
On the Asymptotic Optimality of the Spt Rule for the Flow Shop Average Completion Time Problem

Author(s): Cathy H. Xia, J. George Shanthikumar, Peter W. Glynn

Source: *Operations Research*, Vol. 48, No. 4 (Jul. - Aug., 2000), pp. 615-622

Published by: INFORMS

Stable URL: <http://www.jstor.org/stable/222880>

Accessed: 20/07/2010 02:09

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/action/showPublisher?publisherCode=informs>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



INFORMS is collaborating with JSTOR to digitize, preserve and extend access to *Operations Research*.

ON THE ASYMPTOTIC OPTIMALITY OF THE SPT RULE FOR THE FLOW SHOP AVERAGE COMPLETION TIME PROBLEM

CATHY H. XIA

IBM Thomas J. Watson Research Center, PO Box 704, Yorktown Heights, New York 10598
hhcx@watson.ibm.com

J. GEORGE SHANTHIKUMAR

Department of IEOR and the Walter A. Haas School of Business, University of California at Berkeley, Berkeley, California 94720
jgshant@ieor.berkeley.edu

PETER W. GLYNN

Department of EES & OR, Stanford University, Stanford, California 94305
glynn@leland.stanford.edu

(Received September 1997; revisions received October 1998, December 1998; accepted December 1998)

Consider a flow shop with M machines in series, through which a set of jobs are to be processed. All jobs have the same routing, and they have to be processed in the same order on each of the machines. The objective is to determine such an order of the jobs, often referred to as a permutation schedule, so as to minimize the total completion time of all jobs on the final machine. We show that when the processing times are statistically exchangeable across machines and independent across jobs, the Shortest Processing Time first (SPT) scheduling rule, based on the total service requirement of each job on all M machines, is asymptotically optimal as the total number of jobs goes to infinity. This extends a recent result of Kaminsky and Simchi-Levi (1996), in which a crucial assumption is that the processing times on all M machines for all jobs must be i.i.d.. Our work provides an alternative proof using martingales, which can also be carried out directly to show the asymptotic optimality of the weighted SPT rule for the Flow Shop Weighted Completion Time Problem.

1. INTRODUCTION

Flow shop scheduling, as modeled by a tandem queue, is a rich classical topic motivated mainly by practical needs in the manufacturing context. With recent developments in communication networks, the role of tandem queue scheduling becomes increasingly significant. Research has concentrated on different issues such as buffering, service effort allocation, and completion time minimization. Results such as reversibility, duality, and the bowl-shaped phenomenon are all well known in the literature, as established by Muth (1979), Pinedo (1995a), and Hillier and Boling (1979). A more detailed literature review of these issues can be found in Buzacott and Shanthikumar (1993), Pinedo (1995b), and Weber (1992).

In this paper, we restrict our attention to the issue of completion time. Consider a flow shop with M machines in series, through which a set of jobs are to be processed. All jobs have the same routing, and they have to be processed in the same sequence on each of the machines. The objective is to determine a sequence of jobs, often referred to as a permutation schedule, that minimizes the average—or equivalently, the total—of the completion times of all jobs on the final machine. It is assumed that the buffer space between machines is unlimited and that the jobs are processed individually without preemption. This problem is well known

to be NP-hard even in the two-machine case, as pointed out by Garey et al. (1976).

Research on completion time can be divided into two categories. One focuses on *makespan*, the time needed to complete processing all jobs, which takes a social point of view from the system side. The other concentrates on the *average completion time* of all jobs on the final machine, which focuses more on the individual jobs. Although the latter captures important real-life managerial scheduling concerns that are not reflected in makespan related objectives, the analysis is often very complicated and difficult. So far, most results in this area have focused on minimizing the makespan. For a complete survey of makespan related work see, for example, Pinedo (1995b).

For average completion time related scheduling, most of the previous research has studied deterministic problems using branch-and-bound strategies and is often limited to a small number of jobs and machines. Examples include Krone and Steiglitz (1974) and Van de Velde (1990), along with many others. In this paper, we look at the stochastic case instead and search for an asymptotically optimal schedule. We adopt the notion of asymptotically optimal scheduling that has been recently considered by Shanthikumar and Xu (1997), and Xia (1999), in queueing control problems. The limiting performance of tandem queues

Subject classifications: Asymptotically optimal scheduling; flow shop average completion time problem. Flow shop scheduling; average completion time, asymptotic optimality. Tandem queueing; asymptotic optimality of SPT.

Area of review: STOCHASTIC MODELS.

has been studied extensively by, for example, Glynn and Whitt (1991). In this study, we combine ideas from both and apply them to the *Flow Shop Average Completion Time Problem*. We show that as the number of jobs to be scheduled becomes larger and larger, which is often demanded in real-life applications, the Shortest total Processing Time (SPT) first rule is asymptotically optimal, provided that the processing times are statistically exchangeable across machines and independent across jobs.

