

Dynamically scheduling and maintaining a flexible server

Jefferson Huang^{*1}, Douglas G. Down², Mark E. Lewis³, and Cheng-Hung Wu⁴

¹Operations Research Department, Naval Postgraduate School, Monterey, CA 93943-5098, USA

²Department of Computing and Software, McMaster University, Hamilton, Ontario L8S 4L7,
Canada

³School of Operations Research and Information Engineering, Cornell University, Ithaca, NY
14853-3801, USA

⁴Institute of Industrial Engineering, National Taiwan University, 1, Sec. 4, Roosevelt Rd., Taipei
106, Taiwan

February 26, 2021

*Corresponding Author: jefferson.huang@nps.edu

Abstract

Deciding how to jointly schedule jobs and perform preventive maintenance is a fundamental problem in flexible manufacturing systems, particularly those arising in semiconductor manufacturing. At the same time, past work in this area shows that, even when there is only one station and one type of job, identifying policies that minimize the amount of work-in-process (WIP) is a difficult problem. In this paper, we study a single-station version of this problem with an arbitrary number of job classes, with the objective of minimizing average maintenance costs plus the weighted average amount of WIP. We identify conditions under which it suffices to schedule jobs according to both a server-state-dependent version of the $c\mu$ -rule, and a static $c\mu$ -rule where the average service rates are used. One of these conditions states that the ratio between the service rates should remain constant as the server deteriorates. When this assumption does not hold, scheduling with the $c\mu$ -rule can in fact lead to an unstable system; we illustrate this using a simple example. On the other hand, we also present numerical evidence that $c\mu$ -based scheduling performs well compared to other scheduling rules, and relative to a policy based on solving a Markov decision process.

Keywords: scheduling, priority, maintenance, Markov decision process, queue

1 Introduction

The yield of a manufacturing process, defined as the fraction of output that is of sufficient quality, is a key economic performance indicator; see e.g., [3]. In semiconductor manufacturing, yield improvement has been recognized as an effective means of managing costs and sustaining profitability [2]. In particular, yield increases on the order of even 1-2% can lead to significant savings in wafer manufacturing costs [9].

One of the key determinants of yield is the health of the machines processing the jobs that eventually become finished products. As the underlying condition of a given machine deteriorates, the increased frequency of significant process deviations (as identified by, e.g., statistical process control procedures [23]) leads to more re-work and tuning, which in turn reduces the rate at which good products are produced (i.e., the “service rate” of the machine). Eventually, it can become worthwhile to take the deteriorated machine offline for maintenance, after which the service rate is improved.

The need for good maintenance policies, and the increasing prevalence of sensorized equipment in the semiconductor and other advanced manufacturing industries, has led to the emergence of condition-based maintenance (CBM) as a potentially cost-effective alternative to more commonly used age or job-based maintenance rules [8, 13]. At the same time, almost every machine used to process jobs in the semiconductor manufacturing setting is flexible, in the sense that it can be used to process more than one type of job. Hence a fundamental problem in semiconductor manufacturing, and more generally in flexible manufacturing systems with deteriorating equipment, is how to simultaneously (1) allocate jobs to flexible machines that deteriorate over time, and (2) perform preventive maintenance on those machines.

When there is a single operation to be performed, Kaufman and Lewis [16] show the difficulty in developing control policies that use as inputs both the current workload in the system and the condition of the machine. This leaves a broad class of manufacturing system configurations, which includes configurations arising in the semiconductor industry, without guidance on how to consider the trade-offs between resource allocation and resource maintenance. In this paper, we

consider the question of joint maintenance and scheduling in a parallel queueing setting. When there are a number of job classes, the manager can decide to assign a single resource (henceforth called a *server*) to any of the classes, or to begin preventive maintenance. The goal is to provide adequate service (in the form of minimizing weighted queue-lengths) to each class, while noting that a deteriorated server works slower. Since [16] shows in the single queue setting that the usual monotonicity properties of an optimal control do not hold, there is little hope of finding simple solutions to the scheduling/maintenance pairing. Instead, we seek insights into the following questions:

1. Given a choice between prioritizing scheduling or maintenance, where should a decision-maker focus his/her efforts?
2. Given the complexity of optimal policies in general, are there easier-to-implement heuristics that perform comparably?

We address both questions by presenting conditions under which scheduling with a natural extension of the classic $c\mu$ -rule is without loss of optimality (Theorem 5), and numerical results indicating that this heuristic performs well more generally (Section 5).

1.1 Related Literature

The two types of decisions described above, namely maintaining a deteriorating machine and scheduling jobs in a queueing system, have typically been considered separately in the literature. In particular, the majority of papers in the maintenance literature do not consider the effect that the amount of work in the system (i.e., the queue length) may have on optimal maintenance decisions; see e.g., the surveys [17, 19, 21, 24]. Two papers where such effects are accounted for, via models that are very closely related to the one presented in Section 2, are Kaufman and Lewis [16] and Cai et al. [5], which we describe in more detail below. Moreover, while previous work such as that of Andradóttir et al. [1] and Wu et al. [25, 26] has accounted for server failures in the context of queueing models of flexible manufacturing systems, with the exception of Cai et al. [5] we are

not aware of any other work in this area that combines scheduling with maintenance decisions. We note that there has been recent work on joint scheduling and maintenance in the contexts of deterministic scheduling [14], developing metaheuristics [7], and mixed-integer programming that incorporates constraints specific to semiconductor manufacturing [6, 27].

Kaufman and Lewis [16] analyze the structure of optimal maintenance policies for the server of an $M/M/1$ queue with only one type of job, where the service rate deteriorates according to a pure-death process. In particular, [16, Example 3.6] shows that the optimal policy under the average cost criterion may not be monotone in the queue lengths. For certain deterioration levels it may be optimal to perform maintenance when there are no queued jobs, not perform maintenance when there are few queued jobs, and to perform maintenance for all sufficiently large queue lengths. On the other hand, [16, Theorems 3.2, 4.2; Proposition 4.10] provide conditions under which there is an optimal policy that is monotone in the server's health. This means that there is an optimal policy with the following structure: For each fixed number of queued jobs i , there is a threshold s_i^* such that maintenance is performed if and only if the deterioration level is worse than s_i^* . Finally, numerical experiments are presented [16, Section 5] that illustrate some pitfalls associated with using some simple and natural heuristics, underscoring the difficulty of the problem.

Cai et al. [5] consider an $M/G/1$ queueing model with at most two types of jobs, which is motivated by potential semiconductor manufacturing applications. In this model, the service and deterioration dynamics differ from those in [16]. In particular, jobs cannot be preempted while they are being served, and deterioration events can only occur when a service completion occurs. On the other hand, while in [16] it is assumed that at each deterioration event the server moves to the "next-worse" state with probability 1, the model in Cai et al. [5] allows the server to move, as a result of a single deterioration event, to any state that is worse. For this model, analogous results to the ones in [16] hold. Namely, the optimal policy may not be monotone in the number of jobs [5, Section 5], but under certain conditions there exists an optimal policy that is monotone in the server's health [5, Theorems 3.3, 4.3]. In addition, for the case of two types of jobs, it is shown

that it may be suboptimal to always prioritize a job type that seems, from a cost and deterioration perspective, to be superior to the other [5, Section 4.2]. A monotonicity result [5, Theorem 4.4] on the value function when one job type is superior to the other in the aforementioned sense is also provided. Finally, numerical results [5, Sections 5,6] are presented to illustrate the savings that an optimal joint scheduling and maintenance policy can provide, relative to first-in first-out scheduling and maintenance after a fixed number of jobs have been completed.

In this paper, we consider joint scheduling and maintenance in the context of a G/M/1 queue with an arbitrary number of job classes. As in Kaufman and Lewis [16], and in contrast to Cai et al. [5], we assume that jobs can be preempted by the decision-maker, or interrupted by a failure event. We also consider deterioration dynamics that are more general than those in [16], and which differ from those in [5]. In particular, in [5] it is assumed that deterioration events must coincide with service completions, but that deterioration rates can depend on which type of job is worked on. In contrast, we assume that deterioration events can happen at any time, but that the deterioration rates are the same for both job types. While Cai et al. [5] argue that work-dependent deterioration is important for certain semiconductor manufacturing applications, Sloan and Shankthikumar [22] note that for some wafer fabrication processes, such as in etch operations, it is reasonable to assume that deterioration does not depend on the type of job.

1.2 Contributions and Outline

The main contributions of the paper are as follows. After presenting the scheduling and maintenance model in Section 2, we consider the problem of scheduling in the presence of a deteriorating server in Section 3, without preventive maintenance. We provide a condition (Assumption CR) under which it is optimal to schedule the jobs according to a static priority rule, when the service rates are modulated according to a (possibly non-Markovian) point process; see Theorem 2 and Remark 3. In addition, Example 1 shows that, when the conditions of Theorem 2 do not hold, scheduling according to the aforementioned priority rule can in fact lead to an unstable system, in the sense that the average number of queued jobs grows to infinity. From the perspective of

system design, this provides a strong incentive to invest in ensuring that Assumption **CR** below holds. Next, in Section 4 we return to the joint scheduling and maintenance problem. We use the results in Section 3 to provide conditions under which it suffices to search for an optimal policy among those that schedule according to a priority rule and that are monotone in the server’s health (Theorem 7). In Section 5, we provide numerical results indicating that the priority-rule based scheduling policies considered in Section 4 can perform well across a range of system parameters. The numerical results also underscore the value of good maintenance policies, and of incorporating service rate information in scheduling and maintenance policies. This latter point was also observed by Iravani and Duenyas [15], in the context of a single job type. Finally, conclusions and future research directions are presented in Section 6. Unless otherwise indicated, proofs of stated results are provided in Appendix A.

2 Joint Scheduling and Maintenance Model

Jobs from K classes arrive randomly over time. Each arriving job requires a random amount of work, and all incoming work is processed by a single server. The arrival times of the jobs are modeled by independent point processes on $\mathbb{R}_+ := [0, \infty)$, while the amount of work required by each arriving job is assumed to be exponentially distributed with unit mean, independently of the other jobs. It is assumed that the arrival process is *regular*, in the sense that with probability 1, there can be at most a finite number of arrivals during any finite interval of time.

Jobs of the same class are homogeneous. However, both the cost incurred by a waiting job and the time required to complete a job depend on the job’s class. For $k = 1, \dots, K$, the cost incurred by a waiting job of class k is assumed to accumulate continuously at a constant rate c_k .

As time progresses, the health of the server deteriorates. This is modeled by assuming that when the server is able to perform work, it spends a random amount of time in its current state $s \in \{0, 1, \dots, S\}$ before deteriorating to a state that is at least as bad. In particular, lower numbered states indicate worse health. In addition, given that a deterioration event has occurred, the probability that the server then transitions to state ℓ from its current state s is $q(\ell|s)$, where $q(\ell|s) = 0$ if

$\ell > s$. Once the server's state reaches 0, it undergoes maintenance for a random amount of time, after which it returns to state S . We assume that the times at which the server changes state are independent of both the amount and the nature of the work that the server has completed. In particular, the random times at which the server changes state are modeled by a point process on \mathbb{R}_+ that is independent of the arrival processes and work requirements. Like the arrival processes, this process of deterioration times is also assumed to be regular.

The rate at which the server can complete work of a given class depends on the server's health. If the state of the server is s , then the rate at which it can complete work of class k is μ_k^s . We assume that $\mu_k^0 = 0 < \mu_k^1 \leq \dots \leq \mu_k^S < \infty$ for each class $k = 1, \dots, K$. In other words, the server cannot complete any work while it is undergoing maintenance (i.e., is in state 0), and higher-numbered states indicate less deterioration.

For ease of exposition, the server is referred to as being *online* if its state is not 0, and *offline* if its state is 0. In addition, the server is said to *deteriorate* if it transitions from state $s \geq 2$ to state $\ell \geq 1$, and *fails* if it transitions to state 0 without the influence of the decision-maker. Whenever a failure occurs, *corrective maintenance* (CM) is initiated at cost $K_c \geq 0$. When the server is online, the decision-maker can initiate *preventive maintenance* (PM), which is modeled as an instantaneous transition of the server state to 0 at cost $K_p \geq 0$.

In addition, when the server is online and there is a job in the system, the server may be assigned to work on that job. Since jobs of the same class are assumed to be homogeneous, we equate selecting a job to work on with selecting which class to assign the server to.

The decision-maker is only able to exert control over the server at decision epochs, which occur whenever one of the following events occurs:

- A job arrives and the server is online.
- A job is completed and the server is online.
- The server deteriorates or fails.
- The server comes back online (from state 0).

Accordingly, jobs that are currently in service may be preempted.

At each decision epoch, the decision-maker knows the current state of the system, i.e., the number i_k of jobs of each class $k \in \{1, \dots, K\}$ present and the state $s \in \{0, 1, \dots, S\}$ of the server. Let $\mathbb{X} := \{0, 1, \dots\}^K \times \{0, 1, \dots, S\}$ denote the set of all possible system states. We also denote by \mathbf{i} the vector whose k th entry is i_k . In addition to the current state, the decision-maker also knows the history of the system (i.e., the past queue lengths, server states, and event times) up to the current decision epoch. In deciding whether to serve one of the classes or initiate PM, the decision-maker follows a policy π that prescribes (possibly in a randomized way) the action to take at each decision epoch, given the current state and history of the system. We restrict attention to policies that are *non-idling* (i.e., never call for an online server to idle when there is work to do) and *non-anticipative* (i.e., do not depend on future information). Let Π denote the set of all such policies. Of particular interest are the *deterministic stationary* policies; under such a policy π , the action $\pi(x)$ is taken whenever the system is in state $x \in \mathbb{X}$.

We compare policies on the basis of the long-run average cost per unit time incurred from a given initial state. To define this optimality criterion, fix any $\pi \in \Pi$. Let $Q_k^\pi(t)$ denote the number of jobs in class k , including those in service, at time t and $\mathbf{Q}^\pi(t)$ be the vector whose k^{th} entry is $Q_k^\pi(t)$. Also, let $S^\pi(t)$ denote the state of the server at time t under π , let $M_c^\pi(t)$ (resp. $M_p^\pi(t)$) equal 1 if CM (resp. PM) is initiated at time t under π , and let $M_c^\pi(t)$ (resp. $M_p^\pi(t)$) equal 0 otherwise. Finally, for $n = 1, 2, \dots$ let t_n^π denote the n^{th} decision epoch under π .

If the system is in state $(\mathbf{i}, s) \in \mathbb{X}$ at time 0, then the long-run expected *average cost* per unit time that is incurred by following the policy π is

$$w^\pi(\mathbf{i}, s) := \limsup_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{n: t_n^\pi \leq T} [K_c M_c^\pi(t_n^\pi) + K_p M_p^\pi(t_n^\pi)] + \int_0^T \sum_{k=1}^K c_k Q_k^\pi(t) dt \mid \mathbf{Q}^\pi(0) = \mathbf{i}, S^\pi(0) = s \right].$$

A policy $\pi_* \in \Pi$ is *optimal* if $w^{\pi_*}(\mathbf{i}, s) = \min_{\pi \in \Pi} w^\pi(\mathbf{i}, s)$ for every initial state $(\mathbf{i}, s) \in \mathbb{X}$.

