

## ORIGINAL ARTICLE

# Optimizing resource allocation in service systems via simulation: A Bayesian formulation

Weiwei Chen<sup>1</sup>  | Siyang Gao<sup>2</sup> | Wenjie Chen<sup>3</sup> | Jianzhong Du<sup>4</sup>

<sup>1</sup>Department of Supply Chain Management, Rutgers University, Piscataway, New Jersey, USA

<sup>2</sup>Department of Advanced Design and Systems Engineering, City University of Hong Kong, Kowloon, Hong Kong

<sup>3</sup>Academy for Advanced Interdisciplinary Studies and Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen, China

<sup>4</sup>School of Management, Fudan University, Shanghai, China

## Correspondence

Weiwei Chen, Department of Supply Chain Management, Rutgers University, Piscataway, NJ 08854, USA.

Email: [wchen@business.rutgers.edu](mailto:wchen@business.rutgers.edu)

**Handling Editor:** Michael Pinedo

## Funding information

National Natural Science Foundation of China, Grant/Award Number: 72091211; City University of Hong Kong, Grant/Award Numbers: 7005269, 7005568

## Abstract

The service sector has become increasingly important in today's economy. To meet the rising expectation of high-quality services, efficiently allocating resources is vital for service systems to balance service qualities with costs. In particular, this paper focuses on a class of resource allocation problems where the service-level objective and constraints are in the form of probabilistic measures. Further, process complexity and system dynamics in service systems often render their performance evaluation and optimization challenging and relying on simulation models. To this end, we propose a generalized resource allocation model with probabilistic measures, and subsequently, develop an optimal computing budget allocation (OCBA) formulation to select the optimal solution subject to random noises in simulation. The OCBA formulation minimizes the expected opportunity cost that penalizes based on the quality of the selected solution. Further, the formulation takes a Bayesian approach to consider the prior knowledge and potential performance correlations on candidate solutions. Then, the asymptotic optimality conditions of the formulation are derived, and an iterative algorithm is developed accordingly. Numerical experiments and a case study inspired by a real-world problem in a hospital emergency department demonstrate the effectiveness of the proposed algorithm for solving the resource allocation problem via simulation.

## KEYWORDS

Bayesian model, optimal computing budget allocation, ranking and selection, resource allocation, service systems

## 1 | INTRODUCTION

The past few decades have witnessed a rapid growth of the service economy globally (Buckley & Majumdar, 2018). Nowadays, the service sector plays a vital role in a wide range of industries, including retail, healthcare, hospitality, finance, and IT. Traditional industries such as manufacturing have also seen the trend of servitization (Örşdemir et al., 2019) by introducing more service components in product offerings, requiring manufacturing processes to be more dynamic and flexible (Dmitry et al., 2021). Such service-oriented systems are often tailored to meet different customer requirements and expectations, and are complicated in process dependen-

cies and dynamics. Further, many activities in service systems deal directly with humans rather than machines, thereby introducing various forms of variability and challenges in measuring system performance. This paper focuses on service systems without tractable mathematical structures in view of the operational complexity, process nonlinearity, and performance variability, thus relying on simulation models in system evaluation and optimization.

A notable challenge of the service systems under study is to meet rising expectations on services without overspending on relevant resources. On the one hand, service providers are required or incentivized to meet and improve service quality as measured by certain metrics, such as waiting time (WT) in call centers and hospitals, on-time ratio in online food ordering and delivery platforms, and availability percentage in IT

Accepted by Michael Pinedo, after three revisions.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2022 The Authors. *Production and Operations Management* published by Wiley Periodicals LLC on behalf of Production and Operations Management Society.

services. These service goals are expressly provisioned in a service contract or perceived as crucial to customer satisfaction and market share. On the other hand, service providers work with a given set of resources, such as personnel with specific skills or training, space and equipment availability, and an operating budget. With the limitation on resources, they are motivated to improve service by intelligently allocating resources in place (e.g., functional areas) and time (e.g., shifts). For example, a hospital can improve care delivery by managing bed assignments and nurse schedules (Best et al., 2015; Kim & Mehrotra, 2015); a retail store can optimize product variety and quantity by properly utilizing storage and shelf spaces (Ton & Raman, 2010). This paper studies a class of service problems where one service level is optimized, and several others are imposed as constraints, subject to given resource limitations.

A case in point is a real-world resource allocation problem in a hospital emergency department (ED) (see Section 3.1 for the formulation and Section 5.2 for the case study). Common resources in an ED include specialized medical staff, beds, medical equipment, and salary budgets. The ED has established a set of achievable service goals as measured by WTs for different categories of patients. While the target service levels on urgent patients are high and satisfiable, the ED has recognized the need to improve the service quality on less-urgent patients. Therefore, it is desirable to seek resource allocation scenarios that improve the “nice-to-have” service goal while continuing to meet all other “must-have” requirements. Furthermore, the imposed service levels on WTs are expressed in the form of probabilistic measures; that is, a service level is defined as the percentage of customers not waiting longer than a given threshold.

Such probabilistic constraints (or chance constraints) are often used in measuring quality of service (Hong et al., 2015), and thus problems with similar structures appear in other service systems. Call centers may operate under a delay-percentile contract, where delays for a certain percentage of contract customers should be within a fixed delay bound, and noncontract customers have no guarantees on delay (Milner & Olsen, 2008). A call center may seek an operator schedule to minimize the delay for noncontract customers, while satisfying the requirement for contract customers. Such problems are also prevailing in public service offices, such as governmental departments that issue passports, driver licenses, and ID cards. They need to determine the staffing level for a specific time slot, such that at least a certain percentage of customers would wait less than a given time limit (Taigel et al., 2018). In this paper, we will use the term “resource allocation” for the aforesaid staffing schedules, equipment assignments, and other allocation of resources in place and time.