It should be noted that this research was motivated by a recent result of Kaminsky and Simchi-Levi (1996b), in which it was first pointed out that the SPT rule was asymptotically optimal when the processing times are i.i.d. across machines and jobs. We provide a simplified proof using martingales and extend the result to a more general realm, where the jobs need not be identical—an assumption crucial to their paper. We show that for the SPT rule to be asymptotically optimal, the jobs can have different processing time distributions, as long as they are independent across jobs and the processing times of each job on all machines are statistically exchangeable (not necessarily independent). In addition, our argument can be easily extended to show the asymptotic optimality of SWPT—Shortest Weighted total Process Time first rule—for the *Flow Shop Weighted Completion Time Problem*. This extends another result of Kaminsky and Simchi-Levi (1996a) and simplifies their argument.

The remainder of the paper is organized as follows. Section 2 introduces the basic model. Section 3 reviews the well-known critical path approach using an activity network in determining the final completion time for each job. The total completion time for all jobs is then evaluated in §§4 and 5, where a simple lower bound and an upper bound are presented, respectively. In §6 we show, using martingales, that when scheduling is based on the total service requirement on all machines for each job, the difference between the upper bound and the lower bound converges to zero under an *exchangeability* hypothesis. The main result is then presented in §7.

2. THE MODEL

We consider the traditional Flow Shop Average Completion Time Problem, which is well known in the field of scheduling. In this problem, a flow shop consisting of M machines, each with unlimited buffer space, must sequentially process n jobs. Each machine can handle at most one job at a time, and each job can be processed on only one machine at a time without preemption. Job j , $j = 1, 2, \dots, n$ has a processing time $t_{j,m}$ on machine m , $m = 1, 2, \dots, M$, and the processing times must satisfy the following assumptions.

ASSUMPTION A1. *Processing times of different jobs are independent, i.e., $((t_{j,1}, \dots, t_{j,M}) : j \geq 1)$ is a sequence of independent random vectors.*

ASSUMPTION A2. *For each job, the processing times on the M machines are exchangeable random variables, i.e., for any given j (≥ 1), $t_{j,1}, \dots, t_{j,M}$ are exchangeable.¹*

We restrict our attention to the set of *permutation* scheduling policies, where the order in which the jobs are processed on the first machine is maintained throughout the system; that is, once the order of processing the n jobs on the first machine is scheduled, all other machines process the jobs in a *first come first serve* (FCFS) manner. It should also be noticed that permutation schedules are not necessarily optimal for the general flow shop scheduling problem where reordering is allowed at each machine. For simplicity, we consider only the permutation schedules.

The objective is to determine π —a *schedule*, or sequence of the jobs, such that $Z_M^\pi(n)$, the total completion times of all jobs on the last machine M is minimized. Let $Z_M^*(n)$ denote the optimal objective function value.

This paper provides a proof of the asymptotic optimality of the well-known rule in scheduling theory, known as the *Shortest Processing Time* (SPT) first policy. According to this policy, jobs are sequenced in increasing order of their *total service requirements* on all M machines.

Note that there may exist different SPT schedules associated with the same total processing time sequence (when there is a tie in the total processing time). For the general flow shop problem (without Assumptions A1 and A2), these SPT schedules are not necessarily optimal. See Kaminsky and Simchi-Levi (1996b) for such a counterexample. Nevertheless, given Assumptions A1, A2, and some finite moment conditions on the processing times, we show that SPT is asymptotically optimal in the sense that

$$\frac{Z_M^{\text{SPT}}(n)}{Z_M^*(n)} \rightarrow 1 \quad \text{a.s. as } n \rightarrow \infty,$$

where $Z_M^{\text{SPT}}(n)$ denotes the corresponding total completion time for processing the n jobs under the SPT policy.

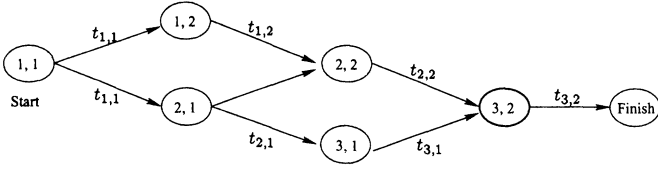
3. CRITICAL PATH

It is well known that an *Activity Network*, composed of nodes and directed arcs, is helpful in calculating completion times in a flow shop. We consider the model in which jobs are processed according to their nominal order $1, 2, \dots, n$. Let node (j, m) be associated with the activity of processing job j on machine m . Activity node (j, m) is connected with its immediate consequent activities, i.e., node $(j + 1, m)$ and node $(j, m + 1)$ with a directed arc of length $t_{j,m}$. This represents the fact that only when job j is completed on machine m , which should take time $t_{j,m}$, can the service of job $j + 1$ on machine m and the service of job j on machine $m + 1$ be initiated. Figure 1, for example, illustrates the activity network graph for processing jobs 1, 2, and 3 sequentially on machines 1 and 2.