3 Scheduling Without Preventive Maintenance

We first consider the case where the decision-maker cannot perform PM, and can only schedule jobs in the presence of a deteriorating server. In this setting, the server can only go offline (i.e., enter state 0) via a failure. Note that, since preventive maintenance is not permitted and the server state evolves independently of the scheduling decisions, the maintenance costs are independent of the policy used. The decision-maker's objective is therefore to find a scheduling policy π that minimizes the weighted long-run average expected number of jobs in the system

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\int_0^T \sum_{k=1}^K c_k Q_k^\pi(t) dt \right]. \quad (1)$$

in the presence of uncontrollable server deterioration.

3.1 $c\mu$ -Rules

Without server deterioration, it follows from [18, Theorem 2.1] that it is optimal to schedule according to the $c\mu$ -rule. According to this rule, if the service rate for class k jobs is μ_k , then priority is given to any class k^* where $c_{k^*} \mu_{k^*} \geq c_k \mu_k$ for every class k .

When the service rate depends on the state of the server, it is natural to consider prioritizing the jobs according to a *state-dependent* $c\mu$ -rule. Namely, if the state of the server is s , prioritize any class k^* where $c_{k^*} \mu_{k^*}^s \geq c_k \mu_k^s$ for every class k . Alternatively, letting $\nu(s)$ denote the long-run expected fraction of time that the server spends in state s , one could employ the following *average* $c\mu$ -rule: Assign priority to any class k^* for which $c_{k^*} \bar{\mu}_{k^*} \geq c_k \bar{\mu}_k$ for every class k , where $\bar{\mu}_k := \sum_{s=0}^S \nu(s) \mu_k^s$ is the *average service rate* for class k jobs.

3.2 Instability of $c\mu$ -Rules

Observe that if $c_{k^*} \mu_{k^*}^s \geq c_k \mu_k^s$ for every server state s , and every class k , then both the state-dependent and average $c\mu$ -rules described above would prioritize class k^* , regardless of the server

state. While it is tempting to conjecture that prioritizing class k^* in this situation is optimal, the following example shows that doing so could in fact be *very* suboptimal.

Example 1. Let $K = 2$. Suppose jobs of classes 1 and 2 arrive according to Poisson processes with rates $\lambda_1 = 5$ and $\lambda_2 = 0.8$, respectively. There are only two server states, $q(1|2) = 1$, and the inter-deterioration times are exponential with rate 1. The service rates are $\mu_1^1 = \mu_1^2 = 10$ and $\mu_2^s = s$ for $s = 1, 2$. Finally, corrective maintenance occurs instantaneously.

Under both the state-dependent and average $c\mu$ -rules, class 1 is given priority regardless of the server state. Note that under this policy, the average service rate for class 2 jobs is $0.5(0.5 \cdot 2 + 0.5 \cdot 1) = 0.75 < 0.8 = \lambda_2$. This indicates that, regardless of the initial state, the system is unstable and has infinite long-run expected average cost. A formal proof that the system is unstable when class 1 is always prioritized is given in Appendix [A.1](#).

On the other hand, consider the policy that prioritizes class s jobs when the server state is s , for $s = 1, 2$. Since $\lambda_2 = 0.8 < (0.5)(2) = 1$, this policy incurs a finite long-run expected average cost regardless of the initial state. This can be proved by showing that the associated fluid model is stable; see Appendix [A.2](#) for details.

3.3 Optimality of $c\mu$ -Rules

The following condition guarantees that it is optimal to prioritize one class over another.

Assumption CR (Constant Ratio). Every state $s \in \{1, \dots, S\}$ (i.e., where the server is online), as well as every pair of classes i, j , satisfies

$$\mu_i^{s-1} \mu_j^s = \mu_i^s \mu_j^{s-1}. \quad (2)$$

Assumption **CR** states that the *ratio* of the service rates for class i and class j jobs remains constant as the server changes state. It can be interpreted as saying that the different service capabilities of the flexible server are *affected equally by deterioration*. Note that Assumption **CR** implies, but is not equivalent to, the condition that $c_1 \mu_1^s \geq c_2 \mu_2^s \geq \dots \geq c_K \mu_K^s$ for every server

state s . For instance, Example 1 satisfies $c_1\mu_1^s \geq c_2\mu_2^s$ for every $s \in \{0, 1, \dots, S\}$, but does not satisfy Assumption CR.

The following theorem, which is the main result of this section, states that Assumption CR guarantees the optimality of the state-dependent and average $c\mu$ -rules described at the beginning of this section.

Theorem 2. *If Assumption CR holds, then the $c\mu$ -rule is optimal.*

Proof. We use Assumption CR to adapt the interchange argument in Nain [18, Proof of Theorem 2.1] to our setting.

The first step is to show that, for every $T \geq 0$, the problem of minimizing the finite-horizon expected weighted queue lengths

$$\mathbb{E} \left[\int_0^T \sum_{k=1}^K c_k Q_k^\pi(t) dt \right] \quad (3)$$

can be reduced to a reward-maximization problem that is amenable to analysis via an interchange argument. To do this, we define some processes of interest. Consider any policy $\pi \in \Pi$ and fixed time $t \in [0, \infty)$. Let $U^\pi(t)$ denote the job class that the server is assigned to at time t under the policy π , and let¹

$$a_k^\pi(t) := \mathbf{1}\{Q_k^\pi(t-) > 0, U^\pi(t) = k\}, \quad k = 1, 2, \dots, K.$$

Also, recalling that we are considering the case of no preventive maintenance, let $S(t)$ denote the state of the server at time t , and let

$$\phi^\pi(t) := \int_0^t \sum_{k=1}^K c_k \mu_k^{S(u)} a_k^\pi(u) du.$$

To reduce the problem of minimizing (1) to that of maximizing

$$\mathbb{E} \left[\int_0^T \phi^\pi(t) dt \right], \quad (4)$$

¹Given a function $f : [0, \infty) \rightarrow \mathbb{R}$, let $f(t-) := \lim_{u \uparrow t} f(u)$ for $t > 0$.

consider the queue-length processes $Q_k^\pi(t)$, $k = 1, 2, \dots, K$, under π , and let $A_k(t)$ denote the cumulative number of class k arrivals during the time interval $[0, t]$. Using an argument analogous to that in [18, Proof of Lemma 2.1] (replace μ_k with $\mu_k^{S(u)}$ and the fact that [4, Partial Result, p. 24] holds for Poisson processes with rates that depend on $S(t)$),

$$\mathbb{E} \left[\int_0^T \sum_{k=1}^K c_k Q_k^\pi(t) dt \right] = \mathbb{E} \left[\int_0^T \sum_{k=1}^K c_k [Q_k(0) + A_k(t)] dt \right] - \mathbb{E} \left[\int_0^T \phi^\pi(t) dt \right]. \quad (5)$$

Since the first term on the right-hand side of (5) does not depend on π , it follows that minimizing (3) is equivalent to maximizing $\mathbb{E} \left[\int_0^T \phi^\pi(t) dt \right]$.

The next step is to show that the $c\mu$ -rule, denoted by $\pi_{c\mu}$, maximizes (4) for every finite horizon $T \geq 0$, i.e., that

$$\mathbb{E} \left[\int_0^T \phi^{\pi_{c\mu}}(t) dt \right] \geq \mathbb{E} \left[\int_0^T \phi^\pi(t) dt \right]. \quad (6)$$

for all $\pi \in \Pi$ and $T \geq 0$. To prove this, we use a slight modification of the sample-path-based construction in [18, Proof of Theorem 2.1] to show that in fact,

$$\int_0^T \phi^{\pi_{c\mu}}(t) dt \geq \int_0^T \phi^\pi(t) dt \quad (7)$$

holds with probability 1 (written w.p.1).

Noting that every policy is optimal if $T = 0$, fix $\pi \in \Pi$ and $T > 0$. Suppose Assumption CR holds, and assume the job classes are numbered so that $c_1 \mu_1^s \geq c_2 \mu_2^s \geq \dots \geq c_K \mu_K^s$ (see the comments following the definition of Assumption CR). Consider the random time

$$\sigma_\pi := \inf\{t > 0 \mid Q_1^\pi(t-) > 0, U^\pi(t) \neq 1\},$$

which denotes the first time that the policy π does not follow the $c\mu$ -rule. If $\sigma_\pi \geq T$ w.p.1, then (7) holds w.p.1., since in this case the policy π follows the $c\mu$ -rule during $[0, T]$.

On the other hand, suppose $\sigma_\pi < T$ with positive probability, in which case there is a positive probability with which π does not follow the $c\mu$ -rule during $[0, T]$. Let Task A denote the class 1 job that is assigned to the server at time σ_π under the $c\mu$ -rule and Task B be the class ℓ job the server is assigned to under π .

Let $\Pi_\infty \supset \Pi$ denote the set of all possibly anticipative policies. We will now construct a policy $\pi_+ \in \Pi_\infty$ that follows the $c\mu$ -rule at time σ_π and satisfies

$$\int_0^T \phi^{\pi_+}(t) dt \geq \int_0^T \phi^\pi(t) dt \quad \text{w.p.1.} \quad (8)$$

First, for all times $t \in [0, \sigma_\pi)$, let π_+ coincide with the $c\mu$ -rule. To define π_+ for times $t > \sigma_\pi$, consider the random variable

$$\sigma_\pi^* := \min\{T_A, \tau_\pi\},$$

where T_A denotes the amount of time required to complete Task A when the server state is $S(\sigma_\pi)$, and $\tau_\pi := \inf\{t \geq \sigma_\pi \mid U^\pi(t) \neq \ell \text{ or } S(t) \neq S(\sigma_\pi)\}$ is the first time after time σ_π that either under π the server stops working on Task B or the server changes state. Assume the policy π_+ works on Task A during $[\sigma_\pi, \sigma_\pi + \sigma_\pi^*)$.

To complete the ‘‘interchange’’ of Tasks A and B, we will complete the definition of π_+ so that after some time ν_π , the number of queued jobs under both π_+ and π agree w.p.1. In particular, let ν_π denote the time when Task A is completed under the policy π . During $[\sigma_\pi + \sigma_\pi^*, \nu_\pi)$, let π_+ mimic the actions taken under π with the following exception: Whenever π works on Task A, but Task A has already been completed under π_+ , the latter policy works on Task B instead. Finally, let π_+ mimic the actions taken under π at all times $t \geq \nu_\pi$.

We claim that at time ν_π , both the queue lengths and the amount of work remaining in the system are the same under both π and π_+ . To verify this, let

$$\kappa := \frac{\mu_1^{s-1}}{\mu_1^s} = \frac{\mu_\ell^{s-1}}{\mu_\ell^s}, \quad s \geq 1,$$

and let $I_n = [\theta_n, \theta'_n)$ denote the n^{th} time interval in $[\sigma_\pi + \sigma_\pi^*, \nu_\pi)$ during which π_+ serves class ℓ while π serves class 1 and the server state does not change. Observe that under π_+ , the amount of work done on Task A during $[\sigma_\pi, \sigma_\pi + \sigma_\pi^*)$ is the same as the amount of work done on this job during $\cup_n I_n$ under π . Note that by Assumption **CR**,

$$\mu_k^r = \kappa^{s-r} \mu_k^s, \quad k = 1, \ell, r, s \geq 1; \quad (9)$$

this is because when $r > s \geq 1$,

$$\mu_k^r = \frac{\mu_k^r}{\mu_k^{r-1}} \dots \frac{\mu_k^{r-(r-s-1)}}{\mu_k^s} = \kappa^{s-r} \mu_k^s,$$

and when $s > r \geq 1$,

$$\mu_k^r = \frac{\mu_k^r}{\mu_k^{r+1}} \dots \frac{\mu_k^{r+(r-s-1)}}{\mu_k^s} = \kappa^{s-r} \mu_k^s.$$

Using (9), the amount of work done on Task A during $[\sigma_\pi, \sigma_\pi + \sigma_\pi^*)$ can be written as

$$\mu_1^{S(\sigma_\pi)} \sigma_\pi^* = \sum_n \mu_1^{S(\theta_n)} (\theta'_n - \theta_n) = \mu_1^{S(\sigma_\pi)} \sum_n \kappa^{S(\sigma_\pi) - S(\theta_n)} (\theta'_n - \theta_n). \quad (10)$$

From (10), we conclude that

$$\sigma_\pi^* = \sum_n \kappa^{S(\sigma_\pi) - S(\theta_n)} (\theta'_n - \theta_n). \quad (11)$$

Hence the amount of work that is done on Task B during $\cup_n I_n$ under π_+ is

$$\sum_n \mu_\ell^{S(\theta_n)} (\theta'_n - \theta_n) = \mu_\ell^{S(\sigma_\pi)} \sum_n \kappa^{S(\sigma_\pi) - S(\theta_n)} (\theta'_n - \theta_n) = \mu_\ell^{S(\sigma_\pi)} \sigma_\pi^*,$$

which is precisely the amount of work done on Task B during $[\sigma_\pi, \sigma_\pi + \sigma_\pi^*)$ under the original policy π . Since π_+ selects exactly the same actions as π at all times $t \in [\sigma_\pi + \sigma_\pi^*, \nu_\pi) \setminus \cup_n I_n$, it follows that both the queue lengths and the remaining amount of work in the system at time ν_π are the same under both π and π_+ .

Since the policies π and π_+ couple at time ν_π , the validity of (8) and the optimality of the $c\mu$ -rule for every finite horizon T can be proved by following [18, Proof of Theorem 2.1] and using the preceding definitions of $\phi^\pi(t)$, σ_π^* , and the intervals $[\theta_n, \theta'_n)$. It follows *a fortiori* that the $c\mu$ -rule is optimal under the average-cost criterion (1). \square

Remark 3. *The proof of Theorem 2 does not rely on the assumption that deterioration events always send the server to a state that is worse. In particular, it holds when the service rates are simply assumed to be modulated (not necessarily in a Markovian way) according to the point process that describes the deterioration process. Hence the proof of Theorem 2 implies that, for a two-class G/M/1 queue with modulated service rates that satisfy Assumption CR, it is optimal to schedule according to the $c\mu$ -rule.*

4 Scheduling with Preventive Maintenance

We now consider the problem of optimally making both scheduling and preventive maintenance decisions. In Section 4.1, we provide conditions under which it suffices to schedule jobs according to the state-dependent $c\mu$ -rule, or the average $c\mu$ -rule, described in Section 3.1. Then, in Section 4.2 we present conditions under which optimal maintenance decisions are monotone in the server's health. When the conditions hold, these results simplify the computation of optimal policies. At the same time, when one or more of the conditions do not hold, they suggest heuristics that may still perform well. The performance of scheduling with the $c\mu$ -rule, when the conditions of Theorem 4 in this section do not hold, is considered numerically in Section 5.