To study these service systems with probabilistic measures, simulation models can be deployed to evaluate the performance of different resource allocation scenarios so as to identify the best one. However, one notable drawback of performing analysis using simulation is the existence of random noises and thus the need to run a large number of

simulation replications in order to evaluate each scenario accurately. With sizable alternatives, it will be very time consuming to run sufficient replications for every scenario, and thus an efficient strategy is necessary to search for the optimal one subject to the computing budget (i.e., the total number of simulation replications). This line of research falls under the domain of ranking-and-selection (R&S). To this end, we will develop an optimal computing budget allocation (OCBA) formulation using the expected opportunity cost (EOC). Further, the proposed formulation will follow the Bayesian framework and impose a correlated prior belief, as treated in Frazier et al. (2009), on the performance of resource allocation scenarios. As such, it utilizes prior knowledge accumulated by the service provider on these scenarios and captures the performance correlation on “neighboring” scenarios.

The remainder of this paper is organized as follows. Section 2 summarizes related literature on applications of stochastic service systems and methodologies of simulation optimization. Section 3 introduces the resource allocation model and the OCBA formulation, with the solution methodology developed in Section 4. Section 5 presents the numerical results, and Section 6 concludes this paper. Proofs of lemmas and theorems, as well as additional analysis and numerical results, are enclosed in the Supporting Information.

## 2 | LITERATURE REVIEW

In this section, we will review relevant literature from both application and methodology viewpoints. As the problem under study considers probabilistic measures, we limit the review to those dealing with stochastic systems.

From the application point of view, most modeling papers studying stochastic resource allocation in service systems focus on problems with tractable mathematical structures. These well-structured problems can be modeled and solved as stochastic optimization problems, such as stochastic programming (Bodur & Luedtke, 2017), robust optimization (Mattia et al., 2017), and Markov decision process (Huh et al., 2013). These problems typically have one or more stochastic constraints and/or objective with others being deterministic. While the stochastic constraints/objective are often in the form of expectations (see, e.g., Bodur & Luedtke, 2017; Gans et al., 2015), some problems study probabilistic constraints as those formulated in this paper. For example, Beraldi et al. (2004) developed a stochastic programming framework to determine the emergency medical service site locations and the number of emergency vehicles, where the service-level constraints are probabilistic constraints; Robbins and Harrison (2010) studied a call center scheduling problem subject to a service constraint that a proportion of calls must be answered within a fixed time.

For stochastic service systems that do not possess tractable structures, simulation has been the tool of choice for analysis. Simulation can be used to evaluate the stochastic feasibility requirements, such as passenger processing targets at an

airport (Mason et al., 1998), and average WT requirements for patients in EDs (Izady & Worthington, 2012). Simulation can also be combined into optimization procedures to search for the optimal solution, known as simulation(-based) optimization, or optimization via simulation. For example, Atlason et al. (2008) and Cezik and L'Ecuyer (2008) solved call center staffing problems with non-closed-form service-level functions evaluated via simulation; Tsai and Zheng (2013) addressed a two-echelon inventory problem, where one constraint requires the expected response time at each depot to be within a threshold; Guo et al. (2017) and Chen et al. (2020) studied staffing problems in hospital EDs while satisfying stochastic service quality requirements. This paper contributes to the above literature on studying a generalized resource allocation problem in service systems, where two unique structures exist: (1) The performance of neighboring resource allocation scenarios may be correlated and the modeler can obtain prior knowledge on their joint distribution; and (2) the objective and multiple constraints are stochastic and in the form of probabilistic measures.

To tackle these special structures, this paper contributes to the R&S literature from the methodology point of view, particularly to those constrained R&S methods for solving problems with stochastic constraints and/or objective. One stream of research on constrained R&S is based on the sequential indifference-zone (IZ) framework and aims to identify the optimal solution with a probability exceeding a stipulated threshold (Andradóttir & Kim, 2010; Batur & Kim, 2010; Healey et al., 2014; Hong et al., 2015). Another type of formulation is based on the OCBA framework and aims to maximize the probability of correct selection (PCS) (Lee et al., 2012). A heuristic procedure was developed in Choi et al. (2021) by considering the precision of the sample means when maximizing PCS, so as to improve the algorithm performance when large simulation noise exists. Finally, large-deviation approaches have also been investigated, where Hunter and Pasupathy (2013) and Pasupathy et al. (2015) seek to maximize the rate of decay of the probability of false selection (PFS). Gao et al. (2019) further incorporated quadratic regression metamodels in the large-deviation framework to improve the search efficiency.

While these existing methods optimize the PCS or PFS measure, this paper develops a new Bayesian OCBA formulation using the EOC measure, motivated by the aforesaid structures of the resource allocation problem under study. First, the OCBA formulations with the PCS or PFS measure aim to identify the best solution with the maximum probability, treating all nonbest solutions equally. However, for the service systems under study, the industry is more interested in the economic value of the selection than its statistical significance. In other words, the decision makers are more concerned with the loss incurred in selecting a solution, rather than the statistical significance of selecting exactly the best one. As such, the OCBA formulation proposed in this paper uses the EOC measure as the objective, which penalizes for the quality of the selection. While EOC has been used as an objective in unconstrained R&S problems (Chick

& Wu, 2005; Gao et al., 2017), it has rarely been studied in the constrained R&S counterparts. Second, in order to capture the prior knowledge on the solution performances and the correlation of neighboring solutions, the proposed formulation uses the Bayesian framework. While Pujowidianto et al. (2013) proposed a frequentist version of EOC for a constrained R&S problem, to the best of our knowledge, this is the first Bayesian EOC formulation developed in the constrained R&S literature that accounts for the prior belief and correlation of solutions. Third, the probabilistic measures in the objective and constraints of the problem provide convenience in the formulation and its theoretical properties, including a strong finite-time convergence property of the proposed algorithm. We mention that although the probabilistic measures can be written as expectations of an indicator function following a Bernoulli distribution, the information on the Bernoulli distribution makes our method more efficient and robust, as pointed out in Hong et al. (2015).

