at station SN. The failure and repair rates of the parallel stations are $\lambda_{ip}$ and $\mu_{ip}$, $i = 1, \ldots, M$, and the failure and repair rates of the stations in series are $\lambda_{is}$ and $\mu_{is}$, $i = 1, \ldots, N$. As explained above, the combined processing rate of the parallel system is 1 volume/time and this rate is equal to the processing rates of the stations in the series arrangement.

A series-parallel system is operational if the parallel system and the series system are operational. Thus, the instantaneous availability of a series-parallel system $\underline{P}_{\mathrm{UP}}^{SP}(t)$ is the product of the instantaneous availabilities of the series and the parallel systems availability:

$$
\begin{aligned}
\underline{P}_{\mathrm{UP}}^{SP}(t) &= \underline{P}_{\mathrm{UP}}^{S}(t)\underline{P}_{\mathrm{UP}}^{P}(t), \\
&= \left( \prod_{n=1}^{N} \left[ \frac{\mu_{ns}}{\lambda_{ns} + \mu_{ns}} + \frac{\lambda_{ns}}{\lambda_{ns} + \mu_{ns}} e^{-(\lambda_{ns}+\mu_{ns})t} \right] \right) \\
&\quad \times \left( 1 - \prod_{m=1}^{M} \left[ \frac{\lambda_{mp}}{\lambda_{mp} + \mu_{mp}} - \frac{\lambda_{mp}}{\lambda_{mp} + \mu_{mp}} e^{-(\lambda_{mp}+\mu_{mp})t} \right] \right).
\end{aligned}
\tag{38}
$$

When the failure and repair rates are the same for all the stations, i.e., $\lambda_{ns} = \lambda$, $\mu_{ns} = \mu$, $n = 1, \ldots, N$ and $\lambda_{mp} = \lambda$, $\mu_{mp} = \mu$, $m = 1, \ldots, M$,

$$
\begin{aligned}
\underline{P}_{\mathrm{UP}}^{SP}(t) &= \left( \frac{\mu^N}{(\lambda+\mu)^M} + \sum_{k=1}^{N} \binom{N}{k} \frac{\mu^{N-k}\lambda^k e^{-(\lambda+\mu)kt}}{(\lambda+\mu)^N} \right) \\
&\quad \times \left( \frac{(\lambda+\mu)^M - \lambda^M}{(\lambda+\mu)^M} + \sum_{m=1}^{M} \binom{M}{m} (-1)^{m+1} \frac{\lambda^M}{(\lambda+\mu)^M} e^{-m(\lambda+\mu)t} \right),
\end{aligned}
\tag{39}
$$

$$P_{\text{UP}}^{SP}(t) = \frac{\mu^N[(\lambda+\mu)^M - \lambda^M]}{(\lambda+\mu)^{M+N}} + \sum_{m=1}^{M} \binom{M}{m}(-1)^{m+1}\frac{\lambda^M \mu^N}{(\lambda+\mu)^{M+N}}e^{-m(\lambda+\mu)t}$$

$$+ \sum_{k=1}^{N} \binom{N}{k}\frac{\mu^{N-k}\lambda^k\left[(\lambda+\mu)^M - \lambda^M\right]}{(\lambda+\mu)^{M+N}}e^{-(\lambda+\mu)kt} \qquad (40)$$

$$+ \sum_{k=1}^{N}\sum_{m=1}^{M} \binom{M}{m}\binom{N}{k}\frac{(-1)^{m+1}\lambda^{M+k}\mu^{N-k}}{(\lambda+\mu)^{M+N}}e^{-(m+k)(\lambda+\mu)t}.$$