4.1 Optimal Scheduling

In this section, a *maintenance policy* is a rule that stipulates, given the current state of the system, whether or not to initiate maintenance. If maintenance is not initiated, a *scheduling policy* determines which customer class (if any) should be served. The set of all stationary deterministic maintenance policies is identified with the set of all functions $f : \{0, 1, \dots\}^K \times \{0, 1, \dots, S\} \rightarrow \{0, 1\}$ where $f(\mathbf{i}, s) = 1$ (resp. $= 0$) if and only if the maintenance policy f calls for maintenance to be initiated (resp. no maintenance) when the state is (\mathbf{i}, s) . Note that $f(\mathbf{i}, 0) = 1$ for all \mathbf{i} .

According to Theorem 2 in Section 3, if Assumption CR holds, then the $c\mu$ -rule is the optimal scheduling policy in the presence of a deteriorating server that cannot be preventively maintained. In the context of joint scheduling and maintenance, Theorem 2 can be generalized to Theorem 4 below. To state this theorem, a maintenance policy f is said to be *queue-oblivious* if there exists a function $g : \{0, 1, \dots, S\} \rightarrow \{0, 1\}$ satisfying

$$f(\mathbf{i}, s) = g(s) \quad \text{for all } (\mathbf{i}, s) \in \{0, 1, \dots\}^K \times \{0, 1, \dots, S\}.$$

In other words, a queue-oblivious maintenance policy is stationary, deterministic, and does not depend on any queue-length information. Examples of queue-oblivious maintenance policies include *server threshold* policies (where the server is maintained if and only if its state is below

a certain threshold), *job-based* policies (where the server is maintained whenever a certain fixed number of jobs have been completed), and *calendar-based* policies (where the server is maintained whenever a certain fixed amount of time has elapsed).

Theorem 4. *Suppose Assumption CR holds. Then under any queue-oblivious maintenance policy, it is optimal to schedule according to the $c\mu$ -rule. In particular, consider any server state $s \geq 1$, and assume that the current state is one in which the maintenance policy calls for no maintenance. If*

$$c_1\mu_1^s \geq c_2\mu_2^s \geq \dots \geq c_K\mu_K^s, \quad (12)$$

then Assumption CR implies that (12) holds for all $s \in \{0, 1, \dots, S\}$, and it is optimal to prioritize the classes in the order $1, 2, \dots, K$.

Proof. Under a queue-oblivious maintenance policy, the evolution of the server state does not depend on how the jobs are served. The theorem then follows from the proof of Theorem 2, which does not require any assumptions on where the server state transitions to when deterioration events occur (see Remark 3). □

Theorem 4 immediately implies the following theorem, which is the main result of this section.

Theorem 5. *If Assumption CR holds, and the decision-maker is restricted to queue-oblivious maintenance policies, then it is without loss of optimality to only consider joint scheduling and maintenance policies where jobs are scheduled according to the (static) priority policy described in Theorem 4.*

4.2 Optimal Maintenance Decisions

Up to this point, we have only assumed that the arrival processes are described by independent point processes on \mathbb{R}_+ . Under Assumption M below, the problem can be formulated as a *semi-Markov decision process (SMDP)*. The main result in this section (Theorem 6) states that under this assumption and Assumption S below, the search for an optimal policy can be restricted to policies that are monotone in the server's health.

Assumption M (Markovian Arrivals and Deterioration).

- (i) The point process modeling the arrival times of jobs of class k is a Poisson process with rate λ_k . Furthermore, the K arrival processes are mutually independent.
- (ii) The server deteriorates according to a continuous-time Markov chain. In particular, if its current state is $s \in \{1, \dots, S\}$, then the time until the next deterioration event is exponentially distributed with rate $\alpha_s > 0$.
- (iii) The maintenance times (i.e., the times that the server spends in the offline state) are independent and identically distributed with common distribution $G(\cdot)$ whose mean $1/\alpha_0 := \int_0^\infty t dG(t)$ satisfies $0 < 1/\alpha_0 < \infty$.

Assumption S (Stability).

- (i) There is a server state $s^* \in \{1, \dots, S\}$ satisfying

$$\sum_{k=1}^K \frac{\lambda_k}{\sum_{s=s^*}^S (\mu_k^s / \alpha_s)} < \frac{1}{(1/\alpha_0) + \sum_{s=s^*}^S (1/\alpha_s)}.$$

- (ii) The server can only deteriorate to the next-worse state, i.e.,

$$q(s-1|s) = 1 \quad \forall s \geq 1.$$

A joint scheduling and maintenance policy is *monotone in the server's health* if, for every fixed number of jobs of each class in the system, PM is initiated whenever the server's health state is sufficiently low. The following proposition states that under Assumptions **M** and **S**, one can restrict the search for an optimal joint scheduling and maintenance policy to deterministic stationary policies that are monotone in the server's health. A proof is provided in Appendix **A.3**.

Proposition 6. *Suppose Assumptions **M** and **S** hold. Then there exists an optimal joint scheduling and maintenance policy that is deterministic, stationary, and monotone in the server's health.*

Combining the conclusions of Propositions **4** and **6** leads to the following theorem, which is the main result in this section.

Theorem 7. *Suppose Assumptions CR, M, and S hold. Then there exists an optimal deterministic stationary policy that is both monotone in the server’s health and schedules jobs according to the $c\mu$ -rule.*

Remark 8. *Under Assumptions M and S, there may not be an optimal policy that is monotone in the queue lengths. In particular, letting $K = 2$, $K_c = K_p = 0$, $c_1 = c_2 = 1$, $\lambda_1 = \lambda_2 = 1$, $S = 4$, $\mu_1^1 = \mu_2^1 = 1/2$, $\mu_1^2 = \mu_2^2 = 1$, $\mu_1^3 = \mu_2^3 = 3/2$, $\mu_1^4 = \mu_2^4 = 2$, and $\alpha_s = 1/5$ for $s = 0, 1, 2, 3, 4$, we obtain the model instance considered in [16, Example 3.6]. It was shown in [16] that the optimal policy for this model is such that, for server state 2, initiating maintenance is optimal when there are no jobs, not optimal when there are 1 to 11 jobs, and optimal when there are more than 11 jobs.*

5 Numerical Experiments

In this section, we numerically examine policy performance as the problem parameters vary. Specifically, we consider variations in holding/maintenance costs (Section 5.1.2), maintenance rate (Section 5.1.3), the degree of service capacity deterioration per server state change (Section 5.1.4), and the variability of inter-event (i.e., arrival, service, deterioration, or maintenance) times (Section 5.1.5). In the simulation model, we relax a number of assumptions used in the preceding sections. This includes departures from Assumption CR, a lack of the non-preemptive service, class-dependent deterioration rates α_k^s , and allowing for non-exponential inter-arrival, service, and inter-deterioration time distributions.

5.1 Simulation Setup

We restrict attention to three job classes, three server states (2 = like-new, 1 = deteriorated, and 0 = failed), and identical fixed predictive and corrective maintenance costs (i.e., $K_p = K_c$). The assumption of equal predictive and corrective maintenance costs was made in part to limit the number of varying parameters, and is reasonable for modeling systems where the additional cost due to unexpected failure is small. Service is assumed to be non-preemptive, and both job classes have unlimited buffers.

5.1.1 Model Parameters

A total of 21 parameters are used to specify the sets of model instances considered in Sections 5.1.2, 5.1.3, 5.1.4, and 5.1.5. Each parameter has three levels, corresponding to “low”, “moderate”, and “high”. The sets of parameter values were selected to balance the objective of examining policy performance as certain parameters vary, with the computational resources available for this study. Sections 5.1.2, 5.1.3, 5.1.4, and 5.1.5 provide details on how the parameter values were selected.

The first group of parameters are the system utilization and arrival rate levels $\tilde{\lambda}_k$ for each class k , in jobs per hour. Here, the “utilization level” ρ should be interpreted as stipulating that there exists a policy such that the utilization level for each job class (i.e., the arrival rate divided by the average service rate) does not exceed ρ .

Level	Utilization	$\tilde{\lambda}_1$	$\tilde{\lambda}_2$	$\tilde{\lambda}_3$
Low	0.4	1	1	1
Moderate	0.6	5	5	5
High	0.8	10	10	10

Table 1: The utilization is the arrival rate divided by the average service rate, taken to be the same for each class. The arrival rate levels $\tilde{\lambda}_k$ are in units of jobs per hour.

Note that whether or not the numerical values of the selected arrival rate levels are consistent with the selected utilization level will depend on the average service rate, which in turn depends on the service, deterioration, and maintenance rates. For example, if $\rho = 0.8$ and $\tilde{\lambda}_1 = \tilde{\lambda}_2 = \tilde{\lambda}_3 = 10$ are selected, then the average service rate should be at least $10/0.8 = 12.5$. To account for this, we select all of the parameter levels first, and then modify the arrival rates to ensure the existence of a policy with sufficiently large average service rates. This is accomplished by solving a linear program (LP). Appendix A.4 contains the formulation of this LP, and specifies how a solution to it is used to determine the modified arrival rates.

The average service rate for class k jobs is μ_k^s jobs per hour, when the server state is s . The levels for the like-new service rates, and the ratios between the like-new and deteriorated service

rates, are given in Table 2.

Level	μ_1^2	μ_2^2	μ_3^2	$\frac{\mu_1^1}{\mu_1^2}$	$\frac{\mu_2^1}{\mu_2^2}$	$\frac{\mu_3^1}{\mu_3^2}$
Low	1	1	1	0.01	0.01	0.01
Moderate	5	5	5	0.5	0.5	0.5
High	10	10	10	0.99	0.99	0.99

Table 2: For each class k , the like-new service rate is μ_k^2 and the deteriorated service rate is μ_k^1 .

The average deterioration rate when the server is in state s and working on a class k job is denoted by α_k^s , and the average maintenance rate is denoted by α_0 . Moreover, deteriorations are assumed to always transition the server to the next-worse state, i.e., $q(1|2) = q(0|1) = 1$. The levels for the like-new deterioration rates, the ratios between the like-new and deteriorated deterioration rates, and the maintenance rate are given in Table 3 in units of deterioration events per hour.

Level	α_1^2	α_2^2	α_3^2	$\frac{\alpha_1^1}{\alpha_1^2}$	$\frac{\alpha_2^1}{\alpha_2^2}$	$\frac{\alpha_3^1}{\alpha_3^2}$	α_0
Low	0.01	0.01	0.01	1.01	1.01	1.01	0.01
Moderate	0.5	0.5	0.5	1.5	1.5	1.5	0.5
High	1	1	1	2	2	2	1

Table 3: For each class k , the deterioration rate when the server is like-new and working on class k jobs is α_k^2 , and is α_k^1 when the server is deteriorated.

Finally, the levels for the holding cost rates c_k and the fixed maintenance cost $K_c = K_p$ are given in Table 4.

Level	c_1	c_2	c_3	$K_c = K_p$
Low	0.5	0.5	0.5	5
Moderate	1	1	1	10
High	2	2	2	20

Table 4: For each class k , the holding cost rate is c_k . The preventive maintenance cost is K_p , and the corrective maintenance cost is assumed to be $K_c = K_p$.

In Sections 5.1.2, 5.1.3, 5.1.4, and 5.1.5, specific combinations of the above parameter levels are used to estimate the effect of varying, respectively, the holding and maintenance costs, the maintenance rate, the degree of capacity deterioration, and the variability of the inter-event distributions.

5.1.2 Holding vs. Maintenance Costs

We first consider the effect of jointly varying the holding cost rates c_k and the fixed maintenance cost $K_p = K_c$. The parameter values were determined by first randomly selecting the levels of the parameters other than the holding cost rates and fixed maintenance cost. Then, the holding cost rates and fixed maintenance cost were jointly varied. The actual parameter values, obtained by scaling the selected arrival rate levels to ensure the existence of a stable policy, are given in Table 5. All inter-event (i.e., inter-arrival, service, inter-deterioration, and maintenance) times are assumed to be exponentially distributed.

#	Utilization	λ_1	λ_2	λ_3	μ_1^2	μ_2^2	μ_3^2	μ_1^1	μ_2^1	μ_3^1	α_1^2	α_2^2	α_3^2	α_1^1	α_2^1	α_3^1	α_0	c_1	c_2	c_3	$K_p = K_c$
1	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	0.5	2	5
2	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	2	2	5
3	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	0.5	2	10
4	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	2	2	10
5	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	0.5	2	20
6	0.6	0.22	1.1	1.1	5	5	10	2.5	4.95	5	0.01	1	0.01	0.0101	1.01	0.0101	0.5	0.5	2	2	20

Table 5: In the study of holding vs. maintenance costs, a total of six sets of parameter values were selected according to the procedure described in Section 5.1.2.

5.1.3 Maintenance Rates

Next, we will consider the effect of varying the maintenance rate α_0 , i.e., the rate at which the server is restored from the failed state 0 to the like-new state 2. The parameter values were determined by first randomly selecting the levels of the parameters other than the maintenance rate. Then, the maintenance rate was varied. The actual arrival rates were obtained by scaling the selected arrival rate levels to ensure the existence of a stable policy, resulting in distinct arrival rates for each possible maintenance rate. The parameter values are given in Table 6. All inter-event (i.e., inter-arrival, service, inter-deterioration, and maintenance) times are assumed to be exponentially distributed.

#	Utilization	λ_1	λ_2	λ_3	μ_1^2	μ_2^2	μ_3^2	μ_1^1	μ_2^1	μ_3^1	α_1^2	α_2^2	α_3^2	α_1^1	α_2^1	α_3^1	α_0	c_1	c_2	c_3	$K_p = K_c$
1	0.4	0.038	0.019	0.019	10	5	5	0.1	2.5	2.5	0.5	0.5	1	1	0.75	1.01	0.01	2	0.5	1	5
2	0.4	0.72	0.36	0.36	10	5	5	0.1	2.5	2.5	0.5	0.5	1	1	0.75	1.01	0.5	2	0.5	1	5
3	0.4	0.89	0.44	0.44	10	5	5	0.1	2.5	2.5	0.5	0.5	1	1	0.75	1.01	1	2	0.5	1	5

Table 6: In the study of varying maintenance rates, a total of three sets of parameter values were selected according to the procedure described in Section 5.1.3.

5.1.4 Degree of Service Capacity Deterioration

We also consider the effect of varying the degree of service rate deterioration for one of the classes. The parameter values were determined by first randomly selecting the levels of the parameters other than $\frac{\mu_3^1}{\mu_3}$. Then, the value of $\frac{\mu_3^1}{\mu_3}$ was varied across its three levels of 0.01, 0.5, and 0.99. The parameter values are given in Table 7. All inter-event (i.e., inter-arrival, service, inter-deterioration, and maintenance) times are assumed to be exponentially distributed.