Let $F_{j,m}$ be the completion time of job j on machine m . A basic recursion for the completion times is

$$F_{k+1,m} = \max\{F_{k,m}, F_{k+1,m-1}\} + t_{k+1,m}. \quad (1)$$

From (1), it is easy to establish the following lemma by induction.

Figure 1. Activity network graph.


LEMMA 1. *The completion time of job k on the last machine M is given by:*

$$F_{k,M} = \max_{1 \leq l_1 \leq l_2 \leq \dots \leq l_{M-1} \leq k} \left\{ \sum_{j=1}^{l_1} t_{j,1} + \sum_{j=l_1}^{l_2} t_{j,2} + \dots + \sum_{j=l_{M-1}}^k t_{j,M} \right\} \quad (2)$$

for $k = 1, \dots, n$.

Lemma 1 is a well-known result; see, for example, Conway et al. (1967), Pinedo (1995b), and Glynn and Whitt (1991). Graphically, based on the activity network, this is equivalent to saying that the completion time of the last activity must be the directed path with maximum length from the starting node to the finishing node, which is often identified as the *Critical Path*.

4. A SIMPLE LOWER BOUND FOR THE TOTAL COMPLETION TIME

Associated with the M -machine Flow Shop Average Completion Time problem introduced in §2 is the following *single machine* Average Completion Time problem: n jobs are to be processed on a single machine where job j has processing time $T_j := \sum_{m=1}^M t_{j,m}$, $j = 1, \dots, n$. The objective is to determine a sequence of the jobs so as to minimize the average completion time of all jobs. Let $\bar{Z}_1^*(n)$ be the corresponding optimal total completion time.

It turns out that the optimal total completion time of a single machine scheduling problem naturally gives a lower bound for the total completion time of the M -machine scheduling problem. This was first observed by Kaminsky and Simchi-Levi (1996b), and we state their result in the following lemma.

LEMMA 2. *Consider the M -machine Flow Shop Average Completion Time Problem and the Associated Single Machine Scheduling Problem. Suppose that the corresponding optimal total completion times are $Z_M^*(n)$ and $\bar{Z}_1^*(n)$, respectively. We then have*

$$\frac{1}{M} \bar{Z}_1^*(n) \leq Z_M^*(n). \quad (3)$$

It is well known that the *Shortest Processing Time* first rule is the optimal solution to the single machine scheduling problem; see for example Conway et al. (1967) and Pinedo

(1995b). Therefore, if we order jobs $1, 2, \dots, n$ in the increasing order of T_j s, such that $T_1 \leq T_2 \leq \dots \leq T_n$, then the optimal total completion time is given by $\bar{Z}_1^*(n) = \sum_{k=1}^n \sum_{j=1}^k T_j$.

In later sections we will see that the optimality of the SPT rule for single machine scheduling is directly related to the asymptotically optimal scheduling of M -machine scheduling. To be more specific, we will show that the lower bound given by (3) is asymptotically tight, which then yields the asymptotic optimality of the SPT rule based on the T_j s for the original problem.

5. UPPER BOUND FOR THE TOTAL COMPLETION TIME

From (2) we obtain

$$\begin{aligned} F_{k,M} &= \sum_{j=1}^k t_{j,1} \\ &\leq \max_{1 \leq l_1 \leq l_2 \leq \dots \leq l_{M-1} \leq k} \left\{ \left[\sum_{j=l_1}^{l_2} (t_{j,2} - t_{j,1}) \right. \right. \\ &\quad \left. \left. + \dots + \sum_{j=l_{M-1}}^k (t_{j,M} - t_{j,1}) \right] + \sum_{m=1}^{M-1} t_{l_m,1} \right\} \\ &\leq \sum_{m=2}^M \max_{1 \leq l_a \leq l_b \leq k} \left| \sum_{j=l_a}^{l_b} (t_{j,m} - t_{j,1}) \right| \\ &\quad + \max_{1 \leq l_1 \leq l_2 \leq \dots \leq l_{M-1} \leq k} \sum_{m=1}^{M-1} t_{l_m,1} \\ &\leq \sum_{m=2}^M 2 \max_{1 \leq l \leq k} \left| \sum_{j=1}^l (t_{j,m} - t_{j,1}) \right| + M \max_{1 \leq l \leq k} t_{l,1}. \end{aligned}$$