Finally, this paper is more tangentially related to the literature on stochastically constrained optimization via simulation (COvS) and stochastically constrained Bayesian optimization (CBO). Both COvS and CBO are constrained optimization problems where samples from the objective and constraint measures can be treated as being generated from a black box. This black box is the simulation model for COvS and is the real system for CBO. In COvS and CBO, the set of candidate solutions is much larger, and search-based heuristic methods need to be developed to find the optimal or near-optimal solutions. Such heuristics can be based on penalty functions (Park & Kim, 2015), gradients (Luo & Lim, 2013; Nagaraj & Pasupathy, 2013), random search strategies (Chen et al., 2020; Gao & Chen, 2016), or acquisition functions (Gardner et al., 2014; Letham et al., 2019; Ungredda & Branke, 2021). A valuable future research is to extend the methodology developed in this paper to solve the COvS and CBO problems.

### 3 | PROBLEM DESCRIPTION

In this section, we first introduce a resource allocation problem in a hospital, and then generalize the model to other service systems. Then, a Bayesian OCBA model is formulated to efficiently identify the optimal solution of the model via simulation.

#### 3.1 | A resource allocation model with probabilistic measures

A motivating example of our model is a resource allocation problem in the ED of a large public hospital in Hong Kong. This hospital had large patient flows with an average of 350 to 430 patients per day to the ED. Sixty-five percent of the patients were walk-in patients, while the others were brought in by ambulance. All patients were labeled by triage categories based on their clinical conditions, including categories

**TABLE 1** Time limits and service levels for an ED in Hong Kong

Patient category	Expected time limit	Desired service level
I	Immediately	Whenever possible
II	WT ≤ 15 min	≥ 95%
III	WT ≤ 30 min	≥ 90%
IV	WT ≤ 2 h	≥ 75%
I–V	LOS ≤ 4 h	≥ 98%

I (critical), II (emergency), III (urgent), IV (semiurgent), and V (nonurgent), of which categories III and IV patients made up the major proportions (about 90%). The hospital committed to provide high-quality services to ED patients, where WT serves as a key metric. An example of such service-level requirements measured by the patient WT and length of stay (LOS) is shown in Table 1.

As seen from Table 1, patients in category I should be treated immediately with any available resources. The WT and service-level expectations for patients in categories II, III, and IV lower as the urgency drops, while no requirement is imposed on category V patients alone. The ED also maintains a goal of treating at least 98% of patients in all categories within 4 h, measured by LOS. While these requirements were established based on benchmarks and attainability, the focus was mainly on urgent patients. The hospital leadership believed that, with an efficient allocation of resources, it is possible to improve the service level for category IV patients (semiurgent and 48.76% of total) while meeting all other service requirements. Indeed, different categories of patients undergo different treatment procedures, and occupy different medical resources such as manpower, machines, and spaces. By properly assigning medical staff to each service area and each shift, the ED can find a scenario to improve the WT of category IV patients while maintaining the service for urgent patients. More specifically, define  $WT_2$ ,  $WT_3$ ,  $WT_4$ , and  $LOS$  as the daily average patient WT for categories II, III, IV, and the daily average LOS for all patients, respectively. Then, the aforementioned problem can be modeled as follows.

$$\max_{\mathbf{x}} \quad P\{WT_4(\mathbf{x}, \omega) \leq 120\} \quad (1)$$

$$\text{s.t.} \quad P\{WT_2(\mathbf{x}, \omega) \leq 15\} \geq 95\%, \quad (2)$$

$$P\{WT_3(\mathbf{x}, \omega) \leq 30\} \geq 90\%, \quad (3)$$

$$P\{LOS(\mathbf{x}, \omega) \leq 240\} \geq 98\%, \quad (4)$$

$$\mathbf{x} \in \mathcal{X} \quad (5)$$

where  $\omega$  represents random noises caused by uncertainties in patient arrivals and service durations, and  $\mathbf{x}$  is a vector of decision variables defining a staff assignment scenario that satisfies all the deterministic (nonprobabilistic) constraints defined by  $\mathcal{X}$ . For the above example, such deterministic con-

straints include the total budget constraint of personnel, and the minimum and maximum number of specialists required at each service area during each shift (e.g., no more than two nurses but no less than one nurse at the admission desk). Note that although the service level of category I patients is not explicitly modeled in the above formulation, any incoming category I patient should be treated immediately with available resources, thereby reducing the resource availability to other patients and their attendant service levels.