#	Utilization	λ_1	λ_2	λ_3	μ_1^2	μ_2^2	μ_3^2	μ_1^1	μ_2^1	μ_3^1	α_1^2	α_2^2	α_3^2	α_1^1	α_2^1	α_3^1	α_0	c_1	c_2	c_3	$K_p = K_c$
1	0.6	1.03	0.21	1.03	5	5	10	2.5	2.5	0.1	0.5	0.5	0.01	1	0.75	0.0101	0.5	1	1	1	10
2	0.6	1.03	0.21	1.03	5	5	10	2.5	2.5	5	0.5	0.5	0.01	1	0.75	0.0101	0.5	1	1	1	10
3	0.6	1.03	0.21	1.03	5	5	10	2.5	2.5	9.9	0.5	0.5	0.01	1	0.75	0.0101	0.5	1	1	1	10

Table 7: In the study of varying degrees of service capacity deterioration, a total of three sets of parameter values were selected according to the procedure described in Section 5.1.4.

5.1.5 Variability of Inter-Event Distributions

Finally, we consider the effect of varying the coefficient of variation of the inter-event times on policy performance, where the inter-event times are assumed to be gamma distributed. The parameter values were determined by first randomly selecting the levels for the 21 parameters indicated in Section 5.1.1, and scaling the $\tilde{\lambda}_k$'s using the solution to the linear program in Appendix A.4. To specify the gamma inter-event distributions, the selected rates (e.g., the μ_k^s 's) are taken to be the reciprocals of the corresponding expected inter-event times. To complete the specification of the relevant gamma distributions, the coefficient of variation (CV) of each inter-event time distribution was systematically varied between 0.1 (low variability) and 2 (high variability). The parameter values are given in Table 8.

#	Utilization	λ_1	λ_2	λ_3	μ_1^2	μ_2^2	μ_3^2	μ_1^1	μ_2^1	μ_3^1	α_1^2	α_2^2	α_3^2	α_1^1	α_2^1	α_3^1	α_0	c_1	c_2	c_3	$K_p = K_c$
1 – 16	0.6	0.09	0.94	0.94	5	10	5	2.5	5	0.05	0.5	0.5	0.5	1	0.505	1	0.5	1	2	0.5	10

#	Inter-Arrival CV	Service CV	Inter-Deterioration CV	Maintenance CV
1	0.1	0.1	0.1	0.1
2	2	0.1	0.1	0.1
3	0.1	2	0.1	0.1
4	0.1	0.1	2	0.1
5	0.1	0.1	0.1	2
6	2	2	0.1	0.1
7	2	0.1	2	0.1
8	2	0.1	0.1	2
9	0.1	2	2	0.1
10	0.1	2	0.1	2
11	0.1	0.1	2	2
12	2	2	2	0.1
13	2	2	0.1	2
14	2	0.1	2	2
15	0.1	2	2	2
16	2	2	2	2

Table 8: In the study of the variability of the inter-event distributions, a total of sixteen sets of parameter values were selected according to the procedure described in Section 5.1.5. Here, CV stands for coefficient of variation, and all of the inter-event times are gamma distributed.

5.2 Policies

We consider the performance of the following five policies.

$c\mu$ -Rule (CMU): This policy involves scheduling according to the $c\mu$ -rule, and only performing maintenance when the server fails.

Longest-Waiting-Time-First (LWF): This policy involves scheduling according to which class has experienced the longest total waiting time so far, and only performing maintenance when the server fails.

$c\mu$ -Rule with Preventive Maintenance (CMUPM): This policy involves scheduling according to the $c\mu$ -rule, and performing preventive maintenance whenever the server transitions from like-new to deteriorated.

LWF with Preventive Maintenance (LWFPM): This policy involves scheduling according to LWF, and performing preventive maintenance whenever the server transitions from like-new to deteriorated.

MDP-Based Policy (DDPM): This is the optimal joint scheduling and maintenance policy for the discrete-time MDP with finite state and action sets obtained by truncating (at a buffer size of 50 per class), assuming that the inter-arrival, service, inter-deterioration, and maintenance times are exponentially distributed, and uniformizing the associated continuous-time MDP where service preemptions are allowed. This is the most computationally-intensive policy.

5.3 Simulation Results

For each of the five heuristics, 30 replications for each parameter set were performed over a simulation time horizon of 1 year (recall that the inter-event rates are per-hour). The results are summarized in the following sub-sections.

A key takeaway from the simulation results is that in all of the parameter sets considered, scheduling according to the $c\mu$ -rule can lead to performance comparable to the MDP-based policy, when preventive maintenance is done properly. We emphasize that in all but one of the parameter sets on which this observation is based, Assumption **CR** does not hold. Moreover, the presence of class-dependent deterioration rates represents a significant complication of the model analyzed in Section 4, where all of the deterioration rates only depend on the current server state.

While an optimal preventive maintenance policy will generally be complicated (see Kaufman and Lewis [16]), the simulation results indicate that fixing the scheduling rule to be the $c\mu$ -rule, and then searching for a good preventive maintenance policy, can lead to good performance. In particular, for all of the parameter sets considered, the best policy with $c\mu$ -based scheduling performed comparably, and often much better than, the best policy with longest-waiting-time-first scheduling. Moreover, the variability of the performance of $c\mu$ -based scheduling was much lower than that of longest-waiting-time-first scheduling; there were several parameter sets, such as those considered in Section 5.3.4, where combining longest-waiting-time-first scheduling with a poor maintenance rule is almost an order of magnitude worse than making the same mistake under $c\mu$ -based scheduling.

5.3.1 Holding vs. Maintenance Costs

For this set of parameter values, the priority ordering of job classes under the $c\mu$ -rule is class 3, then class 2, then class 1, regardless of the server state; see Table 9. Moreover, note that the constant ratio assumption (Assumption **CR**) does not hold for this set of parameter values; in particular,

$$\frac{\mu_1^1}{\mu_1^2} = \frac{\mu_3^1}{\mu_3^2} = 0.5, \text{ while } \frac{\mu_2^1}{\mu_2^2} = 0.99.$$

#	$c_1\mu_1^2$	$c_2\mu_2^2$	$c_3\mu_3^2$	$c_1\mu_1^1$	$c_2\mu_2^1$	$c_3\mu_3^1$	$\frac{\mu_1^1}{\mu_1^2}$	$\frac{\mu_2^1}{\mu_2^2}$	$\frac{\mu_3^1}{\mu_3^2}$
1, 3, 5	2.5	2.5	20	1.25	2.475	10	0.5	0.99	0.5
2, 4, 6	2.5	10	20	1.25	9.9	10	0.5	0.99	0.5

Table 9: This table contains the indices $c_k\mu_k^s$ used by the $c\mu$ -rule for each of the six sets of parameter values used in the study of holding vs. maintenance costs, along with the service rate ratios μ_k^1/μ_k^2 .

Figure 1 summarizes the observed performance of the five policies. Figure 1 indicates that the priority ordering based on the $c\mu$ -rule, without preventive maintenance, performs comparably to the MDP-based policy, and clearly outperforms the other policies. In contrast, LWF-based scheduling can perform poorly if preventive maintenance is done, doing significantly worse than even $c\mu$ -based scheduling with preventive maintenance. For this set of parameter values, the maintenance rate of $\alpha_0 = 0.5$ is low enough that the gains in service rates from doing preventive maintenance are more than offset by the holding costs incurred by queued class 2 and class 3 jobs, which account for the vast majority of the arrivals ($\lambda_2 = \lambda_3 = 1.1$, while $\lambda_1 = 0.22$) and have the highest holding cost rates ($c_1 = 0.5$, while $c_2 \geq 0.5$ and $c_3 = 2$). Notably, Figure 1 indicates that the relative performance of the five policies remains roughly the same as the holding and maintenance costs varied.

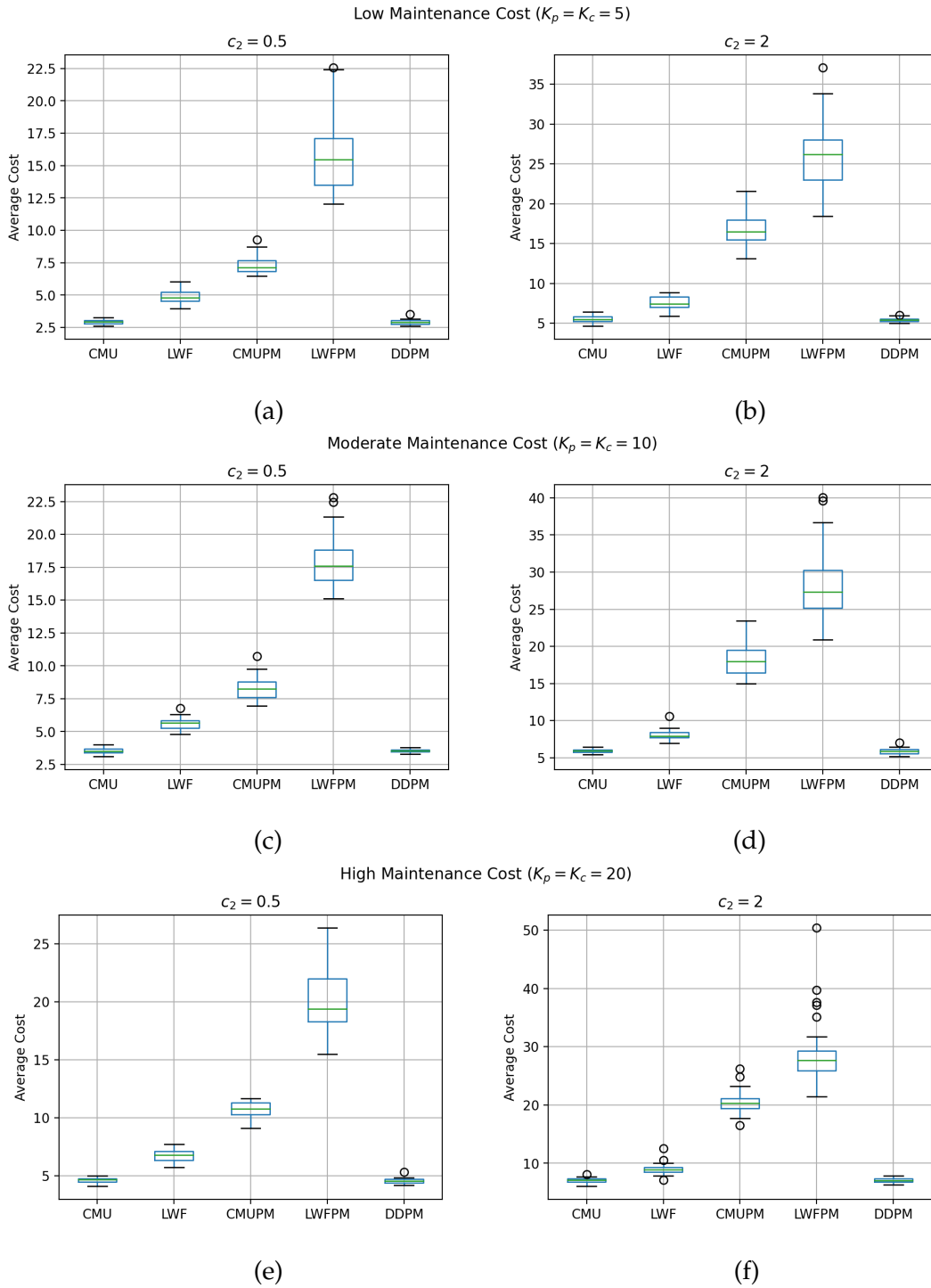


Figure 1: Policy Performance for Holding vs. Maintenance Costs (In all cases, $c_1 = 0.5$ and $c_3 = 2$.)

5.3.2 Maintenance Rates

For this set of parameter values, the priority ordering of job classes under the $c\mu$ -rule changes as the server deteriorates. In particular, class 1 has the highest priority when the server is like-new, but lowest priority when the server is deteriorated; see Table 10. This priority reversal due in part to the severe amount of service capacity loss when deterioration occurs; see Table 6. Specifically, deterioration results in a 99% service rate reduction for class 1 jobs, and a 50% reduction for class 2 and 3. Table 10 also indicates that the constant ratio assumption (Assumption CR) does not hold for this set of parameter values.

#	$c_1\mu_1^2$	$c_2\mu_2^2$	$c_3\mu_3^2$	$c_1\mu_1^1$	$c_2\mu_2^1$	$c_3\mu_3^1$	$\frac{\mu_1^1}{\mu_1^2}$	$\frac{\mu_2^1}{\mu_2^2}$	$\frac{\mu_3^1}{\mu_3^2}$
1, 2, 3	20	2.5	5	0.2	1.25	2.5	0.01	0.5	0.5

Table 10: This table contains the indices $c_k\mu_k^s$ used by the $c\mu$ -rule for each of the three sets of parameter values used in the study of varying maintenance rates, along with the service rate ratios μ_k^1/μ_k^2 .

Figure 2 summarizes the observed performance of the five policies. The results reflect the expectation, in light of the severe degree of deterioration, that the two policies that do not employ preventive maintenance (i.e., CMU and LWF) will perform increasingly poorly as the maintenance rate increases. The relative performance of the remaining three policies was roughly constant with increasing maintenance rate; CMUPM performed comparably with the MDP-based policy (DDPM), while LWFPM tended to perform worse.