Given that the n jobs are processed in an order defined by a permutation σ_n , then the total completion time $Z_M^{\sigma_n}(n) = \sum_{k=1}^n F_{\sigma_n(k),M}$ will satisfy

$$\begin{aligned} &\frac{Z_M^{\sigma_n}(n) - \sum_{k=1}^n \sum_{j=1}^k t_{\sigma_n(j),1}}{n^2} \\ &\leq \frac{1}{n^2} \sum_{k=1}^n \sum_{m=2}^M 2 \max_{1 \leq l \leq k} \left| \sum_{j=1}^l (t_{\sigma_n(j),m} - t_{\sigma_n(j),1}) \right| \\ &\quad + M \frac{1}{n^2} \sum_{k=1}^n \max_{1 \leq l \leq k} t_{\sigma_n(l),1} \\ &\leq 2 \sum_{m=2}^M \frac{1}{n} \max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),m} - t_{\sigma_n(j),1}) \right| \\ &\quad + M \frac{1}{n} \max_{1 \leq l \leq n} t_{\sigma_n(l),1}. \end{aligned} \quad (4)$$

For a sequence of jobs $1, 2, \dots$, let $\mathcal{G} := \sigma(T_1, T_2, \dots)$, i.e., \mathcal{G} is the σ -algebra consisting of all the information associated with the total service requirement on all M machines

for each job. Note that \mathcal{G} includes all the information that is needed for the initial scheduling using SPT. We next show that if σ_n is \mathcal{G} -measurable, meaning we schedule the jobs only based on the T_j s, then the right-hand side of (4) converges to zero as $n \rightarrow \infty$ under some mild conditions.

6. CONVERGENCE THEOREM

LEMMA 3. Given a probability space (Ω, \mathcal{F}, P) , assume that $((t_{j,m} : 1 \leq m \leq M) : j \geq 1)$ satisfies Assumptions A1 and A2. Then

- (i) $E[t_{\sigma_n(j),m} - t_{\sigma_n(j),1} | \mathcal{G}] = 0$.
- (ii) Conditional on \mathcal{G} , $t_{\sigma_n(j),m} - t_{\sigma_n(j),1}$ ($1 \leq j \leq n$) are independent r.v.s.

PROOF. We fix $M = 2$ in this proof. The same argument works for the general case. Note that

$$\begin{aligned} P(t_{j,m} \in dx_{j,m}, 1 \leq j \leq n, m = 1, 2 | \mathcal{G}) \\ &= \prod_{j=1}^n P(t_{j,m} \in dx_{j,m}, m = 1, 2 | \mathcal{G}) \\ &= \prod_{j=1}^n P(t_{j,m} \in dx_{j,m}, m = 1, 2 | T_j). \end{aligned}$$

Because the processing times are exchangeable, i.e., $(t_j^1, t_j^2) \stackrel{d}{=} (t_j^2, t_j^1)$, we have

$$P(t_{j,1} \in dx_{j,1}, t_{j,2} \in dx_{j,2} | T_j) = P(t_{j,2} \in dx_{j,1}, t_{j,1} \in dx_{j,2} | T_j).$$

The results then follow immediately. \square

THEOREM 1. Assume the conditions of Lemma 3. If σ_n is \mathcal{G} -measurable, and

$$\sup_{j \geq 1} \text{var}(t_{j,1}) < \infty, \tag{5}$$

then, for $\delta > 0$,

$$\frac{\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),m} - t_{\sigma_n(j),1}) \right|}{n^{1/2+\delta}} \Rightarrow 0,$$

as $n \rightarrow \infty$, $\forall m = 2, \dots, M$, where \Rightarrow means convergence in distribution.

PROOF. We fix $m = 2$ in this proof. The same argument works for the general case.

Based on Lemma 3, $t_{\sigma_n(j),2} - t_{\sigma_n(j),1}$ ($1 \leq j \leq n$) are conditionally independent random variables with zero mean. Hence, by Kolmogorov's inequality, (refer to, e.g., Williams, 1991) we have:

$$\begin{aligned} P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right| > x | \mathcal{G} \right) \\ \leq \frac{\sum_{j=1}^n E[(t_{\sigma_n(j),2} - t_{\sigma_n(j),1})^2 | \mathcal{G}]}{x^2} \end{aligned}$$

$$\begin{aligned} &= \frac{E[\sum_{j=1}^n (t_{\sigma_n(j),2} - t_{\sigma_n(j),1})^2 | \mathcal{G}]}{x^2} \\ &= \frac{E[\sum_{j=1}^n (t_{j,2} - t_{j,1})^2 | \mathcal{G}]}{x^2} \\ &\leq \frac{4 \sum_{j=1}^n \text{var}(t_{j,1} | \mathcal{G})}{x^2}. \end{aligned}$$

Hence it follows that

$$\begin{aligned} P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right| > x \right) \\ \leq \frac{2}{x^2} \sum_{j=1}^n E[\text{var}(t_{j,1} | \mathcal{G})] \leq \frac{2}{x^2} \sum_{j=1}^n \text{var}(t_{j,1}). \end{aligned}$$

The result follows immediately from this inequality. \square

REMARK. In general, condition (5) is valid when the processing times have uniformly finite expected residual lives, i.e.,

$$R := \sup_j \sup_{t \geq 0} E[t_{j,m} - t | t_{j,m} > t] < \infty. \tag{6}$$

To see this, we note that

$$\begin{aligned} E[t_{j,m}^2] &= 2E \int_0^{t_{j,m}} (t_{j,m} - u) du \\ &= 2E \int_0^\infty (t_{j,m} - u) I(t_{j,m} > u) du \\ &= 2 \int_0^\infty E[t_{j,m} - u | t_{j,m} > u] P(t_{j,m} > u) du \\ &\leq 2R \int_0^\infty P(t_{j,m} > u) du \\ &= 2RE[t_{j,m}] \leq 2R^2, \end{aligned}$$

where the exchange of the expectation and the integral in third equality is justified by Fubini's theorem.