The aforementioned decision problem can be generalized to other service systems, such as call centers and public service offices. To this end, we define a stochastic performance measure  $F_h(\mathbf{x}, \omega)$  for  $h = 0, 1, \dots, H$ , where  $H$  is the total number of constraints and 0 indicates the objective function. Specifically, each  $F_h(\mathbf{x}, \omega)$  is a stochastic metric given a resource allocation scenario  $\mathbf{x}$ . For example,  $F_h(\mathbf{x}, \omega)$  may capture the average daily WT or peak daily LOS in the ED as required, where the value of this metric observed varies due to random noise  $\omega$ . Correspondingly, each performance goal is defined as  $d_h$  for  $h = 0, 1, \dots, H$ , and the desired service level as  $1 - \alpha_h$  for  $h = 1, \dots, H$ . Without loss of generality, we assume that (1) the lower  $d_h$  is, the better performance it has; and (2) the higher  $1 - \alpha_h$  is, the better service level the system provides. Hence, the generalized resource allocation problem for service systems can be formulated as follows:

$$(\mathcal{P}) \quad \min_{\mathbf{x}} \quad P\{F_0(\mathbf{x}, \omega) > d_0\} \quad (6)$$

$$\text{s.t.} \quad P\{F_h(\mathbf{x}, \omega) > d_h\} \leq \alpha_h \quad h = 1, \dots, H, \quad (7)$$

$$\mathbf{x} \in \mathcal{X}, \quad (8)$$

where  $\mathcal{X}$  is a finite discrete solution space of  $\mathbf{x}$ , determined by the deterministic constraints of the service system, such as budget and capacity constraints. In the remainder of this paper, we will focus on the problem  $\mathcal{P}$  defined by Equations (6)–(8) where  $\mathcal{X}$  is finite, and will refer to a resource allocation scenario  $\mathbf{x}$  as a solution of  $\mathcal{P}$ .

A simulation model will be used to evaluate the above probabilistic measures,  $P\{F_h(\mathbf{x}, \omega) > d_h\}$ , for every solution  $\mathbf{x}$ . Since most simulation models are expensive and time consuming to run and the decision may need to be made within a given time period, the challenge of solving problem  $\mathcal{P}$  using a simulation model in this context is the one of allocating the limited computing budget to select the optimal solution under observation noises, thereby falling under the OCBA regime. As mentioned in Section 2, the characteristics of problem  $\mathcal{P}$  warrant a new Bayesian OCBA formulation using the EOC measure, to be introduced in the sequel.

### 3.2 | A Bayesian OCBA formulation

Suppose there is a total of  $K$  alternative solutions (resource allocation scenarios) in  $\mathcal{X}$ , and  $\mathbf{x}_i$  represents one such solution for  $i \in \Theta = \{1, \dots, K\}$ . For a given  $\mathbf{x}_i$ , each simulation

replication  $j$  provides one estimate,  $F_h(\mathbf{x}_i, \omega_{ij})$ , for each performance measure  $h$  where  $\omega_{ij}$  is the random noise observed in the  $j$ -th replication. Here, we assume that for any solution  $\mathbf{x}_i$  and probabilistic measure  $h$ ,  $F_h(\mathbf{x}_i, \omega_{ij})$ 's are identically distributed for  $j = 1, 2, \dots$ . Further define an indicator function

$$X_{hij} \doteq \mathcal{I}(F_h(\mathbf{x}_i, \omega_{ij}) > d_h) = \begin{cases} 1 & \text{if } F_h(\mathbf{x}_i, \omega_{ij}) > d_h \\ 0 & \text{otherwise} \end{cases}. \quad (9)$$

It is seen that  $X_{hij}$  follows a Bernoulli distribution with a success probability  $p_{hi} \doteq P\{F_h(\mathbf{x}_i, \omega_{ij}) > d_h\} = \mathbb{E}[X_{hij}]$ . Apparently,  $p_{hi}$  is unknown in practice, and can be best estimated by its sample average  $\bar{p}_{hi} \doteq \frac{1}{n} \sum_{j=1}^n X_{hij}$ . To solve problem  $\mathcal{P}$  via simulation given the computing budget  $T$ , one needs to determine the number of replications allocated to evaluate each solution  $\mathbf{x}_i$  (denoted by  $N_i$ ), so as to best identify the true optimal solution  $\mathbf{x}_b$  where  $b = \arg \min_{i \in \Theta: p_{hi} \leq \alpha_h, h=1, \dots, H} p_{0i}$ .

To capture the prior knowledge on solutions and correlated samples of neighboring solutions, we proceed to develop an OCBA formulation from a Bayesian perspective. Specifically, since  $p_{hi}$  is practically unknown, the algorithm starts with a prior belief over each  $p_{hi}$  and updates its posterior distribution as simulation samples are gathered. To this end, let  $\Psi_h = (\Psi_{h1}, \dots, \Psi_{hK})$  be a random variable following a multivariate distribution of  $\pi_h^t(\psi_h)$ , that is,  $\Psi_h \sim \pi_h^t(\psi_h)$ , where  $\psi_h = (\psi_{h1}, \dots, \psi_{hK}) \in [0, 1]^K$  is a realization of  $\Psi_h$  and  $\pi_h^t(\psi_h)$  is the Bayesian belief over  $(p_{h1}, \dots, p_{hK})$  for service level  $h$  after  $t$  simulation replications are gathered. Thus,  $\pi_h^0(\psi_h)$  and  $\pi_h^T(\psi_h)$  represent the prior distribution before simulation and the posterior distribution after the computing budget  $T$  exhausts, respectively. Note that  $\pi_h^0(\psi_h)$  allows  $\Psi_{hi}$  and  $\Psi_{hi'}$  to be correlated based on prior knowledge on  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$ . For each realization  $\psi_h$  of  $\Psi_h$  for  $h = 0, 1, \dots, H$ , define  $\bar{b} = \arg \min_{i \in \Theta: \psi_{hi} \leq \alpha_h, h=1, \dots, H} \psi_{0i}$  as the best solution under realization  $\psi_h$ . Noting that  $(p_{h1}, \dots, p_{hK})$  is one possible realization of  $\Psi_h$ , if  $\psi_{hi} = p_{hi}$  for  $h = 0, 1, \dots, H$  and  $i = 1, \dots, K$ , then  $\bar{b} = b$ .