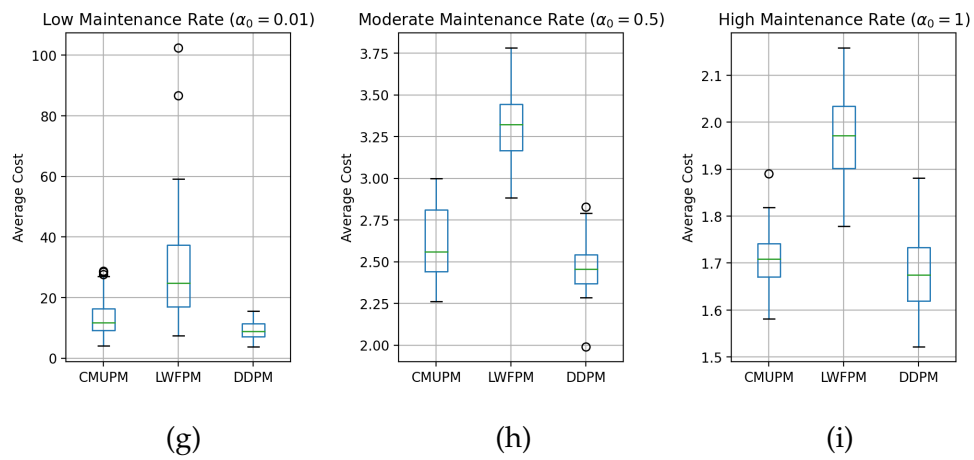
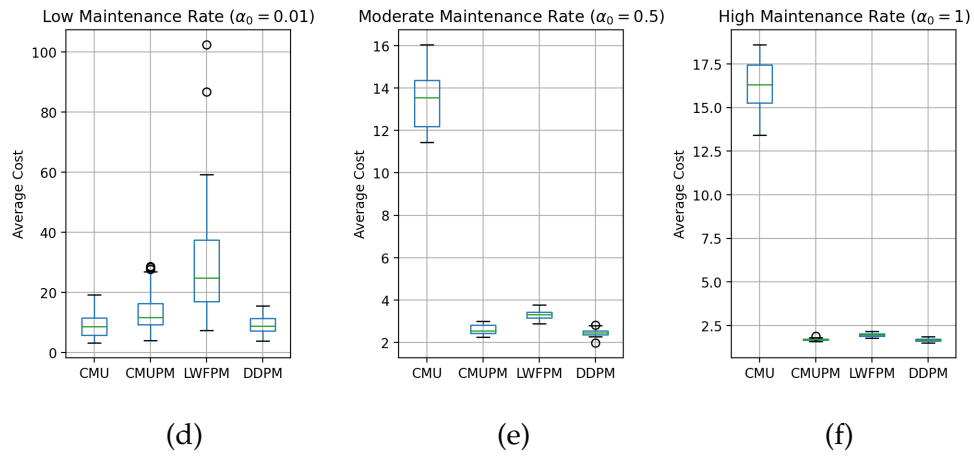
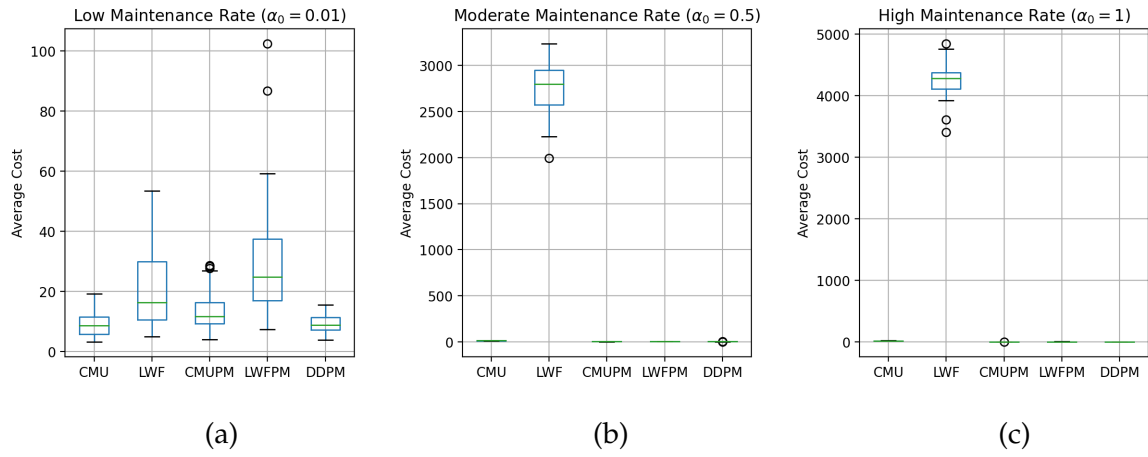


Figure 2: Policy Performance: Varying Maintenance Rates

5.3.3 Degree of Service Capacity Deterioration

For this set of parameter values, the priority ordering of job classes under the $c\mu$ -rule changes as the server deteriorates; the priority indices are given in Table 11. The ratios between deteriorated and like-new service rates are also given in Table 11; note that Assumption CR only holds for the second set of parameter values.

#	$c_1\mu_1^2$	$c_2\mu_2^2$	$c_3\mu_3^2$	$c_1\mu_1^1$	$c_2\mu_2^1$	$c_3\mu_3^1$	$\frac{\mu_1^1}{\mu_1^2}$	$\frac{\mu_2^1}{\mu_2^2}$	$\frac{\mu_3^1}{\mu_3^2}$
1	5	5	10	2.5	2.5	0.1	0.5	0.5	0.01
2	5	5	10	2.5	2.5	5	0.5	0.5	0.5
3	5	5	10	2.5	2.5	9.9	0.5	0.5	0.99

Table 11: This table contains the indices $c_k\mu_k^s$ used by the $c\mu$ -rule for each of the three sets of parameter values used in the study of varying degrees of service capacity deterioration, along with the service rate ratios μ_k^1/μ_k^2 .

Figure 3 summarizes the observed performance of the five heuristic policies. As expected, the performance of the policies that eschew preventive maintenance (i.e., CMU and LWF) worsens as the loss in service capacity increases (i.e., as μ_3^1/μ_3^2 decreases). This is especially true for LWF, which becomes unstable when $\mu_3^1/\mu_3^2 = 0.01$. On the other hand, the performance of the other three policies (i.e., CMUPM, LWFP, and DDPM) is fairly insensitive to changes in μ_3^1/μ_3^2 , with CMUPM tending to outperform LWFP.

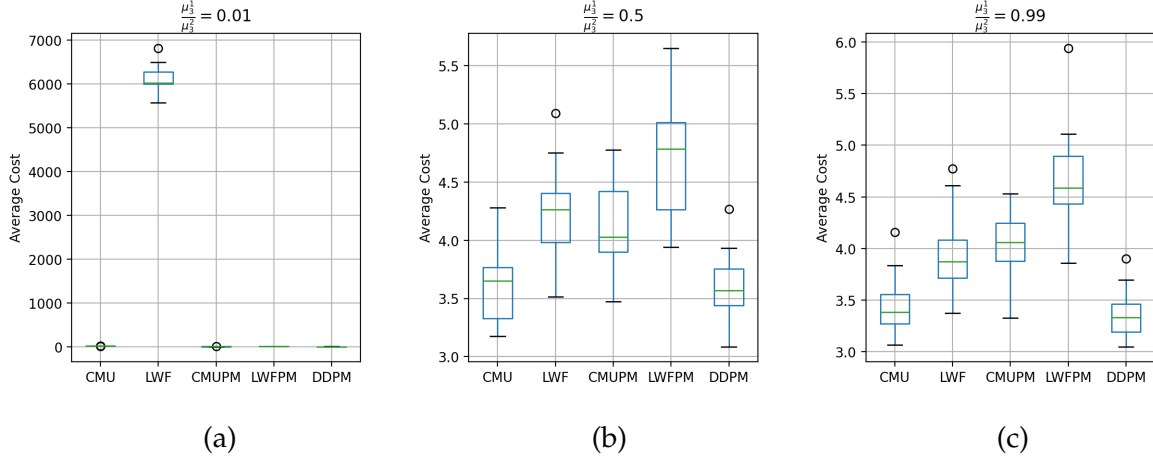


Figure 3: Policy Performance: Degree of Service Capacity Deterioration

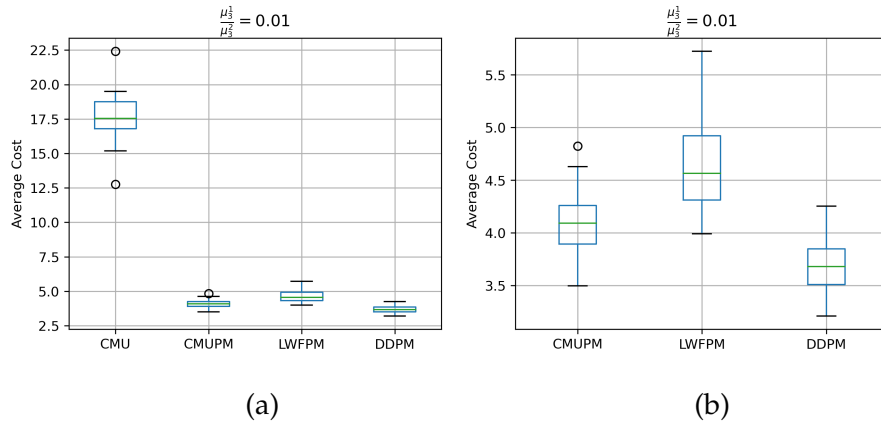


Figure 4: Policy Performance: Degree of Service Capacity Deterioration ($\frac{\mu_3^1}{\mu_3^2} = 0.01$ only)

5.3.4 Variability of Inter-Event Distributions

Recall that, for this set of parameter values, the inter-event times are gamma distributed. The priority ordering of job classes under the $c\mu$ -rule is class 2, then class 1, then class 3, regardless of the server state; see Table 12. The ratios between deteriorated and like-new service rates are also

given in Table 12; note that Assumption CR does not hold.

#	$c_1\mu_1^2$	$c_2\mu_2^2$	$c_3\mu_3^2$	$c_1\mu_1^1$	$c_2\mu_2^1$	$c_3\mu_3^1$	$\frac{\mu_1^1}{\mu_1^2}$	$\frac{\mu_2^1}{\mu_2^2}$	$\frac{\mu_3^1}{\mu_3^2}$
1 – 16	5	20	2.5	2.5	10	0.025	0.5	0.5	0.01

Table 12: This table contains the indices $c_k\mu_k^s$ used by the $c\mu$ -rule for each of the sixteen sets of parameter values used in the study of variability in the inter-event distributions, along with the service rate ratios μ_k^1/μ_k^2 .

Figures 5, 6, 7, and 8 summarize the observed performance of the five heuristic policies; in the following, CV stands for coefficient of variation. Each boxplot represents the observed performance of a particular policy over all of the parameter sets, with the CV of one of the inter-event time distributions held constant. For example, in Figure 5, the boxplot for DDPM under the plot titled “Inter-Arrival CV = 0.1” was generated using the simulation performance data for DDPM on those parameter sets for which the coefficient of variation of the interarrival times is equal to 0.1 (i.e., parameter sets 1, 3, 4, 5, 9, 10, 11, and 15).

The relative performance of the policies is insensitive to changes in variability. In all cases, longest-waiting-time-first scheduling leads to poor performance relative to the MDP-based policy (DDPM), especially when no preventive maintenance is performed. The poor performance without preventive maintenance is not surprising, in light of the significant amount of service rate deterioration for class 3 jobs, and applies to $c\mu$ -based scheduling without preventive maintenance as well. In contrast, $c\mu$ -based scheduling with preventive maintenance (CMUPM) matches the MDP-based policy (DDPM) in performance.

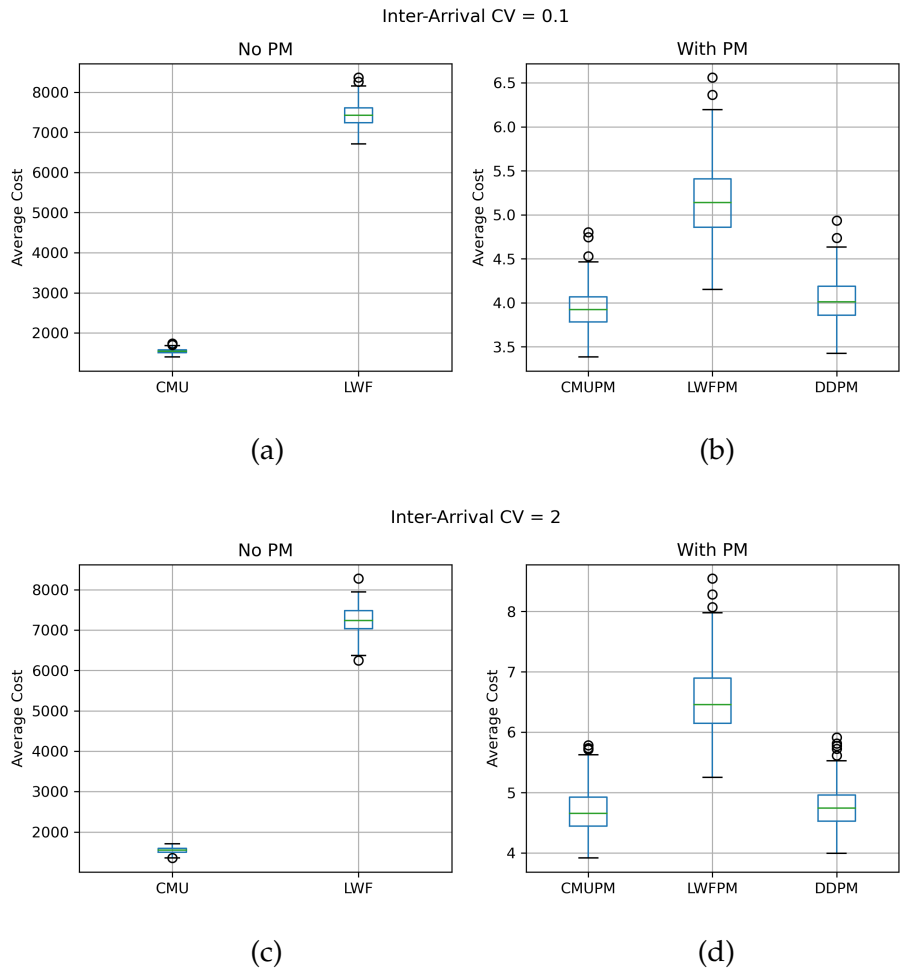


Figure 5: Policy Performance: Varying Inter-Arrival Time Coefficient of Variation (CV)

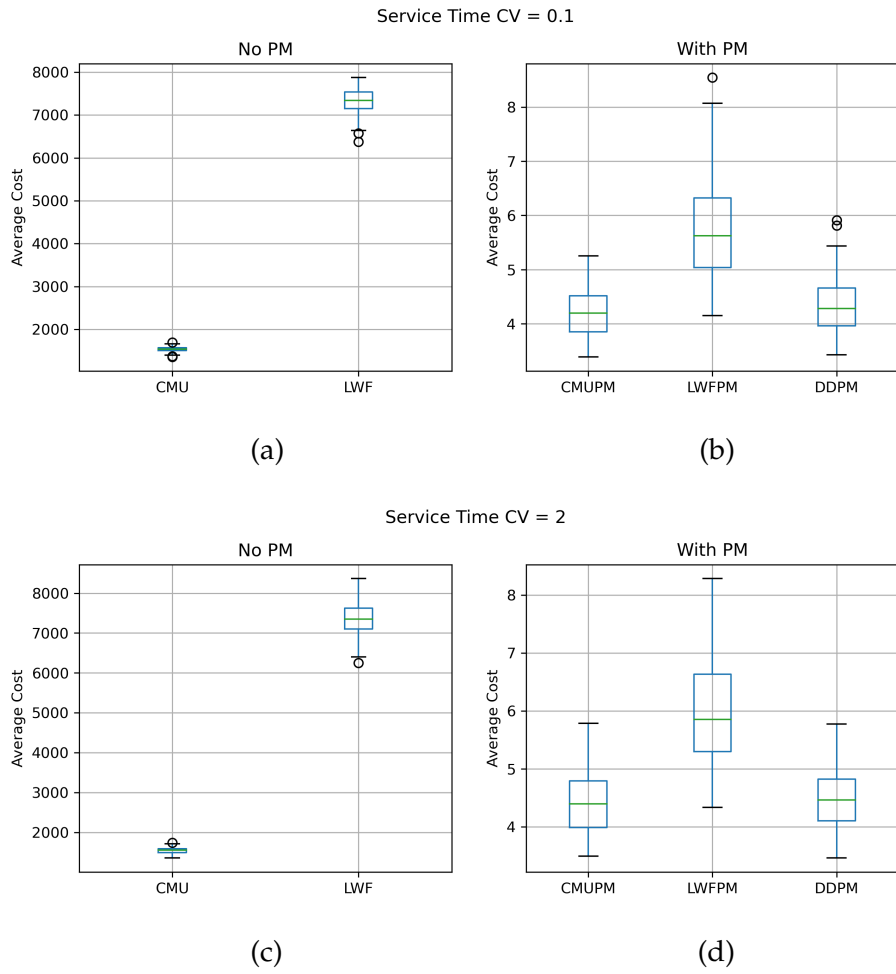


Figure 6: Policy Performance: Varying Service Time Coefficient of Variation (CV)