Finally, the expectation $E[\pi_{N,M}^{SP}]$ and variance $\text{Var}[\pi_{N,M}^{SP}]$ of the output per unit time of a homogeneous series-parallel system with $M$ parallel stations in series with $N$ stations are given as

$$E[\pi_{N,M}^{SP}] = \frac{\mu^N}{(\lambda+\mu)^N}\frac{(\lambda+\mu)^M - \lambda^M}{(\lambda+\mu)^M}, \qquad (41)$$

$$\text{Var}[\pi_{N,M}^{SP}] = 2\frac{\lambda^M \mu^{2N}\left[(\lambda+\mu)^M - \lambda^M\right]}{(\lambda+\mu)^{2(M+N)+1}}$$

$$\times \left[\sum_{m=1}^{M}\binom{M}{m}\frac{(-1)^{m+1}}{m} + \sum_{k=1}^{N}\binom{N}{k}\frac{\left((\lambda+\mu)^M - \lambda^M\right)(\lambda/\mu)^k}{\lambda^M k}\right. \qquad (42)$$

$$\left. + \sum_{k=1}^{N}\sum_{m=1}^{M}\binom{M}{m}\binom{N}{k}\frac{(-1)^{m+1}(\lambda/\mu)^k}{m+k}\right].$$

The equations derived for the series-parallel system also give the equations for the series system. It can be shown after some algebraic manipulations that setting $M = 1$ in the above equations yields the equations for a series system with $N + 1$ stations, and as $M$ approaches infinity, they yield the equations for a series system with $N$ stations.

The equations given above can be used to express the mean and variance of the output per unit time of a production system with given numbers of parallel and series stations. Table 1 gives the expectation and variance of the output per unit time of a series-parallel system with different numbers of series and parallel stations. Since the closed-form expressions for the mean and variance are available, coefficient of variation of the output per unit time $CV[\pi]$ can also be given in closed form:

$$CV[\pi] = \frac{\sqrt{\text{Var}[\pi]}}{E[\pi]}. \qquad (43)$$

Table 1. The expected value and variance of the output per unit time for a series-parallel system with $M$ parallel stations and $N$ stations in series.

| Number of Stations in Series $N$ | Number of Parallel Stations $M$ | Throughput $E[\pi]$ | Variance of the Output per Unit Time $\text{Var}[\pi]$ |
|:---:|:---:|:---:|:---:|
| 1 | 1 | $\frac{\mu\lambda}{(\lambda+\mu)^2}$ | $\frac{\mu^2\lambda(4\mu+\lambda)}{(\lambda+\mu)^5}$ |
| 1 | 2 | $\frac{\mu\lambda(\lambda+2\mu)}{(\lambda+\mu)^3}$ | $\frac{1}{3}\frac{\mu^2(2\lambda+\mu)\lambda(21\mu\lambda+6\mu^2+4\lambda^2)}{(\lambda+\mu)^7}$ |
| 2 | 2 | $\frac{\mu^2\lambda(\lambda+2\mu)}{(\lambda+\mu)^4}$ | $\frac{1}{6}\frac{\mu^3(2\lambda+\mu)(72\lambda\mu^2+28\mu\lambda^2+24\mu^3+5\lambda^3)}{(\lambda+\mu)^9}$ |
| 3 | 2 | $\frac{\mu^3\lambda(\lambda+2\mu)}{(\lambda+\mu)^5}$ | $\frac{1}{30}\frac{\mu^4(2\lambda+\mu)\lambda(540\lambda\mu^3+320\lambda^2\mu^2+115\lambda^3\mu+180\mu^4+18\lambda^4)}{(\lambda+\mu)^{11}}$ |

# 7. NUMERICAL EXPERIMENTS ON THE VARIABILITY OF PERFORMANCE

The methodology described in this study is implemented in software. The analytic results obtained for production systems with identical stations yield very efficient computations since the closed-form expressions are available. Although it is possible to obtain closed-form expressions for production systems with nonidentical stations, the analytic results for such systems are extremely lengthy to be of any practical use. Therefore an algorithm that is very similar to the one outlined in [7] is used to obtain performance measures. Numerous experiments are conducted to investigate the interrelationships among different parameters of the production systems and the performance measures. In this study, only the results for production system with identical stations are reported.

In the numerical experiments, effects of the number of stations in series, the number of stations in parallel, and effects of failure and repair time parameters on the mean, variance, coefficient of variation of the output per unit time, and also on the due-date performance are investigated.