Here, we first make a mild assumption on the prior distribution (Russo, 2020) in Assumption 1.

**Assumption 1.** The prior joint distributions  $\pi_h^0(\psi_h)$  ( $h = 0, 1, \dots, H$ ) are continuous and uniformly bounded away from 0 and  $\infty$ , that is, there exist  $0 < \underline{c} < \bar{c} < \infty$  such that  $\underline{c} \leq \pi_h^0(\psi_h) \leq \bar{c}$  for all  $\psi_h \in [0, 1]^K$  and all  $h = 0, 1, \dots, H$ .

*Remark 1.* Assumption 1 is a mild assumption that covers a large class of distributions, such as the multivariate normal distribution (Howard, 1998) that we will use in the numerical experiments in Section 5. In particular, for any continuous function  $r(\psi)$  satisfying  $0 < \underline{r} \leq r(\psi) \leq \bar{r} < \infty$  for  $\psi \in [0, 1]^K$ , we can convert  $r(\psi)$  into a prior  $\pi_h^0(\psi_h) = \frac{r(\psi_h)}{\int_{[0,1]^K} r(\psi_h') d\psi_h'}$  that satisfies Assumption 1. When no prior

knowledge of the system is available, the modeler can use the uncorrelated uniform prior instead of the correlated prior, or estimate the prior distribution using the maximum likelihood estimation through initial samples (Rasmussen & Williams, 2006). In addition, it should be noted that the lower bound  $\underline{c}$  in Assumption 1 does not have to hold for the entire region of  $[0, 1]^K$  for  $\psi_h$ . In fact, we can relax this condition and only restrict  $\underline{c} \leq \pi_h^0(\psi_h)$  to hold for certain subregions of  $[0, 1]^K$  depending on  $h$  and solution  $\mathbf{x}_i$ . This will become more apparent in the proof of Theorem 1, and we only mention this fact here without stating the detailed conditions. However, since such tighter conditions would overcomplicate the elaboration of Lemma 2, we prefer the succinct statement in Assumption 1 and will stick to it throughout the paper.

Let  $N_i^{(t)}$  be the number of simulation replications allocated to solution  $\mathbf{x}_i$  when  $t$  replications are gathered, and  $L_t(\psi_{hi})$  be the corresponding likelihood function (of observations). Note that  $N_i^{(0)} = 0$ ,  $N_i^{(T)} = N_i$ , and  $\sum_{i \in \Theta} N_i^{(t)} = t$ . Further,  $L_t(\psi_{hi}) = (\psi_{hi})^{\sum_{j=1}^{N_i^{(t)}} X_{hij}} (1 - \psi_{hi})^{\sum_{j=1}^{N_i^{(t)}} (1 - X_{hij})}$ . The posterior distribution can be updated according to the Bayes' rule as follows:

$$\pi_h^t(\psi_h) = \frac{\pi_h^0(\psi_h) \prod_{i=1}^K L_t(\psi_{hi})}{\int_{[0,1]^K} \pi_h^0(\psi_h') \prod_{i=1}^K L_t(\psi_{hi}') d\psi_h'}, \quad h = 0, 1, \dots, H, \quad t \leq T. \quad (10)$$

The posterior distribution in Equation (10) may not have a closed form if the prior distribution is correlated. In this case, we need numerical integration methods to solve integrals in Equation (10). In the case of high-dimensional integrals if the total number of solutions,  $K$ , is large, Monte Carlo estimation and the importance sampling method can provide robust estimates with high accuracy.

We are now ready to introduce the OCBA formulation that minimizes the Bayesian EOC. More specifically, based on  $T$  simulation samples, the selected solution  $\mathbf{x}_{\hat{b}}$  is the one estimated with the minimum posterior mean in the objective function while satisfying all constraints; that is,  $\hat{b} = \arg \min_{i \in \Theta: \mathbb{E}_T[\Psi_{hi}] \leq \alpha_h, h=1, \dots, H} \mathbb{E}_T[\Psi_{0i}]$ , where

$$\mathbb{E}_t[\Psi_{hi}] = \int_{[0,1]^K} \psi_{hi} \pi_h^t(\psi_h) d\psi_h, \quad h = 0, 1, \dots, H, \quad i = 1, \dots, K, \quad t \leq T \quad (11)$$

is the posterior mean. The Bayesian EOC is then defined as

$$EOC_{Bayes} = \int \dots \int_{[0,1]^{K \times (H+1)}} \sum_{h=0}^H OC_{hb} \pi_0^T(\psi_0) \pi_1^T(\psi_1) \dots \pi_H^T(\psi_H) d\psi_0 d\psi_1 \dots d\psi_H, \quad (12)$$

where  $OC_{hb} = \varpi_h \max(\psi_{hb} - \alpha_h, 0)$ ,  $h = 1, \dots, H$ , and

$$OC_{0b} = \begin{cases} \max\{\psi_{0b} - \psi_{0\tilde{b}}, 0\}, & \text{if } \tilde{b} = \arg \min_{i \in \Theta: \psi_{hi} \leq \alpha_h, h=1, \dots, H} \psi_{0i} \text{ exists,} \\ 1, & \text{if } \arg \min_{i \in \Theta: \psi_{hi} \leq \alpha_h, h=1, \dots, H} \psi_{0i} \text{ does not exist.} \end{cases} \quad (13)$$

Here, the weight  $0 < \varpi_h < \infty$  for  $h = 1, \dots, H$  reflects the unit cost of violating constraint  $h$  relative to the unit cost of the objective. Specifically, the opportunity cost  $OC_{0b}$  is the penalty for the objective of the selected solution  $\mathbf{x}_{\tilde{b}}$ , and the cost  $OC_{hb}$  penalizes the loss incurred if  $\mathbf{x}_{\tilde{b}}$  violates the constraint on service level  $h$ . Note that when  $\arg \min_{i \in \Theta: \psi_{hi} \leq \alpha_h, h=1, \dots, H} \psi_{0i}$  does not exist, we still recommend the solution  $\mathbf{x}_{\tilde{b}}$  when the constrained optimization problem is infeasible. Hence, we set the attendant opportunity cost to be 1 to penalize such selection, noting that the opportunity cost of  $\max\{\psi_{0b} - \psi_{0\tilde{b}}, 0\} \leq 1$  when  $\tilde{b}$  exists. A large weight  $\varpi_h$  can be imposed to penalize against selecting an infeasible solution without affecting the asymptotic results to be derived later.