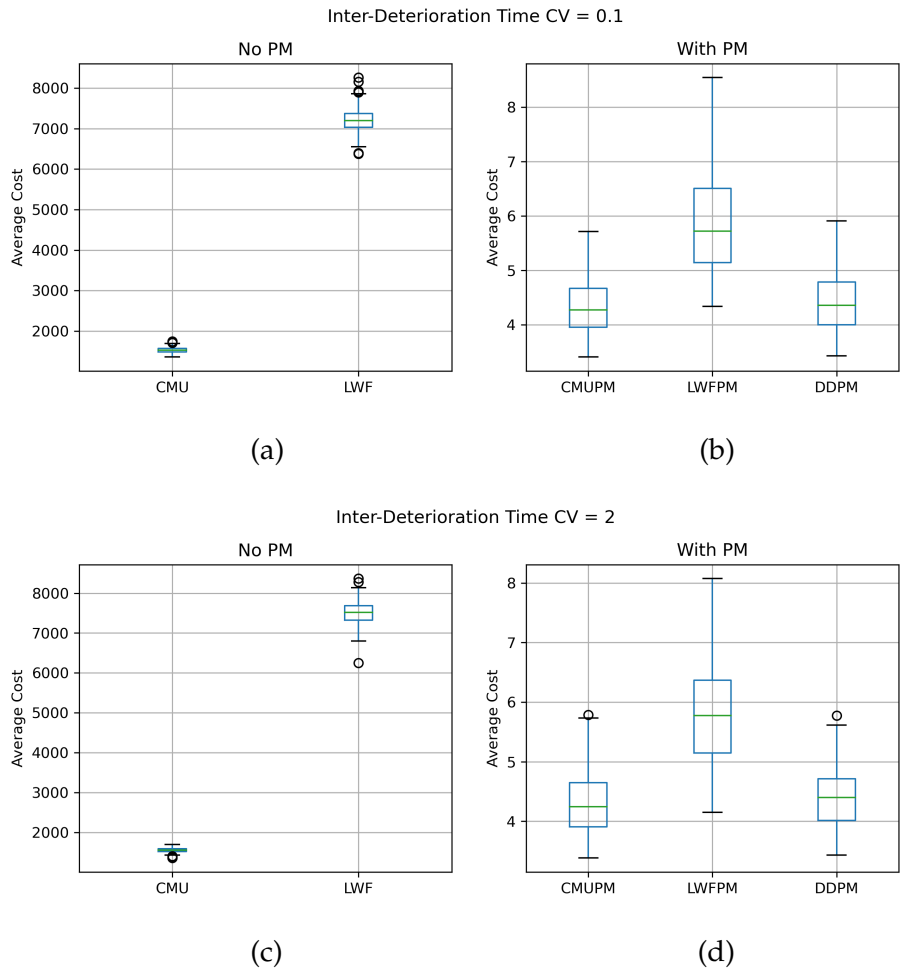


Figure 7: Policy Performance: Varying Inter-Deterioration Time Coefficient of Variation (CV)

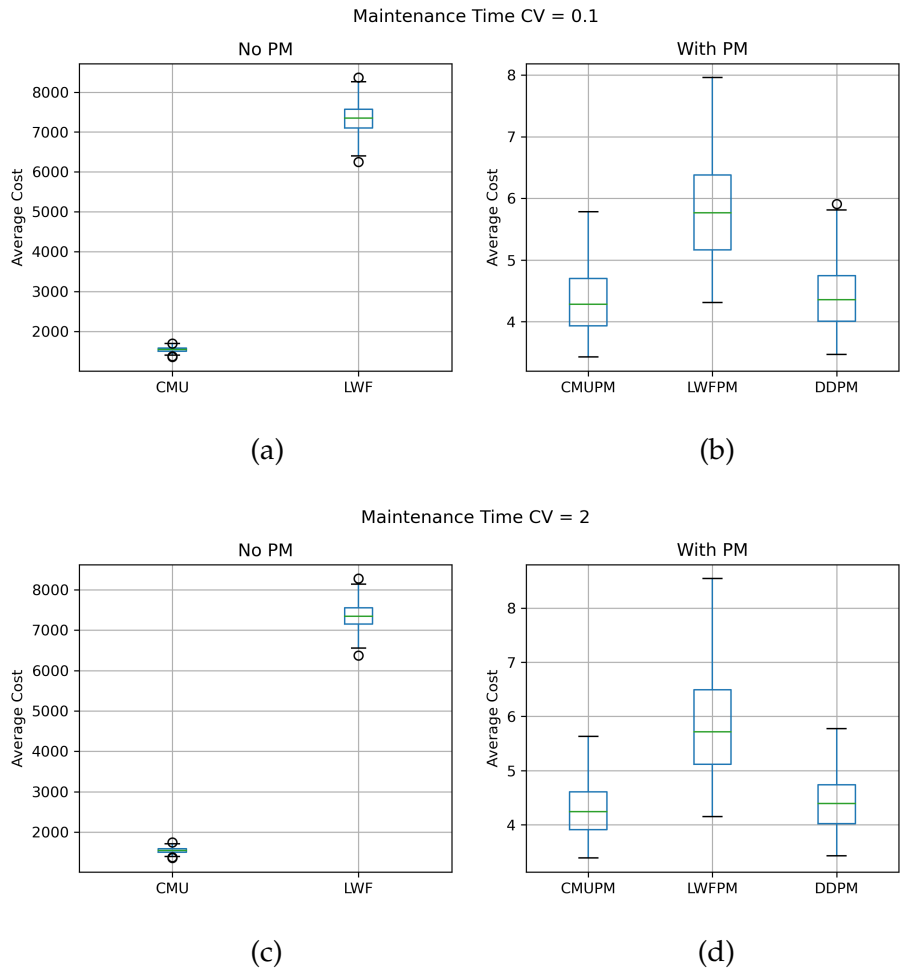


Figure 8: Policy Performance: Varying Maintenance Time Coefficient of Variation (CV)

6 Conclusion

In this work, we used a queueing control model to study the problem of how to jointly allocate work and perform preventive maintenance for a flexible server. We identified a condition (Assumption CR) under which it is optimal to schedule according to a state-dependent $c\mu$ -rule, as well as an average $c\mu$ -rule where the mean service rates are used, when preventive maintenance is not possible (Theorem 2). When Assumption CR does not hold, using these $c\mu$ -based schedul-

ing rules may result in an unstable system (Example 1), but our numerical results indicate that it is still possible for such scheduling rules to perform well without Assumption CR.

We then used Theorem 2 to show that, when the preventive maintenance policies are restricted to be age-based, calendar-based, or more generally independent of the queue lengths, it is without loss of optimality to use the aforementioned $c\mu$ -based scheduling rules (Theorem 5). In the context of semiconductor manufacturing, the implementation of condition-based maintenance is still very much on the cutting-edge of current research (see e.g., Djurdjanovic [13]), and that age/job-based preventive maintenance policies remain very relevant to practice (see e.g., Yao et al. [27]). Regarding the structure of preventive maintenance policies, we were able to prove that under assumptions analogous to those considered in Kaufman and Lewis [16] for one job class, the monotonicity property of optimal maintenance policies identified in [16, Theorem 3.2] is preserved when there are an arbitrary number of job types. In particular, there exists an optimal joint scheduling and maintenance policy where, for each fixed number of jobs of each class, the maintenance decisions are based on a threshold on the server state.

Finally, we presented the results of numerical experiments that compared the performance of $c\mu$ -based scheduling with a more naïve scheduling rule (longest-queue-first) and with scheduling based on solving an MDP (Section 5). We observed that, regardless of whether preventive maintenance was performed, the $c\mu$ -based scheduling rules are competitive with MDP-based scheduling. Moreover, the worse and more variable performance of longest-queue-first scheduling illustrated the value of incorporating service-rate information. On the other hand, the numerical experiments did not suggest an easy distinction between situations where scheduling has more of a performance impact than preventive maintenance, or vice versa. A better understanding of this distinction, as well as other research directions described in Section 6.1 below, is left for future work.

6.1 Future Work

Our work suggests a number of promising research directions.

Optimality Conditions for the $c\mu$ -Rule: Assumption CR, which only involves the service rates, does not depend on the deterioration dynamics of the server. For situations where Assumption CR is too strong, it would be worthwhile to identify conditions on the deterioration process under which $c\mu$ -based scheduling remains optimal. Moreover, it may be possible to relax the queue-obliviousness of maintenance policies in Theorem 4. Finally, it would be interesting to determine whether there are any guarantees on the optimality of $c\mu$ -based scheduling as a function of some measure of the degree to which Assumption CR is violated.

Good and Implementable Maintenance Policies: The focus of this paper has been on identifying conditions under which it suffices to follow a simple policy for the scheduling decisions. This of course leaves open the question of how maintenance policies should be derived. As was pointed out in Kaufman and Lewis [16], the optimal MDP-based policies can be very complicated. It would therefore be worthwhile to develop maintenance heuristics that both perform well across system parameters of interest, and that are easy to implement.

Relaxing Modeling Assumptions: In many applications, including some in semiconductor manufacturing [5, 7], the assumption that the server deteriorates independently of the work it performs is too strong. It would also be of interest to consider multiple servers and/or stations, or to assume that the server state is only partially observable.

Acknowledgments. The authors thank Chiao-Ju Sun and James Wu for their help with the numerical experiments. The work of the second author is supported by the Natural Sciences and Engineering Research Council of Canada.

References

- [1] S. Andradóttir, H. Ayhan, and D. G. Down. Compensating for failures with flexible servers. *Operations Research*, 55(4):753–768, 2007.

- [2] K. Backer, R. J. Huang, M. Lertchaitawee, M. Mancini, and C. Tan. Taking the next leap forward in semiconductor yield improvement. *McKinsey & Company*, May 2018.
- [3] R. E. Bohn and C. Terwiesch. The economics of yield-driven processes. *Journal of Operations Management*, 18(1):41–59, 1999.
- [4] P. Brémaud. *Point Processes and Queues: Martingale Dynamics*. Springer-Verlag New York, 1981.
- [5] Y. Cai, J. J. Hasenbein, E. Kutanoglu, and M. Liao. Single-machine multiple-recipe predictive maintenance. *Probability in the Engineering and Informational Sciences*, 27(2):209–235, 2013.
- [6] Y. Cai, E. Kutanoglu, J. Hasenbein, and J. Qin. Single-machine scheduling with advanced process control constraints. *Journal of Scheduling*, 15(2):165–179, 2012.
- [7] M. Celen and D. Djurdjanovic. Integrated maintenance decision-making and product sequencing in flexible manufacturing systems. *Journal of Manufacturing Science and Engineering*, 137(4):041006–041006–15, 2015.
- [8] M. E. Cholette, M. Celen, D. Djurdjanovic, and J. D. Rasberry. Condition monitoring and operational decision making in semiconductor manufacturing. *IEEE Transactions on Semiconductor Manufacturing*, 26(4):454–464, 2013.
- [9] Integrated Circuit Engineering Corporation. Yield and Yield Management. In *Cost Effective IC Manufacturing*. Scottsdale, AZ, 1997.
- [10] J. G. Dai. On positive Harris recurrence of multiclass queueing networks: A unified approach via fluid limit models. *The Annals of Applied Probability*, 5(1):49–77, 1995.
- [11] J. G. Dai. A fluid limit model criterion for instability of multiclass queueing networks. *The Annals of Applied Probability*, 6(3):751–757, 1996.

- [12] J. G. Dai and S. P. Meyn. Stability and convergence of moments for multiclass queueing networks via fluid limit models. *IEEE Transactions on Automatic Control*, 40(11):1889–1904, November 1995.
- [13] D. Djurdjanovic. Condition monitoring and operational decision-making in modern semiconductor manufacturing systems. In *Proceedings of the Pacific Rim Statistical Conference for Production Engineering*, ICSA Book Series in Statistics, pages 41–66. Springer, Singapore, 2018.
- [14] L. Grigoriu and D. Briskorn. Scheduling jobs and maintenance activities subject to job-dependent machine deteriorations. *Journal of Scheduling*, 20(2):183–197, 2017.
- [15] S. M. R. Iravani and I. Duenyas. Integrated maintenance and production control of a deteriorating production system. *IIE Transactions*, 34(5):423–435, 2002.
- [16] D. L. Kaufman and M. E. Lewis. Machine maintenance with workload considerations. *Naval Research Logistics*, 54(7):750–766, 2007.
- [17] J. J. McCall. Maintenance policies for stochastically failing equipment: A survey. *Management Science*, 11(5):493–524, 1965.
- [18] P. Nain. Interchange arguments for classical scheduling problems in queues. *Systems & Control Letters*, 12(2):177–184, 1989.
- [19] W. P. Pierskalla and J. A. Voelker. A survey of maintenance models: The control and surveillance of deteriorating systems. *Naval Research Logistics*, 23(3):353–388, 1976.
- [20] L. I. Sennott. Average cost semi-Markov decision processes and the control of queueing systems. *Probability in the Engineering and Informational Sciences*, 3(2):247–272, 1989.
- [21] Y. S. Sherif and M. L. Smith. Optimal maintenance models for systems subject to failure—A review. *Naval Research Logistics*, 28(1):47–74, 1981.

- [22] T. W. Sloan and J. G. Shanthikumar. Combined production and maintenance scheduling for a multiple-product, single-machine production system. *Production and Operations Management*, 9(4):379–399, 2000.
- [23] C. J. Spanos. Statistical process control in semiconductor manufacturing. *Microelectronic Engineering*, 10(3):271–276, 1991.
- [24] C. Valdez-Flores and R. M. Feldman. A survey of preventive maintenance models for stochastically deteriorating single-unit systems. *Naval Research Logistics*, 36(4):419–446, 1989.
- [25] C. Wu, D. G. Down, and M. E. Lewis. Heuristics for allocation of reconfigurable resources in a serial line with reliability considerations. *IIE Transactions*, 40(6):595–611, 2008.
- [26] C. Wu, M. E. Lewis, and M. Veatch. Dynamic allocation of reconfigurable resources in a two-stage tandem queueing system with reliability considerations. *IEEE Transactions on Automatic Control*, 51(2):309–314, 2006.
- [27] X. Yao, E. Fernández-Gaucherand, M. C. Fu, and S. I. Marcus. Optimal preventive maintenance scheduling in semiconductor manufacturing. *IEEE Transactions on Semiconductor Manufacturing*, 17(3):345–356, 2004.