Examples satisfying condition (6) include cases where:

1. The $t_{j,m}$ s are uniformly bounded.
2. $t_{j,m} \sim \exp(\mu_j)$ with $\inf_j \mu_j > 0$.
3. The $t_{j,m}$ s are new better than used in expectation (NBUE) with $\sup_j E[t_{j,m}] < \infty$.

Our next theorem deals with the almost sure convergence to zero of the right-hand side of (4) as n goes to infinity, which requires the following stronger condition.

THEOREM 2. Assume the condition of Lemma 3. If

$$\sup_{j \geq 1} E[|t_{j,1}|^p] < \infty, \tag{7}$$

for some $p > 2$, and $(\sigma_n, n \geq 1)$ is a sequence of \mathcal{G} -measurable permutations, then

$$\frac{\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),m} - t_{\sigma_n(j),1}) \right|}{n} \rightarrow 0 \quad \text{a.s. as } n \rightarrow \infty.$$

PROOF. Again, we look only at $m = 2$ and extend the same argument to general m .

For $1 \leq l \leq n$, let $\mathcal{H}_l := \mathcal{G} \vee \sigma(t_{\sigma_n(j),2} - t_{\sigma_n(j),1} : 1 \leq j \leq l)$. Then $(\sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) : 0 \leq l \leq n)$ is a martingale with respect to $(\mathcal{H}_l : 1 \leq l \leq n)$. This is obvious because, given \mathcal{G} , the increments are conditionally independent with zero conditional expectation, as established in Lemma 3. Hence for $p \geq 1$,

$$\left(\left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|^p : 0 \leq l \leq n \right)$$

is a submartingale, provided that

$$E \left(\left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|^p \right) < \infty, \tag{8}$$

for $1 \leq j \leq n$.

In order to verify (8), note that under condition (7) we have

$$\begin{aligned} \infty > E \left(\sum_{j=1}^n |t_{j,2} - t_{j,1}|^p \right) &= E \left(\sum_{j=1}^n |t_{\sigma_n(j),2} - t_{\sigma_n(j),1}|^p \right) \\ &= \sum_{j=1}^n E(|t_{\sigma_n(j),2} - t_{\sigma_n(j),1}|^p). \end{aligned}$$

So, it follows that $E(|t_{\sigma_n(j),2} - t_{\sigma_n(j),1}|^p) < \infty$ for $1 \leq j \leq n$. Note that $|x|^p$ is a convex function in x ; therefore, by Jensen's inequality,

$$\begin{aligned} E \left(\left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|^p \right) \\ \leq \sum_{j=1}^l E(|t_{\sigma_n(j),2} - t_{\sigma_n(j),1}|^p) < \infty. \end{aligned}$$

We now apply the *maximal inequality for submartingales* (see Chung 1974, p. 303):

$$\begin{aligned} P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right| > x \right) \\ = P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|^p > x^p \right) \\ \leq x^{-p} E \left(\left| \sum_{j=1}^n (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|^p \right) \\ = x^{-p} E \left(\left| \sum_{j=1}^n (t_{j,2} - t_{j,1}) \right|^p \right). \end{aligned}$$

Hence,

$$\begin{aligned} P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right| > x \right) \\ \leq x^{-p} E \left(\left| \sum_{j=1}^n (t_{j,2} - t_{j,1}) \right|^p \right). \end{aligned}$$

Under condition (7), we know that for $p > 2$,

$$\sup_{j \geq 1} E[|t_{j,2} - t_{j,1}|^p] < \infty.$$

This implies, using the *Burkholder-Gundy square inequality* for martingales (see Burkholder 1972), that

$$E \left[\left| \sum_{j=1}^n (t_{j,2} - t_{j,1}) \right|^p \right] = O(n^{p/2}). \tag{9}$$

Therefore, $\forall \varepsilon > 0$,

$$\begin{aligned} \sum_{n=1}^{\infty} P \left(\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right| > \varepsilon n \right) \\ \leq \varepsilon^{-p} \sum_{n=1}^{\infty} n^{-p} E \left(\left| \sum_{j=1}^n (t_{j,2} - t_{j,1}) \right|^p \right) \\ = \varepsilon^{-p} \sum_{n=1}^{\infty} n^{-p} O(n^{p/2}) < \infty. \end{aligned}$$