Figures 5 and 6 depict the mean and variance of the output of a production system with $N$ ($N = 1, 2, \ldots, 40$) stations in series and $M$ ($M = 1, 2, 3, 4$) stations in parallel. All the stations in the production system are identical with $\lambda = 0.1$ and $\mu = 0.9$. As the number of stations in series increases, the throughput decreases and approaches zero as the number of stations in series approaches infinity. On the other hand, the dependence of the variance of the output per unit time on the number of stations is not monotonic. As $N$ increases, $\text{Var}[\pi]$ may increase first depending on failure and repair parameters and then decreases to zero. This nonmonotonic dependence of $\text{Var}[\pi]$ is discussed in detail in [7].

As the number of stations in parallel increases, the availability of the part of the production system in parallel approaches one, and therefore, throughput is affected only by the number of stations in series. Since variance is affected by the changes in scale, coefficient of variation conveys more information. Figure 7 shows the coefficient of variation of the output per unit time of the same production system depicted in Figures 5 and 6. As the number of stations in series increases, the coefficient of variation of the output per unit time also increases, i.e., the variability of the performance increases.

Figure 8 depicts the interrelationship between the utilization factor of the stations, $\rho = \lambda/\mu$, and the throughput for a production system with $N = 1, 5, 10, 20$ stations in series and $M = 2$ stations in parallel. As $\rho$ increases, the mean number of failures increases more than the mean number of repairs increases. Therefore the total up time and thus the mean throughput decreases. The dependence of $E[\pi]$ on $\rho$ is not affected by $\lambda$ or by $\mu$. That is, $E[\pi]$ is the same for the case when $\lambda = 0.1$ and $\mu = 0.9$ ($\rho = 1/9$) and when $\lambda = 0.2$ and $\mu = 1.8$ ($\rho = 1/9$). On the other hand, the variance of the output per unit time is affected by the values of $\lambda$ and $\mu$ for the same $\rho$. Figure 9 depicts the variance of the output per unit time of a production system with $N = 1, 5, 10, 20$ stations in series and $M = 2$ stations in parallel. All the stations are identical with the same $\rho = 1$. Figure 9 shows that as $\lambda$ increases while $\rho$ stays the same (thus while $\mu$ also increases), the variance of the output per unit time decreases.

Figures 10–12 investigate the due-date performance of the production system depending on the system parameters. Probability of meeting a customer order on time is used as the main performance measure. In all the cases, the probability of meeting a customer order that is 80% of the expected amount of materials that the production system can produce in 100 time units is calculated, i.e., $V^* = 0.8E[\pi]T$ and $T = 100$. Note that with this setting

$$P[V \geq V^*] = P\left[z \geq \frac{V^* - E[\pi]T}{\sqrt{\text{Var}[\pi]T}}\right] = \Phi\left(0.2\sqrt{T}\frac{1}{CV[\pi]}\right). \tag{44}$$