A Appendix

A.1 Instability of Statically Prioritizing Class 1 in Example 1

Consider the (non-idling) policy that always prioritizes class 1 when the server is online. To show that this policy is *unstable*, in the sense that it incurs an infinite long-run expected average cost regardless of the initial state, we consider its associated fluid model.

Let $T_{k,s}(t)$ denote the total amount of time during $[0, t]$ that the server has spent serving class k jobs while it is in state s , and suppose $Q_1(0) = Q_2(0) = 0$. Arguments analogous to those in [11, p. 753] (replace k with k, s) imply that for every sequence $\{q_n, n \geq 0\}$ such that $q_n \rightarrow \infty$

there exists a subsequence $\{q_m, m \geq 0\}$ such that $\lim_{m \rightarrow \infty} T_{k,s}(q_m t)/q_m =: \bar{T}_{k,s}(t)$ exists for $k = 1, 2, s = 1, 2$, and $t \geq 0$. According to [11, Proposition 3.1], the associated scaled queue lengths $\bar{Q}_k(t) := \lim_{m \rightarrow \infty} Q_k(q_m t)/q_m, k = 1, 2$, satisfy

$$\bar{Q}_k(t) = \lambda_k t - \mu_k^1 \bar{T}_{k,1}(t) - \mu_k^2 \bar{T}_{k,2}(t), \quad k = 1, 2, t \geq 0, \quad (13)$$

where $\lambda_1 = 5, \lambda_2 = 0.8, \mu_1^1 = \mu_1^2 = 10, \mu_2^1 = 1, \text{ and } \mu_2^2 = 2$. In what follows, we will require derivatives of $\bar{T}_{k,s}(t)$ and $\bar{Q}_k(t)$, for $k = 1, 2$ and $s = 1, 2$. For $t \geq s \geq 0, \bar{T}_{k,s}(t) - \bar{T}_{k,s}(s) \leq t - s$, so $\bar{T}_{k,s}(t)$ is Lipschitz continuous. Hence, by (13), $\bar{Q}_k(t)$ is also Lipschitz continuous. As a result, the required derivatives exist almost everywhere.

Since class 1 is prioritized in states $s = 1, 2$, and the server is always online (corrective maintenance occurs instantaneously), the server is always busy at class 1 whenever class 1 jobs are present. As a result,

$$\bar{Q}_1(t) > 0 \implies \frac{d}{dt} \bar{T}_{1,1}(t) + \frac{d}{dt} \bar{T}_{1,2}(t) = 1 \implies \frac{d}{dt} \bar{Q}_1(t) = -5 < 0. \quad (14)$$

Note that for a nonnegative function $f(t)$, if $\frac{d}{dt} f(t) < 0$ whenever $f(t) > 0$, then if $f(t_0) = 0$ for some $t_0 \geq 0, f(t) = 0$ for all $t \geq t_0$. As $\bar{Q}_1(0) = 0$, it then follows from (14) that $\bar{Q}_1(t) = 0$ for all $t \geq 0$ which, according to (13), implies that $\frac{d}{dt} \bar{T}_{1,1}(t) + \frac{d}{dt} \bar{T}_{1,2}(t) = \frac{1}{2}$. But since the deterioration rates for states $s = 1, 2$ are equal and $q(1|2) = 1$, this means

$$\frac{d}{dt} \bar{T}_{1,1}(t) = \frac{d}{dt} \bar{T}_{1,2}(t) = \frac{1}{4}. \quad (15)$$

The equality of the deterioration rates and $q(1|2) = 1$ yield that the limiting proportion of times spent in $s = 1$ and $s = 2$ are equal. This and the fact that the server cannot be busy more than 100 percent of the time imply that

$$\frac{d}{dt} \bar{T}_{1,s}(t) + \frac{d}{dt} \bar{T}_{2,s}(t) \leq \frac{1}{2}, \quad s = 1, 2. \quad (16)$$

Combining (13), (15), and (16), we conclude that

$$\begin{aligned} \frac{d}{dt} \bar{Q}_2(t) &= 0.8 - (1) \frac{d}{dt} \bar{T}_{2,1}(t) - (2) \frac{d}{dt} \bar{T}_{2,2}(t) \\ &\geq 0.8 - (1) \left(\frac{1}{2} - \frac{d}{dt} \bar{T}_{1,1}(t) \right) - (2) \left(\frac{1}{2} - \frac{d}{dt} \bar{T}_{1,2}(t) \right) = 0.05 > 0. \end{aligned}$$

According to [11, Theorem 3.2], this implies that statically prioritizing class 1 is unstable.

A.2 Existence of a Stable Policy in Example 1

Consider the policy that prioritizes class s when the server state is s , for $s = 1, 2$. To show that this policy incurs a finite long-run expected average cost regardless of the initial state, by [12, Theorem 4.1] it suffices to show that its associated fluid model is stable in the sense that it drains and remains empty after a finite amount of time [10].

To define the fluid model, again consider the function $T_{k,s}(t)$ defined in Appendix A.1, and let $q = Q_1(0) + Q_2(0)$. Any limit point as $q \rightarrow \infty$ of the scaled process

$$\left(\frac{Q_1(qt)}{q}, \frac{Q_2(qt)}{q}, \frac{T_{1,1}(qt)}{q}, \frac{T_{2,1}(qt)}{q}, \frac{T_{1,2}(qt)}{q}, \frac{T_{2,2}(qt)}{q} \right)$$

is called a *fluid limit* of the original system. Every fluid limit

$$(\bar{Q}_1(t), \bar{Q}_2(t), \bar{T}_{1,1}(t), \bar{T}_{1,2}(t), \bar{T}_{2,1}(t), \bar{T}_{2,2}(t))$$

satisfies a set of differential equations known as the *fluid model*. For the system in Example 1 under the proposed policy, the fluid model is:

$$\frac{d}{dt} \bar{Q}_1(t) = \lambda_1 - \mu_1^1 \frac{d}{dt} \bar{T}_{1,1}(t) - \mu_1^2 \frac{d}{dt} \bar{T}_{1,2}(t), \quad (17)$$

$$\frac{d}{dt} \bar{Q}_2(t) = \lambda_2 - \mu_2^1 \frac{d}{dt} \bar{T}_{2,1}(t) - \mu_2^2 \frac{d}{dt} \bar{T}_{2,2}(t), \quad (18)$$

where $\lambda_1 = 5$, $\lambda_2 = 0.8$, $\mu_1^1 = \mu_1^2 = 10$, $\mu_2^1 = 1$, and $\mu_2^2 = 2$.

We now show that, under the proposed policy, every fluid limit is stable. In other words, for every fluid limit there exists a finite time $t_e \geq 0$ such that $\bar{Q}_1(t) = \bar{Q}_2(t) = 0$ for all $t \geq t_e$. First, recall that the deterioration rates are equal to 1, and that CM occurs instantaneously. This implies that the limiting proportions of time spent in states $s = 1$ and $s = 2$ are equal and the server is always online and busy if there are jobs present. Thus,

$$\bar{Q}_1(t) + \bar{Q}_2(t) > 0 \implies \frac{d}{dt} [\bar{T}_{1,s}(t) + \bar{T}_{2,s}(t)] = \frac{1}{2} \quad \forall s \in \{1, 2\}. \quad (19)$$

Recall that class 2 is prioritized when $s = 2$. Hence, according to (19),

$$\bar{Q}_2(t) > 0 \implies \frac{d}{dt} \bar{T}_{2,2}(t) = \frac{1}{2}. \quad (20)$$

Combining (18) with (20), and recalling that $\frac{d}{dt} \bar{T}_{2,1}(t) \geq 0$ for all t , we conclude that

$$\bar{Q}_2(t) > 0 \implies \frac{d}{dt} \bar{Q}_2(t) \leq \lambda_2 - 1 < 0, \quad (21)$$

since $\lambda_2 = 0.8 < 1$. So, as $\bar{Q}_2(t_e) = 0$ for $t_e = \bar{Q}_2(0)/(1 - \lambda_2)$, this with (21) yields

$$\bar{Q}_2(t) = 0, \quad t \geq t_e. \quad (22)$$

Next, we consider what happens to the fluid in queue 1 after queue 2 has drained. In general, since class 1 is prioritized when the server state $s = 1$, we know from (19) that $\frac{d}{dt} \bar{T}_{1,1}(t) = \frac{1}{2}$ whenever $\bar{Q}_1(t) > 0$. According to (17), this means

$$\bar{Q}_1(t) > 0 \implies \frac{d}{dt} \bar{Q}_1(t) = 5 - 10 \cdot \frac{d}{dt} \bar{T}_{1,2}(t). \quad (23)$$

On the other hand, suppose $t \geq t_e$. From (22), we know that $\frac{d}{dt} \bar{Q}_2(t) = 0$. Moreover, since class 1 is prioritized when $s = 1$, we also know that $\frac{d}{dt} \bar{T}_{2,1}(t) = 0$. In light of (18), these two observations imply that $\frac{d}{dt} \bar{T}_{2,2}(t) = \frac{\lambda_2}{2}$. According to (19) and the fact that $\lambda_2 < 1$, this means

$$\frac{d}{dt} \bar{T}_{1,2}(t) = \frac{1 - \lambda_2}{2} > 0. \quad (24)$$

We therefore conclude from (23) that

$$t \geq t_0 \text{ and } \bar{Q}_1(t) > 0 \implies \frac{d}{dt} \bar{Q}_1(t) < 0. \quad (25)$$

In summary, (22) and (25) imply that both queues drain and remain empty after a finite amount of time, i.e., that the fluid model is stable.

A.3 Proof of Theorem 6

In this section, we assume that Assumptions **M** and **S** hold. Under Assumption **M**, the joint scheduling and maintenance model described in Section 2 is a *semi-Markov decision process (SMDP)*; for background on SMDPs, see e.g., Sennott [20] and the references therein.

An SMDP is defined by the following objects:

1. the state set \mathbb{X} ,
2. sets of available actions $A(x)$ for each $x \in \mathbb{X}$,
3. transition probabilities $p(y|x, a)$ for each $x, y \in \mathbb{X}$ and $a \in A(x)$,
4. distributions $F(\cdot|x, a, y)$ for the time spent in each state $x \in \mathbb{X}$ given that action $a \in A(x)$ is taken and the next state of the process is $y \in \mathbb{X}$,
5. immediate costs $D(x, a)$ and cost rates $d(x, a)$ for each $x \in \mathbb{X}$ and $a \in A(x)$.

Recalling that we are only considering nonidling policies, for the joint scheduling and maintenance problem the above objects are defined as follows.

1. $\mathbb{X} = \{0, 1, \dots\}^K \times \{0, 1, \dots, S\}$;
2. letting $k = 0, 1, \dots, K$ denote idling ($k = 0$) or serving class k ($k > 0$), and letting PM and CM respectively denote initiating preventive and corrective maintenance, for $(\mathbf{i}, s) \in \mathbb{X}$ let

$$A(\mathbf{i}, s) = \begin{cases} \{\text{CM}\}, & \text{if } s = 0; \\ \{0, \text{PM}\}, & \text{if } \sum_{k=1}^K i_k = 0, s \geq 1; \\ \{\ell : i_\ell > 0\} \cup \{\text{PM}\}, & \text{if } \sum_{k=1}^K i_k > 0, s \geq 1; \end{cases} \quad (26)$$

3. for $(\mathbf{i}, s), y \in \mathbb{X}$ and $a \in A(\mathbf{i}, s)$, letting \mathbf{e}_k be the vector in \mathbb{R}^{K+1} where the k^{th} entry is a 1 and all others are zero, $\mathbf{n} = (n_1, \dots, n_K)$ and $\mu_0^s \equiv 0$, and recalling that by Assumption **M**(iii) the maintenance times are iid with distribution $G(\cdot)$,

$$p(y|(\mathbf{i}, s), a) = \begin{cases} \int_0^\infty \prod_{\ell=1}^K \frac{e^{-\lambda_\ell t} (\lambda_\ell t)^{n_\ell}}{n_\ell!} & \\ \quad \text{if } s = 0, a = \text{CM}, y = (\mathbf{i} + \mathbf{n}, S) & \\ \quad \text{or } s \geq 1, a = \text{PM}, y = (\mathbf{i} + \mathbf{n}, S); & \\ \frac{\lambda_k}{\sum_{\ell=1}^K \lambda_\ell + \mu_k^s + \alpha_s} & \text{if } s \geq 1, a = k, y = (\mathbf{i} + \mathbf{e}_k, s); \\ \frac{\alpha_s}{\sum_{\ell=1}^K \lambda_\ell + \mu_k^s + \alpha_s} & \text{if } s \geq 1, a = k, y = (\mathbf{i}, s - 1); \\ \frac{\mu_k^s}{\sum_{\ell=1}^K \lambda_\ell + \mu_k^s + \alpha_s} & \text{if } s \geq 1, a = k, y = (\mathbf{i}, s) - \mathbf{e}_k; \end{cases}$$

4. for $(\mathbf{i}, s), \mathbf{y} := (\mathbf{j}, u) \in \mathbb{X}$ and $\mathbf{a} \in A(\mathbf{x})$ (where $\mathbf{j} = (j_1, \dots, j_K)$),

$$F(t | (\mathbf{i}, s), \mathbf{a}, \mathbf{y}) = \begin{cases} G(t) & \text{if } s = 0, \mathbf{a} = \text{CM}, \mathbf{j} \geq \mathbf{i}, u = S \\ & \text{or } s \geq 1, \mathbf{a} = \text{PM}, \mathbf{j} \geq \mathbf{i}, u = S; \\ 1 - e^{-\lambda_k t} & \text{if } s \geq 1, \mathbf{a} \in \{0, 1, \dots, K\}, (\mathbf{j}, u) = (\mathbf{i} + \mathbf{e}_k, s); \\ 1 - e^{-\alpha_s t} & \text{if } s \geq 1, \mathbf{a} \in \{0, 1, \dots, K\}, (\mathbf{j}, u) = (\mathbf{i}, s - 1); \\ 1 - e^{-\mu_k^\dagger t} & \text{if } s \geq 1, \mathbf{a} = k, (\mathbf{j}, u) = (\mathbf{i} - \mathbf{e}_k, s); \end{cases}$$

where $\mathbf{j} \geq \mathbf{i}$ is interpreted componentwise.