Subsequently, the Bayesian OCBA formulation of problem  $\mathcal{P}$  is defined as follows:

$$\begin{aligned} (\text{OCBA}-\mathcal{P}) \quad & \min_{N_1, \dots, N_K} EOC_{Bayes} \\ \text{s.t.} \quad & \sum_{i=1}^K \frac{N_i}{T} = 1, \\ & N_i \geq 0, \quad i = 1, \dots, K. \end{aligned} \quad (14)$$

We assume that in problem  $\mathcal{P}$ , the constraint measures are not exactly equal to the constraint limit on the right-hand side; that is,  $p_{hi} \neq \alpha_h$  for all constraints in any solution. This is a common assumption made in the OCBA literature, which ensures that  $EOC_{Bayes}$  approaches 0 as the computing budget increases.

The  $\text{OCBA}-\mathcal{P}$  formulation distinguishes from other OCBA models, such as the one in Lee et al. (2012) for constrained R&S, not only in the use of the EOC objective, but also in capturing the performance correlation of solutions rather than assuming independence in solution measures. Intuitively speaking, capturing such correlation enables more effective learning of solution performances from samples, since the performance of a solution is not only learned using its own samples, but also using samples from its neighboring solutions. On the other hand, it brings challenges in developing theoretical results, which will be addressed next.

## 4 | SOLUTION METHODOLOGY

Although  $EOC_{Bayes}$  in Equation (12) is very complex to analyze (Chick et al., 2010), we will show that as  $T \rightarrow \infty$ ,  $\tilde{b}$  exists almost surely, and  $EOC_{Bayes}$  is well defined and decays in a deterministic rate. Therefore, we will next solve  $\text{OCBA}-\mathcal{P}$  by optimizing the approximated rate of decay of  $EOC_{Bayes}$ .

### 4.1 | Rate of decay of Bayesian EOC

Considering the definition of the selected best  $\tilde{b}$ ,  $\tilde{b}$  does not always exist due to simulation noises, even if the optimal solution exists for problem  $\mathcal{P}$ . Particularly, if  $N_i$  remains fixed ( $N_i < \infty$ ) for some solution  $\mathbf{x}_i$  as  $T \rightarrow \infty$ ,  $\tilde{b}$  may not exist since  $\{i \in \Theta : \mathbb{E}_T[\Psi_{hi}] \leq \alpha_h, h = 1, \dots, H\}$  could be an empty set. To avoid this unfavorable situation, we make the following assumption:

**Assumption 2.** As  $T \rightarrow \infty$ ,  $N_i \rightarrow \infty$  and  $\rho_i \doteq \lim_{T \rightarrow \infty} \frac{N_i}{T}$  exists for all  $i = 1, \dots, K$ .

*Remark 2.* In fact, under Assumption 2, we will show in the sequel that  $\tilde{b}$  exists given that  $\mathcal{P}$  is feasible and  $\lim_{T \rightarrow \infty} \frac{1}{T} \log EOC_{Bayes} \leq 0$  as  $T \rightarrow \infty$ . Meanwhile, it can be shown that if  $N_i < \infty$  for some  $\mathbf{x}_i$  as  $T \rightarrow \infty$ , either  $\tilde{b}$  does not exist or  $\lim_{T \rightarrow \infty} \frac{1}{T} \log EOC_{Bayes} = 0$ , which is an inferior allocation compared to the allocation under Assumption 2.

To facilitate the analysis, we first divide the set of indices for nonbest solutions into two exhaustive and mutually exclusive sets,  $\Theta_O$  and  $\Theta_F$ , such that  $\Theta = \Theta_O \cup \Theta_F \cup \{b\}$ .  $\Theta_O$  is the set of indices for nonbest solutions that are feasible, whereas  $\Theta_F$  is the set of indices for nonbest solutions with at least one constraint in Equation (7) violated. Further defining the set of constraints violated for any solution  $\mathbf{x}_i$  as  $\Pi_i = \{h : p_{hi} > \alpha_h, h \in \{1, \dots, H\}\}$ , we have

$$\Theta_O \doteq \{i \in \Theta : i \neq b, \Pi_i = \emptyset\}, \quad \Theta_F \doteq \{i \in \Theta : \Pi_i \neq \emptyset\}. \quad (15)$$

Furthermore, for any  $i \in \Theta_F$ , we define  $h_i^* = \arg \max_{h \in \Pi_i} (p_{hi} - \alpha_h)$  as the “most violated” constraint. Lemma 1 then provides an approximate and upper bound of  $EOC_{Bayes}$  (called  $AEOC_{Bayes}$ ).