Figure 10 depicts the effects of the number of stations in series and the number of stations in parallel on the due-date performance. As the number of stations in series increases, the
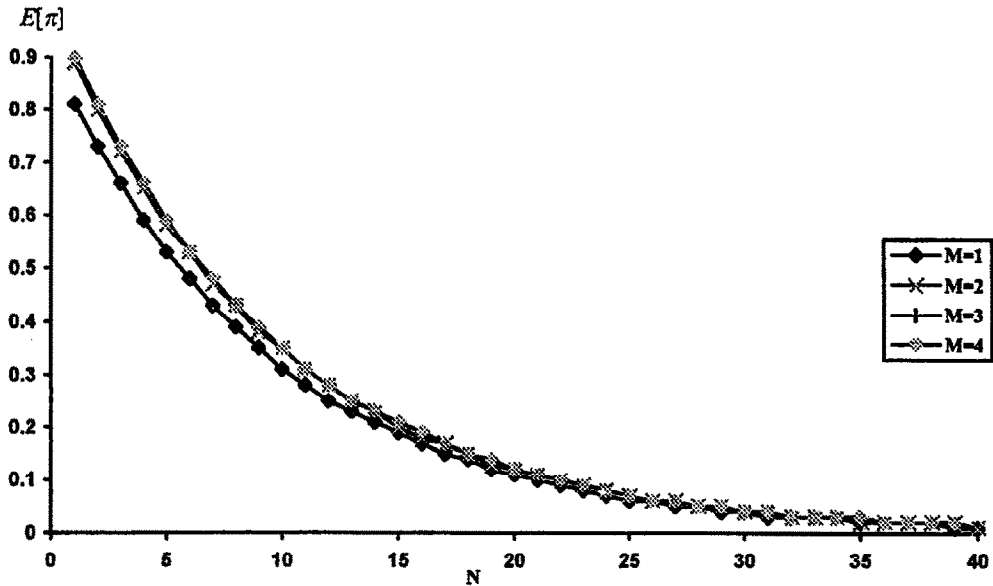
Figure 5. The throughput of a production system with $N$ stations in series and $M$
stations in parallel. All the stations are identical ($\lambda = 0.1$, $\mu = 0.9$).
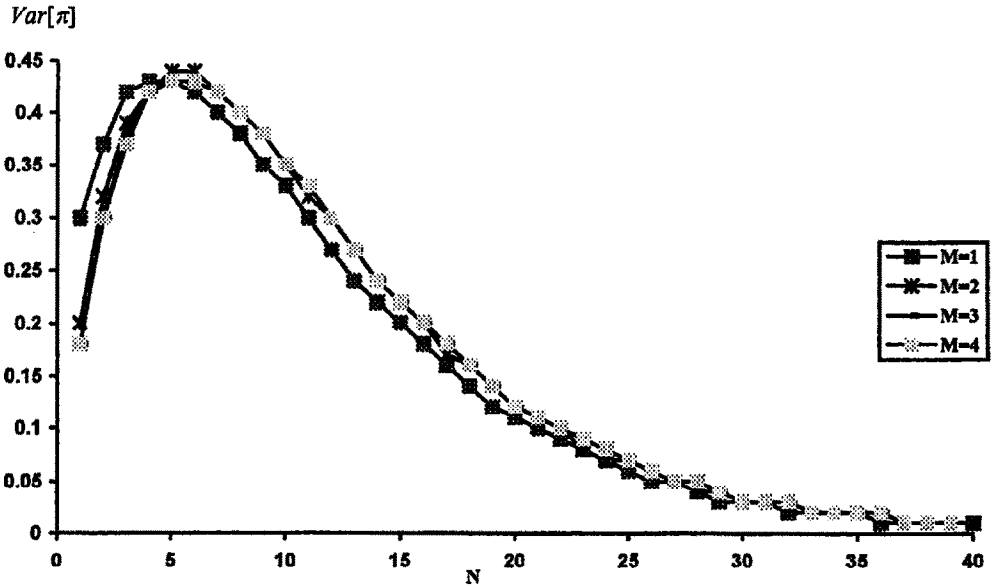


Figure 6. Variance of the output per unit time of a production system with $N$ stations
in series and $M$ stations in parallel. All the stations are identical ($\lambda = 0.1$, $\mu = 0.9$).

due-date performance deteriorates as a result of increased variation in the system, and as the
number of parallel stations increases, the due-date performance slightly improves. The due-date
performance is mainly dictated by $N$ as opposed to $M$.

Figures 11 and 12 depict the effects of the failure and repair rates and the number of stations
on a production system with $N = 1, 2, \ldots, 40$ stations in series and $M = 2$ stations in parallel.
Figure 11 shows that as the failure rate increases, the probability of meeting a customer order on
time decreases. On the contrary, Figure 11 shows that as the repair rate increases, the due-date
performance improves.