Consequently, the *Borel-Cantelli lemma* establishes that if σ_n is \mathcal{G} -measurable, then

$$\frac{\max_{1 \leq l \leq n} \left| \sum_{j=1}^l (t_{\sigma_n(j),2} - t_{\sigma_n(j),1}) \right|}{n} \rightarrow 0 \quad \text{a.s. as } n \rightarrow \infty.$$

□

REMARK. In general, condition (7) is valid given that the processing times have uniformly finite $(p - 1)$ th moment for the residual lives, i.e.,

$$\sup_j \sup_{t \geq 0} E[(t_{j,m} - t)^{p-1} | t_{j,m} > t] < \infty.$$

Examples include:

1. The $t_{j,m}$ s are uniformly bounded.
2. $t_{j,m} \sim \exp(\mu_j)$ with $\inf_j \mu_j > 0$.

7. ASYMPTOTIC OPTIMALITY

Based on (4), Theorems 1 and 2 then establish, respectively, the weak convergence and almost sure convergence of

$$\frac{Z_M^{\sigma_n}(n) - \sum_{k=1}^n \sum_{j=1}^k t_{\sigma_n(j),1}}{n^2} \rightarrow \text{zero as } n \rightarrow \infty.$$

Returning to §5, if instead we reformulate (2) as

$$F_{k,M} \leq \sum_{j=1}^k t_{j,m} + \max_{1 \leq l_1 \leq l_2 \leq \dots \leq l_{M-1} \leq k} \left\{ \left[\sum_{s=1, \dots, M; s \neq m} \sum_{j=l_{s-1}}^{l_s} (t_{j,s} - t_{j,m}) \right] + \sum_{s=1}^{M-1} t_{l_s, m} \right\},$$

where $l_0 = 1$ and $l_M = k$, similar arguments can then be carried out to obtain the convergence of

$$\frac{Z_M^{\sigma_n}(n) - \sum_{k=1}^n \sum_{j=1}^k t_{\sigma_n(j), m}}{n^2} \rightarrow 0 \text{ as } n \rightarrow \infty,$$

for $m = 1, 2, \dots, M$.

The next result then follows immediately.

THEOREM 3. *In an M -machine flow shop, assume the conditions of Lemma 3. If jobs $1, 2, \dots, n$ are processed according to a \mathcal{G} -measurable permutation σ_n , and $Z_M^{\sigma_n}(n)$ is the corresponding total completion time, then*

(i) *under condition (5),*

$$\frac{Z_M^{\sigma_n}(n) - \frac{1}{M} \sum_{k=1}^n \sum_{j=1}^k T_{\sigma_n(j)}}{n^2} \Rightarrow 0, \tag{10}$$

as $n \rightarrow \infty$,

(ii) *under condition (7),*

$$\frac{Z_M^{\sigma_n}(n) - \frac{1}{M} \sum_{k=1}^n \sum_{j=1}^k T_{\sigma_n(j)}}{n^2} \rightarrow 0 \text{ a.s. as } n \rightarrow \infty. \tag{11}$$

If we process the n jobs based on the SPT rule, i.e., the permutation σ_n is defined such that

$$T_{\sigma_n(1)} \leq T_{\sigma_n(2)} \leq \dots \leq T_{\sigma_n(n)}, \tag{12}$$

and call the corresponding performance $Z_M^{\text{SPT}}(n)$, then

$$\frac{1}{M} \bar{Z}_1^*(n) \leq Z_M^*(n) \leq Z_M^{\text{SPT}}(n),$$

where $\bar{Z}_1^*(n) = \bar{Z}_1^{\text{SPT}}(n) = \sum_{k=1}^n \sum_{j=1}^k T_{\sigma_n(j)}$. Therefore,

$$\begin{aligned} \left| \frac{Z_M^{\text{SPT}}(n) - Z_M^*(n)}{n^2} \right| &\leq \left| \frac{Z_M^{\text{SPT}}(n) - \frac{1}{M} \bar{Z}_1^*(n)}{n^2} \right| \\ &= \left| \frac{Z_M^{\text{SPT}}(n) - \frac{1}{M} \sum_{k=1}^n \sum_{j=1}^k T_{\sigma_n(j)}}{n^2} \right|. \end{aligned}$$

From Theorem 3 it then follows immediately that, as $n \rightarrow \infty$,

$$\left| \frac{Z_M^{\text{SPT}}(n) - Z_M^*(n)}{n^2} \right| \Rightarrow 0 \text{ under condition (5),}$$

and

$$\left| \frac{Z_M^{\text{SPT}}(n) - Z_M^*(n)}{n^2} \right| \rightarrow 0 \text{ a.s. under condition (7).}$$

LEMMA 4. *If there exists a $\delta > 0$ such that*

$$\inf_{j \geq 1} P\{t_{j,m} \geq \delta\} > 0, \tag{13}$$

then

$$\liminf_{n \rightarrow \infty} \frac{\bar{Z}_1^*(n)}{n^2} > 0 \text{ a.s.}$$

PROOF. For notation simplicity, let $\alpha_j := P\{T_j \geq \delta\}$ and $\gamma := \inf_j \alpha_j$. Clearly (13) implies $0 < \gamma \leq 1$.