**Lemma 1.**  $EOC_{Bayes}$  in Equation (12) is bounded from above by the following  $AEOC_{Bayes}$  almost surely:

$$\begin{aligned} EOC_{Bayes} &\leq AEOC_{Bayes} \\ &= \sum_{h=1}^H \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h \\ &\quad + \sum_{i \in \Theta_O} \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} (\psi_{0b} - \psi_{0i}) \pi_0^T(\psi_0) d\psi_0 \\ &\quad + \sum_{i \in \Theta_F} \int_{\{\psi_{h_i^*}: \psi_{h_i^*i} \leq \alpha_{h_i^*}\}} \pi_{h_i^*}^T(\psi_{h_i^*}) d\psi_{h_i^*}, \end{aligned} \quad (16)$$

where  $h_i^* = \arg \max_{h \in \Pi_i} (p_{hi} - \alpha_h)$  for any  $i \in \Theta_F$ .

*Proof.* The proof is given in the Supporting Information EC.1.1 to this paper.  $\square$

**Remark 3.** Although  $AEOC_{Bayes}$  is an upper bound of  $EOC_{Bayes}$ , it is constructed to be as tight as possible such that it closely approximates  $EOC_{Bayes}$ . The proof of Lemma 1 in the Supporting Information EC.1.1 shows why this bound is tight, where the Bonferroni inequality is applied to reach Equation (EC.4), and subspace is constructed to provide a tight bound on the EOC in Equation (EC.5) for both cases of  $i \in \Theta_O$  and  $i \in \Theta_F$ . Further,  $AEOC_{Bayes}$  becomes tighter as the total simulation budget becomes larger, and  $AEOC_{Bayes} \rightarrow 0$  as  $T \rightarrow \infty$  (Russo, 2020).

We now proceed to analyze the rate of decay of  $AEOC_{Bayes}$ . We first present Lemma 2 and Lemma 3, and then develop Theorem 1 on the rate of decay using the lemmas.

**Lemma 2.** For any  $0 \leq a_1 < a_2 \leq 1$ ,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_{hi}: a_1 \leq \psi_{hi} \leq a_2\}} \pi_h^T(\psi_{hi}) d\psi_{hi} = -\rho_i \inf_{\psi_{hi} \in [a_1, a_2]} I_{hi}(\psi_{hi}), \quad (17)$$

where

$$I_{hi}(\psi_{hi}) = p_{hi} \log \frac{p_{hi}}{\psi_{hi}} + (1 - p_{hi}) \log \frac{1 - p_{hi}}{1 - \psi_{hi}}. \quad (18)$$

*Proof.* The proof is given in the Supporting Information EC.1.2 to this paper.  $\square$

**Lemma 3** (Ganesh et al. (2004)). Consider positive sequences  $a_j(n)$ ,  $j = 1, \dots, m$ . If  $\lim_{n \rightarrow \infty} \frac{1}{n} \log a_j(n)$  exists for all  $j$ , then  $\lim_{n \rightarrow \infty} \frac{1}{n} \log(\sum_{j=1}^m a_j(n)) = \max_{j \in \{1, \dots, m\}} \{ \lim_{n \rightarrow \infty} \frac{1}{n} \log a_j(n) \}$ .

**Theorem 1.** The rate of decay of  $AEOC_{Bayes}$  in Equation (16) is bounded by the following AEOC-B:

$$\begin{aligned} & -\lim_{T \rightarrow \infty} \frac{1}{T} \log AEOC_{Bayes} \geq AEOC-B \\ & = \min \left\{ \min_{h=1, \dots, H} \rho_b I_{hb}(\alpha_h), \min_{i \in \Theta_O} \rho_b I_{0b}(c_i), \right. \\ & \quad \left. \min_{i \in \Theta_O} \rho_i I_{0i}(c_i), \min_{i \in \Theta_F} \rho_i I_{hi}^*(\alpha_{hi}^*) \right\}, \end{aligned} \quad (19)$$

where  $p_{0b} < c_i < p_{0i}$  for any  $i \in \Theta_O$ , and  $I_{hi}(x)$  is given by Equation (18) for any  $h = 0, 1, \dots, H$  and  $i = 1, \dots, K$ .

*Proof.* By Equation (16) and Lemma 3, we have

$$\begin{aligned} & -\lim_{T \rightarrow \infty} \frac{1}{T} \log AEOC_{Bayes} \\ & = -\max \left\{ \max_{h=1, \dots, H} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h, \right. \\ & \quad \left. (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h, \right. \end{aligned} \quad (20)$$

$$\max_{i \in \Theta_O} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} (\psi_{0b} - \psi_{0i}) \pi_0^T(\psi_0) d\psi_0, \quad (21)$$

$$\max_{i \in \Theta_F} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_{hi}^*: \psi_{hi}^* \leq \alpha_{hi}^*\}} \pi_{hi}^T(\psi_{hi}^*) d\psi_{hi}^* \}. \quad (22)$$

Consider the term in Equation (20) for any  $h = 1, \dots, H$ . Since  $1 \leq 1 + \varpi_h(\psi_{hb} - \alpha_h) \leq 1 + \varpi_h$  for the domain  $\{\psi_h: \psi_{hb} > \alpha_h\}$ , we have

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h \\ & \geq \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} \pi_h^T(\psi_h) d\psi_h \end{aligned} \quad (23)$$

and

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h \\ & \leq \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h) \pi_h^T(\psi_h) d\psi_h \\ & = \lim_{T \rightarrow \infty} \frac{1}{T} \log(1 + \varpi_h) + \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} \pi_h^T(\psi_h) d\psi_h \\ & = \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} \pi_h^T(\psi_h) d\psi_h \end{aligned} \quad (24)$$

resulting

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} (1 + \varpi_h(\psi_{hb} - \alpha_h)) \pi_h^T(\psi_h) d\psi_h \\ & = \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} \pi_h^T(\psi_h) d\psi_h. \end{aligned} \quad (25)$$