5. for $(\mathbf{i}, s) \in \mathbb{X}$ and $\mathbf{a} \in A(\mathbf{i}, s)$,

$$D((\mathbf{i}, s), \mathbf{a}) = \begin{cases} K_c & \text{if } \mathbf{a} = \text{CM}; \\ K_p & \text{if } \mathbf{a} = \text{PM}; \\ 0 & \text{otherwise}; \end{cases}$$

and

$$d((\mathbf{i}, s), \mathbf{a}) = \sum_{k=1}^K c_k i_k.$$

It is useful to consider discounting the expected total cost incurred over an infinite horizon. In particular, given a *discount rate* $\beta > 0$, the expected β -discounted cost incurred from the initial state $(\mathbf{i}, s) \in \mathbb{X}$ under the policy $\pi \in \Pi$ is

$$v_\beta^\pi(\mathbf{i}, s) := \mathbb{E} \left[\sum_{n: t_n^\pi \leq t} e^{-\beta t_n^\pi} [K_c M_c^\pi(t_n^\pi) + K_p M_p^\pi(t_n^\pi)] + \int_0^\infty e^{-\beta t} \sum_{k=1}^K c_k Q_k^\pi(t) dt \mid \mathbf{Q}^\pi(0) = \mathbf{i}, S^\pi(0) = s \right].$$

Moreover, a policy π_* is β -optimal if $v_\beta^{\pi_*}(\mathbf{x}) = \inf_{\pi \in \Pi} v_\beta^\pi(\mathbf{x}) =: v_\beta(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{X}$.

Definition 9. A function $v : \mathbb{X} \rightarrow \mathbb{R}$ is monotone in the server's health if

$$\mathbf{i} \leq \mathbf{i}', s \geq s' \implies v(\mathbf{i}, s) \leq v(\mathbf{i}', s').$$

A straightforward adaptation of the sample-path argument in [16, Proof of Proposition 3.3] can be used to prove the following useful monotonicity property of v_β .

Proposition 10. *The value function v_β is monotone in the server's health.*

Lemma 11. *Assumptions **M** and **S** imply that the hypotheses of [20, Theorem 2, Proposition 4] hold.*

Proof. The hypotheses of [20, Theorem 2] consist of [20, Assumptions 1-5].

1. For $t \geq 0$, $x, y \in \mathbb{X}$, and $a \in A(x)$, let

$$H(t|x, a) := \sum_{y \in \mathbb{X}} p(y|x, a) F(t|x, a, y).$$

The first assumption states that there exist $\epsilon, \delta > 0$ such that

$$1 - H(\delta|x, a) \geq \epsilon \quad \forall x \in \mathbb{X}, a \in A(x). \quad (27)$$

First, recall that according to Assumption **M**(iii), $1/\alpha_0 = \int_0^\infty t dG(t) > 0$. This implies that there exists a $\delta^* > 0$ such that $1 - G(\delta^*) > 0$. Moreover, letting

$$\bar{\gamma} := \max\{\lambda_1, \dots, \lambda_K, \alpha_1, \dots, \alpha_B, \mu_1^1, \dots, \mu_1^S, \dots, \mu_K^1, \dots, \mu_K^S\} > 0$$

and

$$\epsilon^* := \min\{1 - G(\delta^*), e^{-\bar{\gamma}\delta^*}\} > 0,$$

it follows that (27) holds with $\epsilon = \epsilon^*$ and $\delta = \delta^*$.

2. For $x \in \mathbb{X}$ and $a \in A(x)$, let

$$\tau(x, a) := \sum_{y \in \mathbb{X}} p(y|x, a) \int_0^\infty t dF(t|x, a, y). \quad (28)$$

The second assumption states that there exists a constant $B < \infty$ such that

$$\tau(x, a) \leq B \quad \forall x \in \mathbb{X}, a \in A(x). \quad (29)$$

Letting

$$\underline{\gamma} := \min\{\lambda_1, \dots, \lambda_K, \alpha_1, \dots, \alpha_B, \mu_1^1, \dots, \mu_1^S, \dots, \mu_K^1, \dots, \mu_K^S\} > 0$$

and

$$B^* := \max\{1/\alpha_0, 1/\underline{\gamma}\} < \infty,$$

it follows that (29) holds with $B = B^*$.

3. The third assumption states that

$$v_\beta(x) < \infty \quad \forall \beta > 0, x \in \mathbb{X}. \quad (30)$$

According to [20, Remark 1], a sufficient condition for (30) to hold is the existence of a policy π such that

$$w^\pi(x) < \infty \quad \forall x \in \mathbb{X}. \quad (31)$$

Let s^* be a state that satisfies Assumption S(i). By analyzing a fluid model analogous to the one in [16, Proof of Proposition 3.1], it can be shown that (31) is satisfied by any policy that initiates PM whenever the server state is less than s^* , and otherwise does not idle an online server if the system is nonempty. Hence (30) holds.

4. Let $\mathbf{0} := (0, \dots, 0, S)$, and

$$h_\beta(x) := v_\beta(x) - v_\beta(\mathbf{0}), \quad x \in \mathbb{X}.$$

The fourth assumption states that there exists a $\beta_0 > 0$ and $M : \mathbb{X} \rightarrow [0, \infty)$ such that

$$h_\beta(x) \leq M(x) \quad \forall \beta \in (0, \beta_0), x \in \mathbb{X} \quad (32)$$

and

$$\exists a(x) \in A(x) \text{ such that } \sum_{y \in \mathbb{X}} p(y|x, a(x))M(y) < \infty \quad \forall x \in \mathbb{X}. \quad (33)$$

Let $X^\pi(t) := (\mathbf{Q}^\pi(t), S^\pi(t))$ denote the state of the system at time t under the policy π . For $z \in \mathbb{X}$, let $\tau_z^\pi := \inf\{t > 0 \mid X^\pi(t) = z\}$ and, for $x, y \in \mathbb{X}$, let

$$C^\pi(x, y) := \mathbb{E} \left[\sum_{n: t_n^\pi \leq \tau_y^\pi} [K_c M_c^\pi(t_n^\pi) + K_p M_p^\pi(t_n^\pi)] + \int_0^{\tau_y^\pi} \sum_{k=1}^K c_k Q_k^\pi(t) dt \mid X^\pi(0) = x \right].$$

denote the expected total cost incurred up to a first passage from x to y under the policy π . According to [20, Remark 1], a sufficient condition for (32) and (33) to hold for some $\beta_0 > 0$, $M : \mathbb{X} \rightarrow [0, \infty)$ is the existence of a stationary policy φ such that

$$C^\varphi(x, \mathbf{0}) < \infty \quad \forall x \in \mathbb{X}. \quad (34)$$

Let s^* be a state that satisfies Assumption **S(i)**, and let φ_{s^*} be any stationary policy that initiates PM whenever the server state is less than s^* , and otherwise does not idle an online server if the system is nonempty. By analyzing a fluid model analogous to the one in [16, Proof of Proposition 3.1], it can be shown that φ_{s^*} satisfies (31), and that the embedded state process under φ_{s^*} is a unichain Markov chain where the set of states $\{(\mathbf{i}, s) \mid s^* \leq s \leq S\}$ is the ergodic class and the remaining states are transient. By [20, Lemma 2], it follows that (34) holds with $\varphi = \varphi_{s^*}$. Hence there exist $\beta_0 > 0$ and $M : \mathbb{X} \rightarrow [0, \infty)$ such that (32) and (33) hold.

5. The fifth assumption states that there exist a $\beta_0 > 0$ and $N \geq 0$ such that

$$-N \leq h_\beta(x) \quad \forall \beta \in (0, \beta_0), x \in \mathbb{X}. \quad (35)$$

Since $h_\beta(x) = v_\beta(x) - v_\beta(\mathbf{0})$ for $x \in \mathbb{X}$, and $\mathbf{0} = (0, \dots, 0, S)$, it follows from Proposition 10 that (35) holds with $N = 0$ and any $\beta_0 > 0$.

Next, the hypotheses of [20, Proposition 4] consist of [20, Assumptions 1-5] and the following assumption: there exist $\epsilon > 0$ and a finite set $G \subset \mathbb{X}$ such that

$$\min_{\mathbf{a} \in A(x)} d(x, \mathbf{a}) \geq \frac{B(g + \epsilon)}{\inf_{x, \mathbf{a}} \tau(x, \mathbf{a})} \quad \forall x \in \mathbb{X} \setminus G \quad (36)$$

where g is a constant from [20, Theorem 2], and $\inf_{x, \mathbf{a}} \tau(x, \mathbf{a}) > 0$ by [20, Lemma 1]. Recalling that

$$d((\mathbf{i}, s), \mathbf{a}) = \sum_{k=1}^K c_k i_k, \quad (\mathbf{i}, s) \in \mathbb{X}, \mathbf{a} \in A(\mathbf{i}, s),$$

consider any $\epsilon^* > 0$ and let

$$G^* := \left\{ (i_1, 0, \dots, 0, s) \in \mathbb{X} \mid i_1 < \left\lfloor \frac{B(g + \epsilon^*)}{c_1 \inf_{x, \mathbf{a}} \tau(x, \mathbf{a})} \right\rfloor \right\}.$$

Then $|G^*| < \infty$, and (36) holds with $\epsilon^* = \epsilon$ and $G^* = G$. □

For $f : \mathbb{X} \rightarrow \mathbb{R}$, $x \in \mathbb{X}$, and $a \in A(x)$, let

$$\begin{aligned} T_{\beta}^a f(x) &:= D(x, a) + d(x, a) \sum_{y \in \mathbb{X}} p(y|x, a) \int_0^{\infty} \int_0^t e^{-\beta u} du dF(t|x, a, y) \\ &\quad + \sum_{y \in \mathbb{X}} p(y|x, a) \int_0^{\infty} e^{-\beta t} dF(t|x, a, y) f(y) \end{aligned}$$

Theorem 12. *Under Assumptions **M** and **S**, the following statements hold.*

(i) *The value function v_{β} satisfies the discounted-cost optimality equation (DCOE)*

$$v_{\beta}(x) = \min_{a \in A(x)} T_{\beta}^a v_{\beta}(x) \quad \forall x \in \mathbb{X}. \quad (37)$$

(ii) *For every $\beta > 0$ there exists a β -optimal deterministic stationary policy π_{β} .*

(iii) *A deterministic stationary policy π is β -optimal if and only if*

$$\pi(x) \in \arg \min_{a \in A(x)} T_{\beta}^a v_{\beta}(x) \quad \forall x \in \mathbb{X}.$$

(iv) *Every β -optimal deterministic stationary policy is monotone in the server's health.*

Proof. According to Lemma 11, [20, Assumptions 1,3] hold, which implies that statements (i)-(iii) hold by [20, Theorem 1].

Next, suppose that it is not β -optimal to perform PM in state (\mathbf{i}, s) . Then by statement (iii),

$$T_{\beta}^{\text{PM}} v_{\beta}(\mathbf{i}, s) > v_{\beta}(\mathbf{i}, s).$$

Since PM incurs the same fixed cost whenever it is initiated, and the subsequent maintenance times are iid, it follows from Proposition 10 that

$$T_{\beta}^{\text{PM}} v_{\beta}(\mathbf{i}, s + 1) = T_{\beta}^{\text{PM}} v_{\beta}(\mathbf{i}, s) > v_{\beta}(\mathbf{i}, s) \geq v_{\beta}(\mathbf{i}, s + 1).$$

By statement (iii), this implies that it is also not β -optimal to perform PM in state $(\mathbf{i}, s + 1)$. Hence statement (iv) holds. \square

Proof of Theorem 6. Lemma 11 implies that [20, Theorem 2, Proposition 4] hold for the SMDP formulated in this section. In particular, [20, Theorem 2, Proposition 4] state that there exists a deterministic stationary optimal policy π_* that is a limit point of a sequence of β -optimal deterministic stationary policies. Since the action sets are finite, it follows that π_* is actually β -optimal for some $\beta > 0$. According to Theorem 12(iv), π_* is monotone in the server's health. \square

A.4 Linear Program for Arrival Rate Scaling

In this section, we formulate the linear program used in Section 5 to scale the selected arrival rate levels $\tilde{\lambda}_k$ to actual arrival rates λ_k to ensure the existence of a stable policy.

Data

- K = number of job classes
- S = number of server states
- $\tilde{\lambda}_k$ = arrival rate level for class k jobs; see Section 5.1.1
- μ_k^s = service rate for class k jobs, when the server state is s
- α_k^s = deterioration rate when the server is in state s and working on a class k job
- α_0 = maintenance rate

Decision Variables

- x_k^s = long-run fraction of time that the server is in state s and serving class k jobs
- x^0 = long-run fraction of time that the server is in state 0
- y^s = long-run average rate at which maintenance is initiated when the server is in state s

Optimization Problem

$$\begin{aligned}
& \text{maximize} && m \\
& \text{subject to} && \sum_{s=1}^S \mu_k^s x_k^s \geq (1+m)\bar{\lambda}_k && k = 1, \dots, K \\
& && \sum_{s=1}^S x_k^s \leq 1 && k = 1, \dots, K \\
& && \sum_{k=1}^K \alpha_k^1 x_k^1 + \sum_{s=1}^S y^s = \alpha_0 x^0 \\
& && \sum_{k=1}^K \alpha_k^S x_k^S = \alpha_0 x^0 \\
& && \sum_{k=1}^K \alpha_k^{s+1} x_k^{s+1} = \sum_{k=1}^K \alpha_k^s x_k^s + y^s && s = 1, \dots, S-1 \\
& && y^s \leq \sum_{k=1}^K \alpha_k^{s+1} x_k^{s+1} && s = 1, \dots, S-1 \\
& && \sum_{k=1}^K \sum_{s=1}^S x_k^s + x^0 = 1 \\
& && 0 \leq x_k^s \leq 1 && k = 1, \dots, K, \quad s = 1, \dots, S-1 \\
& && 0 \leq x^0 \leq 1 \\
& && y^s \geq 0 && s = 1, \dots, S-1
\end{aligned}$$

The first constraint ensures that the average rate at which class i jobs are served is at least equal to the inflated arrival rate for that class i . The second constraint ensures that the proportion of time that the server is serving class i jobs is at most 1. (Redundant with the seventh constraint and non-negativity of the x_i^k 's.) The third constraint is the rate balance constraint for server state 0. It states that the average rate at which the server deteriorates while in state 1 and serving class i jobs, plus the total rate at which maintenance is initiated, equals the average rate at which maintenance/repair is completed. The fourth constraint is the rate balance constraint for server state K . It states that the rate at which the server deteriorates while in state K and serving class

i jobs equals the rate at which maintenance/repair is completed. The fifth constraint is the rate balance constraint for server state k . It states that the average rate at which the server deteriorates while in state $k + 1$ and serving class i jobs equals the rate at which the server deteriorates while in state k and serving class i jobs plus the rate at which maintenance is initiated from state k . The sixth constraint states that the average rate at which maintenance is initiated while the server is in state k is at most the average rate at which the server deteriorates from state $k + 1$ while serving class i jobs. The seventh constraint states that the proportion of time that the server is doing something (i.e., serving or being maintained) equals 1.

Given a solution to the linear program, and a target utilization ρ , the actual arrival rate λ_k for class k jobs is given by ρ times the average service rate for class k jobs:

$$\lambda_k = \rho \sum_{s=1}^S \mu_k^s x_k^s, \quad \text{for } k = 1, \dots, K.$$