For a given increasing sequence $a_1 \leq a_2 \leq \dots \leq a_n$, suppose a_l is the first a_i larger than δ , then $n + 1 - l$ is the total number of a_i s larger than δ , and we have

$$\begin{aligned} \sum_{k=1}^n \sum_{j=1}^k a_j &= \sum_{j=1}^n (n + 1 - j) a_j \geq \delta \sum_{j=l}^n (n + 1 - j) \\ &= \frac{\delta}{2} (n + 1 - l)(n + 2 - l). \end{aligned}$$

Applying this to the increasing sequence in (12) then gives

$$\bar{Z}_1^*(n) = \sum_{k=1}^n \sum_{j=1}^k T_{\sigma_n(j)} \geq \frac{\delta}{2} N_n(N_n + 1),$$

where $N_n := \sum_{j=1}^n I_{\{T_j \geq \delta\}}$.

Let $Y_j = I_{\{T_j \geq \delta\}} I_{\{U_j \leq \gamma/\alpha_j\}}$, where $\{U_j : j \geq 1\}$ is a sequence of independent uniform random variables, independent of the processing time processes. Note that $P(Y_j = 1) = \gamma$ for all j s, therefore the Y_j s are i.i.d. r.v.s. Hence, based on the SLLN,

$$\frac{N_n}{n} \geq \frac{\sum_{j=1}^n Y_j}{n} \rightarrow E[Y_j] = \gamma > 0.$$

It then follows that

$$\frac{\bar{Z}_1^*(n)}{n^2} \geq \frac{N_n(N_n + 1)}{2n^2} \delta \rightarrow \frac{1}{2} \gamma^2 \delta > 0,$$

as $n \rightarrow \infty$. Thus the claim in Lemma 4 follows. \square

REMARK. If $\inf_{j \geq 1} E[t_{j,m}] > 0$ and condition (5) is assumed, condition (13) holds.

To see this, let c be the upper bound on the variances implied by (5), and let b be the assumed lower bound on the means of the $t_{j,m}$ s. Then, for $0 < \varepsilon < 1$, Chebyshev's inequality implies that

$$\begin{aligned} b \leq E[t_{j,m}] &= \int_0^\infty P(t_{j,m} > u) du \\ &= \int_0^\varepsilon P(t_{j,m} > u) du + \int_\varepsilon^{1/\varepsilon} P(t_{j,m} > u) du \\ &\quad + \int_{1/\varepsilon}^\infty P(t_{j,m} > u) du \end{aligned}$$

$$\begin{aligned} &\leq \varepsilon + \left(\frac{1}{\varepsilon} - \varepsilon\right) P(t_{j,m} > \varepsilon) + \int_{1/\varepsilon}^{\infty} \frac{\text{var}(t_{j,m})}{u^2} du \\ &\leq \left(\frac{1}{\varepsilon} - \varepsilon\right) P(t_{j,m} > \varepsilon) + (1 + c)\varepsilon. \end{aligned}$$

Hence,

$$P(t_{j,m} > \varepsilon) \geq \frac{b - (1 + c)\varepsilon}{1/\varepsilon - \varepsilon}.$$

By choosing ε sufficiently small, we find the required positive lower bound on $P(t_{j,m} > \varepsilon)$. Examples of distributions satisfying the lower bound property are:

1. the $t_{j,m}$ s are uniformly bounded away from zero, and
2. $t_{j,m} \sim \exp(\mu_j)$ with $\sup_j \mu_j < \infty$.

We then conclude with the following main result.

THEOREM 4. *For the M -machine flow shop total completion time scheduling problem, the SPT rule is asymptotically optimal in the sense that*

- (i) *under the conditions of Lemma 3 and (5),*

$$\frac{Z_M^{\text{SPT}}(n) - Z_M^*(n)}{n^2} \Rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

- (ii) *under the conditions of Lemma 3 and (7),*

$$\frac{Z_M^{\text{SPT}}(n) - Z_M^*(n)}{n^2} \rightarrow 0 \quad \text{a.s. as } n \rightarrow \infty.$$

If, in addition, all the processing times $t_{j,m}$ s are uniformly bounded away from zero, then

$$\frac{Z_M^{\text{SPT}}(n)}{Z_M^*(n)} \rightarrow 1 \quad \text{a.s. as } n \rightarrow \infty.$$

REMARK. The same argument can be applied to the *Flow Shop Weighted Completion Time Problem*. In this problem, the processing times of job j , $j = 1, \dots, n$, are given an associated weight ω_j , with $0 < \omega_j < 1$, and the objective is to determine a permutation schedule that minimizes the total weighted completion time of all jobs on the final machine. For example, if the n jobs are processed in the nominal order $1, \dots, n$, then the total weighted completion time is $Z = \sum_{j=1}^n \omega_j F_{j,M}$. Under the assumption that the processing times are i.i.d., Kaminsky and Simchi-Levi (1996a) have established the asymptotic optimality of the *Shortest Weighted Completion Time (SWPT)* first rule, which sequences the jobs in the increasing order of the weighted total processing times $\omega_j T_j$, $j = 1, \dots, n$.