Note that all the results depicted in the figures in this section can also be proven by using the
analytical results given in the previous sections. For example, it can be shown from equation (37)
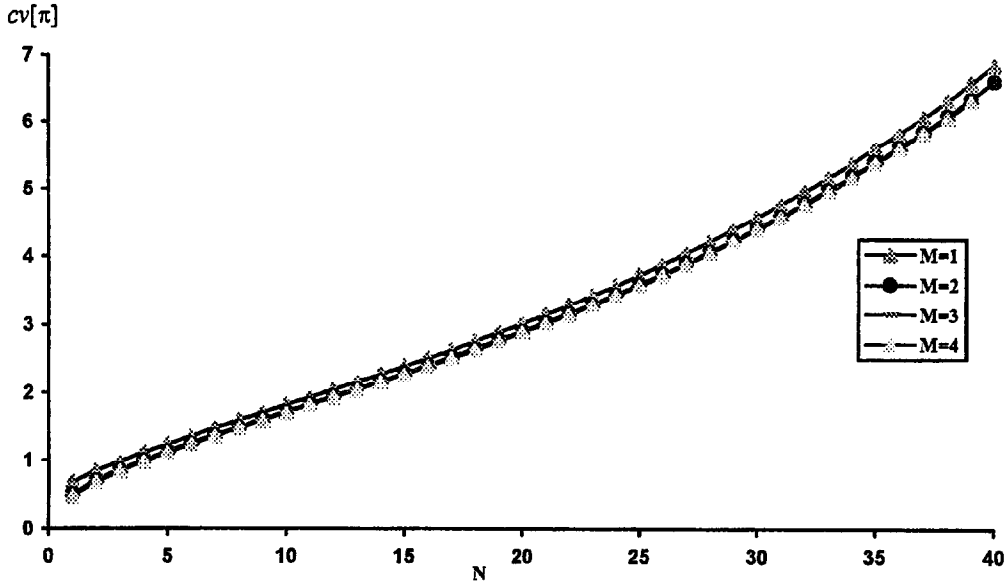
Figure 7. Coefficient of variation of the output per unit time of a production system with $N$ stations in series and $M$ stations in parallel. All the stations are identical ($\lambda = 0.1$, $\mu = 0.9$).
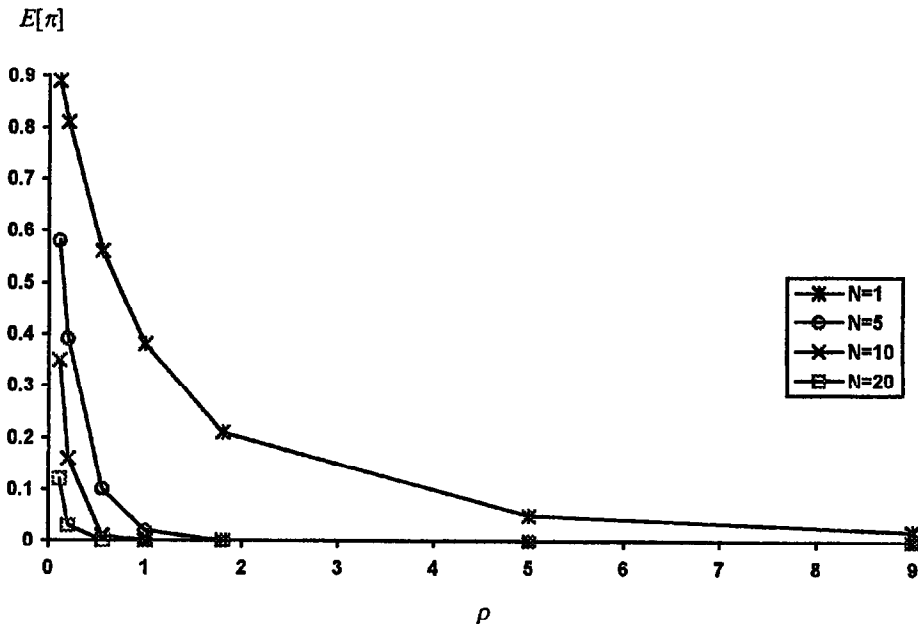


Figure 8. The throughput of a production system with $N$ stations in series and $M = 2$ stations in parallel. All the stations are identical with $\rho = \lambda/\mu$.

that as $N \to \infty$, $P[V^* \geq V] \to 1/2$ as Figures 10–12 suggest. In that respect, the figures convey the same information the analytical results have in a different format, namely, they convey the same information pictorially.