By applying the convergence theorems in §6 to the weighted sum and with Assumptions A1 and A2, we may prove theorems analogous to Theorems 3 and 4, showing that the SWPT rule is asymptotically optimal. This extends Kaminsky and Simchi-Levi (1996a) to a more general setting than i.i.d. processing times.

8. CONCLUSION

In this paper we have provided a martingale approach that proves the asymptotic optimality of the SPT rule for the *Flow Shop Average Completion Time Problem* when the processing times on the machines are i.i.d., or more generally, statistically exchangeable across machines and independent across jobs. The same argument can also be carried out to show the asymptotic optimality of the SWPT rule for the *Flow Shop Weighted Completion Time Problem*. This extends the recent results of Kaminsky and Simchi-Levi (1996a, 1996b).

END NOTES

Random variables X_1, \dots, X_k are said to be *exchangeable* if $(X_1, \dots, X_k) \stackrel{d}{=} (X_{\sigma_k(1)}, \dots, X_{\sigma_k(k)})$ for any deterministic permutation σ_k of $\{1, \dots, k\}$.

ACKNOWLEDGMENTS

This work was supported in part by the U.S. Army Research Office under contract/grant DAAG 55-97-1-0377.

REFERENCES

- Burkholder, D., D. Davis, R. Gundy. 1972. Integral inequalities for convex functions of operators on martingales. *Proc. Sixth Berkeley Sympos. Math. Statist. Probab., Vol. 2*. Univ. of California Press, Berkeley, CA, 223–240.
- Buzacott, J., J.G. Shanthikumar. 1993. *Stochastic Models of Manufacturing Systems*. Prentice Hall, Englewood Cliffs, NJ.
- Chung, K. L. 1974. *A Course in Probability Theory*. Academic Press, New York.
- Conway, R. W., W. L. Maxwell, L. W. Miller. 1967. *Theory of Scheduling*. Addison-Wesley, Reading, MA.
- Garey, M. R., D. S. Johnson, R. Sethi. 1976. The complexity of flowshop and jobshop scheduling. *Math. Oper. Res.* **1** 117–129.
- Glynn, P. W., W. Whitt. 1991. Departures from many queues in series. *Ann. Appl. Probab.* **4** 546–572.
- Hillier, F. S., R. M. Boling. 1979. On the optimal allocation of work in symmetric balanced production line systems with variable operation times. *Management Sci.* **25** 721–728.
- Kaminsky, P., D. Simchi-Levi. 1996a. Probabilistic analysis and practical algorithms for the flow shop weighted completion time problem. *Oper. Res.* **46** 872–882.
- , ———. 1996b. The asymptotic optimality of the SPT rule for the flow shop mean completion time problem. To appear in *Oper. Res.*
- Krone, M. J., K. Steiglitz. 1974. Heuristic programming solutions of a flowshop scheduling problem. *Oper. Res.* **22** 629–638.
- Muth, E. J. 1979. The reversibility property of production lines. *Management Sci.* **25** 152–158.
- Pinedo, M. 1995a. Minimizing the expected makespan in stochastic flow shops. *Oper. Res.* **30** 148–162.
- . 1995b. *Scheduling: Theory, Algorithms and Systems*. Prentice Hall, Englewood Cliffs, NJ.

- Shanthikumar, J. G., S. Xu. 1997. Asymptotically optimal routing and service rate allocation in a multi-server queueing system. *Oper. Res.* **45** 464–469.
- Van de Velde, S. L. 1990. Minimizing the sum of the job completion times in the two-machine flow shop by Lagrangean relaxation. *Ann. Oper. Res.* **26** 257–268.
- Weber, R. R. 1992. The interchangeability of tandem queues with heterogeneous customers and dependent service times. *Adv. Appl. Prob.* **24** 727–737.
- Williams, D. 1991. *Probability with Martingales*. University Press, Cambridge, UK.
- Xia, C. H. 1999. Dynamic scheduling of queueing systems with applications to computer networks and flexible manufacturing. Ph.D. Thesis. Stanford University.
- Xia, C. H., J. G. Shanthikumar, 1998. Optimal and asymptotically optimal scheduling—an extension of the $c\mu$ rule, Working Paper, Stanford University.