## 8. CONCLUSIONS

In this paper, we present closed-form exact formulae for the variance of the throughput of series, parallel, series-parallel continuous material flow production systems with no interstation buffers and time-dependent failures. The methodology uses a general result derived for continuous time Markov chains. Namely, the limiting mean, variance, and covariance of total sojourn times in
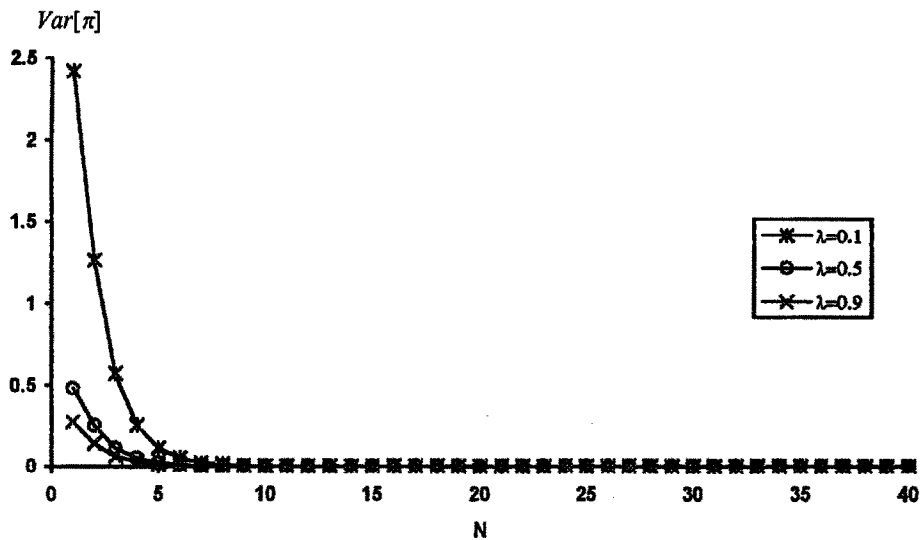
Figure 9. Variance of the output per unit time of a production system with $N$ stations in series and $M = 2$ stations in parallel. All the stations are identical with $\rho = 1$.
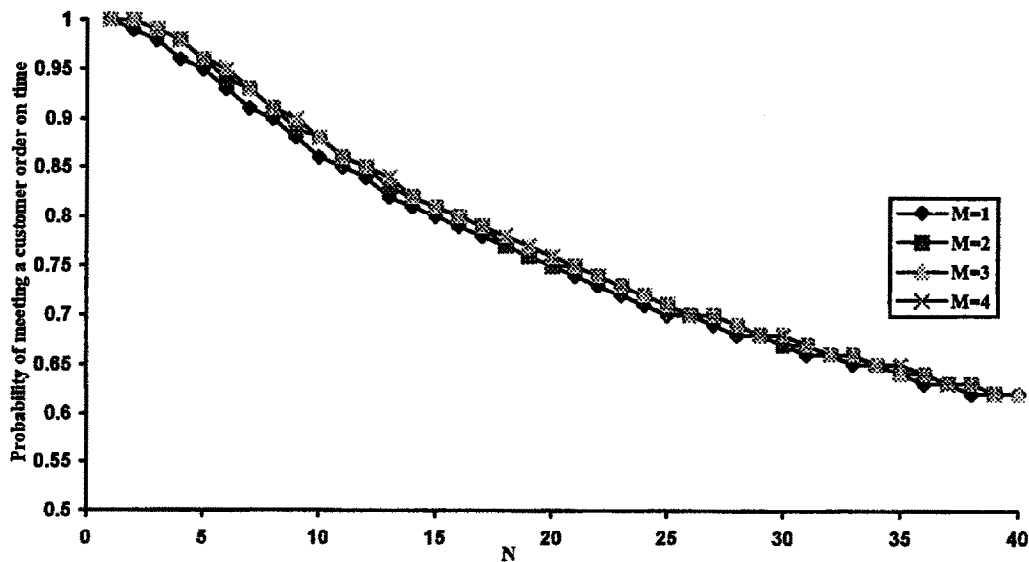


Figure 10. Probability of meeting a customer order of 80% of expected production in $T = 100$ time units ($V^* = 0.8E[\pi]T$) a production system with $N$ stations in series and $M$ stations in parallel. All the stations are identical ($\lambda = 0.1$, $\mu = 0.9$).

specific states of a CTMC can be obtained in closed form from its transient probabilities by using this general result. The analytic result derived in this study determines the asymptotic distribution of the amount of products processed in a fixed time interval. By using this asymptotic distribution, other performance measures such as the probability of meeting a customer order on time can be derived. These results can be used both in the design and also in the control of manufacturing systems.

The numerical results on the due-time performance of the production system are quite intuitive. As the production system gets more complex, or as the stations become more unreliable, the probability of meeting a customer order on time decreases. The corresponding curves can be used either to estimate the due-date performance for a given system or to design a system that meets a given service level, i.e., a given probability of meeting a customer order on time.
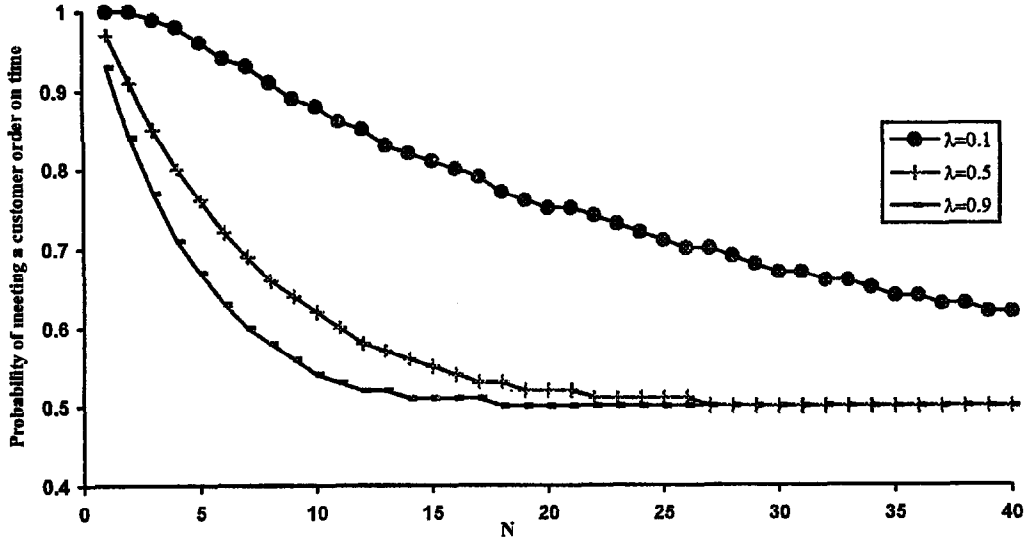
Figure 11. Probability of meeting a customer order of 80% of expected production in $T = 100$ time units $(V^* = 0.8E[\pi]T)$ a production system with $N$ stations in series and $M$ stations in parallel. All the stations are identical ($\mu = 0.9$).
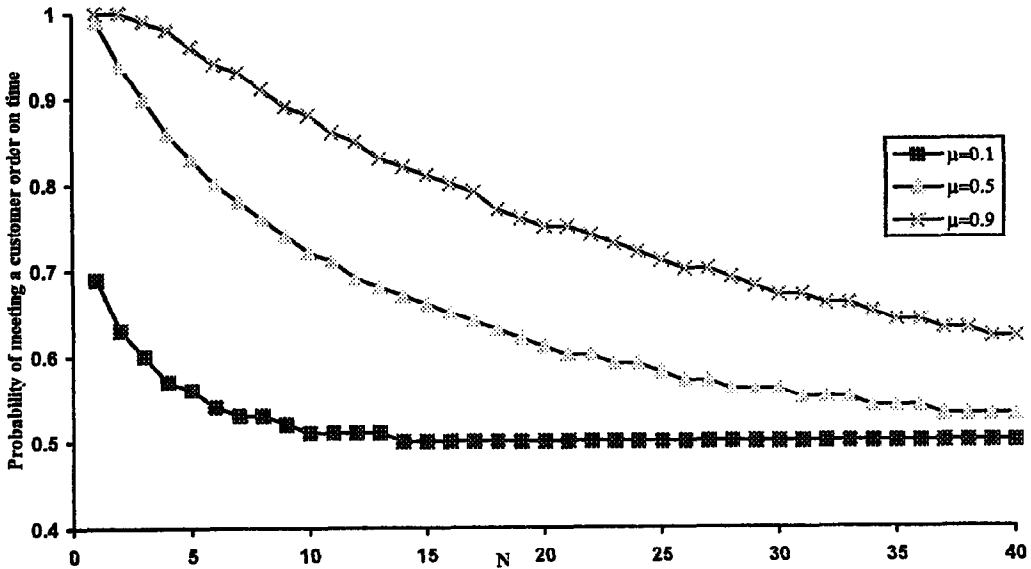


Figure 12. Probability of meeting a customer order of 80% of expected production in $T = 100$ time units $(V^* = 0.8E[\pi]T)$ a production system with $N$ stations in series and $M$ stations in parallel. All the stations are identical ($\lambda = 0.1$).

The system considered in this study is quite simple and the closed-form expressions are obtained under some restrictive assumptions such as time-dependent failures, exponential failure and repair times, etc. We believe even under these restrictive assumptions, a simple analytical model can provide a very useful tool to explore the relationships between parameters and the performance measures. There is a trade-off between the tractability of a model and the range of assumptions incorporated. That is, there exist numerical methods that are very flexible but computationally not efficient, and closed-form expressions for some restrictive cases that are not that flexible but computationally very efficient. This trade-off will be exploited in future research to develop computationally efficient and flexible methods.

# REFERENCES

1. Y. Dallery and S.B. Gershwin, Manufacturing flow line systems: A review of models and analytical results, *Queueing Systems Theory and Applications,* Special Issue on "Queueing Models of Manufacturing Systems" **12** (1–2), 3–94 (1992).
2. H.T. Papadopoulos, An algorithm for calculating the mean sojourn time of $K$-station production lines with no intermediate buffers, Working Paper No. 96-1, Department of Mathematics, University of the Aegean, Samos, Greece, (1996).
3. C. Commault and Y. Dallery, Production rate of transfer lines without buffer storage, *IIE Transactions* **22** (4), 315–329 (1990).
4. S. Gershwin, Variance of the output of a tandem production system, *Queueing Networks with Finite Capacity. Proceedings of the Second International Conference on Queueing Networks with Finite Capacity,* (Edited by R. Onvural and I. Akyıldız), Elsevier, Amsterdam, (1993).
5. I. Duenyas, W.J. Hopp and M.L. Spearman, Characterizing the output process of a CONWIP line with deterministic processing and random outages, *Management Science* **39** (8), 975–988 (1993).
6. B. Tan, Effects of variability on the due-date performance of production lines, *Koç University Working Paper Series,* No. 97-3, Istanbul, Turkey, (1997).
7. B. Tan, Variance of the throughput of an $N$-station production line with no intermediate buffers and time dependent failures, *European Journal of Operational Research* **101** (3), 560–576 (1997).
8. G.J. Miltenburg, Variance of the number of units produced on a transfer line with buffer inventories during a period of length $T$, *Naval Research Logistics* **34**, 811–822 (1987).
9. K.B. Hendrics, The output processes of serial production lines of exponential machines with finite buffers, *Operations Research* **40** (6), 1139–1147 (1992).
10. K.B. Hendrics and J.O. McClain, The output processes of serial production lines of general machines with finite buffers, *Management Science* **39** (10), 1194–1201 (1993).
11. M. Carrascosa, Variance of the output in a deterministic two-machine line, M.S. Thesis, Massachusetts Institute of Technology, Cambridge, MA, (1995).
12. I. Duenyas and W.J. Hopp, Estimating variance of output from cyclic exponential queueing systems, *Queueing Systems* **7**, 337–354 (1990).
13. W. Grassman, Means and variances of time averages in Markovian environments, *European Journal of Operational Research* **31**, 132–139 (1987).
14. A. Høyland and M. Rausand, *System Reliability Theory: Models and Statistical Methods,* John Wiley & Sons, New York, (1994).
15. J.H. Matis, T.E. Wehrly and C.M. Metzler, On some stochastic formulations and related statistical moments of pharmokinetic models, *Journal of Pharmokinetics and Biopharmaceutics* **11** (1), 77–92 (1983).
16. J. Keilson and S.S. Rao, A process with chain dependent growth rate, *Journal of Applied Probability* **7**, 699–711 (1970).
17. W.J. Stewart, *Introduction to the Numerical Solution of Markov Chains,* Princeton University Press, New Jersey, (1994).
18. J.A. Buzacott and J.G. Shanthikumar, *Stochastic Models of Manufacturing Systems,* Prentice-Hall, Englewood Cliffs, NJ, (1993).
19. H.T. Papadopoulos, The throughput rate of multistation reliable production lines with no intermediate buffers, *Operations Research* **43** (4), 712–715 (1995).