Next, consider the term in Equation (21) for any  $i \in \Theta_O$ . For the domain  $\{\psi_0: \psi_{0i} < \psi_{0b}\}$ , there exists a  $\delta > 0$  such that  $\delta \leq \psi_{0b} - \psi_{0i} \leq 1$ . Following the similar proof as above, we have

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} (\psi_{0b} - \psi_{0i}) \pi_0^T(\psi_0) d\psi_0 \\ & = \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0. \end{aligned} \quad (26)$$

Considering the integrand on the right-hand side of Equation (26), define a constant  $c_i$  where  $p_{0b} < c_i < p_{0i}$ . Since  $\{\psi_{0i} < \psi_{0b}\} \subset \{\psi_{0i} \leq c_i\} \cup \{c_i \leq \psi_{0b}\}$ , we have

$$\begin{aligned} & \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0 \leq \int_{\{\psi_0: \psi_{0i} \leq c_i\}} \pi_0^T(\psi_0) d\psi_0 \\ & + \int_{\{\psi_0: c_i \leq \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0. \end{aligned} \quad (27)$$

From Lemma 2,  $\lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: a_1 \leq \psi_{hi} \leq a_2\}} \pi_h^T(\psi_h) d\psi_h$  exists for any  $0 \leq a_1 < a_2 \leq 1$ . Thus, by Equation (26), Equation (27), and Lemma 3,

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} (\psi_{0b} - \psi_{0i}) \pi_0^T(\psi_0) d\psi_0 \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} < \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0 \\ &\leq \max \left\{ \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} \leq c_i\}} \pi_0^T(\psi_0) d\psi_0, \right. \\ & \quad \left. \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: c_i \leq \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0 \right\}. \end{aligned} \quad (28)$$

Substituting Equations (25) and (28) into Equations (20)–(22), we have

$$\begin{aligned} & - \lim_{T \rightarrow \infty} \frac{1}{T} \log AEOC_{Bayes} \\ & \geq - \max \left\{ \max_{h=1, \dots, H} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_h: \psi_{hb} > \alpha_h\}} \pi_h^T(\psi_h) d\psi_h, \right. \\ & \quad \max_{i \in \Theta_O} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: \psi_{0i} \leq c_i\}} \pi_0^T(\psi_0) d\psi_0, \\ & \quad \max_{i \in \Theta_O} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_0: c_i \leq \psi_{0b}\}} \pi_0^T(\psi_0) d\psi_0 \\ & \quad \left. \max_{i \in \Theta_F} \lim_{T \rightarrow \infty} \frac{1}{T} \log \int_{\{\psi_{h_i^*}: \psi_{h_i^*} \leq \alpha_{h_i^*}\}} \pi_{h_i^*}^T(\psi_{h_i^*}) d\psi_{h_i^*} \right\}. \end{aligned} \quad (29)$$

Applying Lemma 2 to the right-hand side of Equation (29), Equation (19) readily follows, since  $I_{hi}(x)$  is a convex function with the minimum achieved at  $x = p_{hi}$ .  $\square$

We mention that the only approximation used in Theorem 1 to arrive at  $AEOC-B$  is Equation (28). The value of  $c_i$  should depend only on  $p_{0b}$  and  $p_{0i}$  and affects the tightness of  $AEOC-B$ . We will discuss the choice of  $c_i$  in the Supporting Information EC.2 to this paper.

## 4.2 | Optimality conditions

Since neither minimizing  $EEOC_{Bayes}$  in problem  $\mathcal{OCBA-P}$  nor maximizing its rate of decay is mathematically tractable, we instead maximize the bound of the rate of decay of  $AEOC_{Bayes}$ , that is,  $AEOC-B$  defined in the right-hand side of Equation (19). Letting  $z = \min\{\min_{h=1, \dots, H}$

$\rho_b I_{hb}(\alpha_h), \min_{i \in \Theta_O} \rho_b I_{0b}(c_i), \min_{i \in \Theta_O} \rho_i I_{0i}(c_i), \min_{i \in \Theta_F} \rho_i I_{h_i^*}(\alpha_{h_i^*})\}$ , the resulting problem can be expressed as

$$(\mathcal{OCBA-P-B}) \quad \max_{z, \rho_1, \dots, \rho_K} z \quad (30)$$

$$\text{s.t.} \quad z \leq \rho_b \min \left\{ \min_{h=1, \dots, H} I_{hb}(\alpha_h), \min_{k \in \Theta_O} I_{0b}(c_k) \right\}, \quad (31)$$

$$z \leq \rho_i I_{0i}(c_i), \quad \forall i \in \Theta_O, \quad (32)$$

$$z \leq \rho_i I_{h_i^*}(\alpha_{h_i^*}), \quad \forall i \in \Theta_F, \quad (33)$$

$$\sum_{i=1}^K \rho_i = 1, \quad (34)$$

$$\rho_i \geq 0, \quad \forall i \in \Theta. \quad (35)$$

Problem  $\mathcal{OCBA-P-B}$  in Equations (30)–(35) is a linear optimization problem, whose optimal solution stated in Theorem 2 provides an asymptotic optimality condition for problem  $\mathcal{OCBA-P}$ .

**Theorem 2.** Let  $\eta_i$  be as follows:

$$\eta_i = \begin{cases} \min \left\{ \min_{h=1, \dots, H} I_{hb}(\alpha_h), \min_{k \in \Theta_O} I_{0b}(c_k) \right\}, & i = b \\ I_{0i}(c_i), & i \in \Theta_O \\ I_{h_i^*}(\alpha_{h_i^*}), & i \in \Theta_F, \end{cases} \quad (36)$$

where  $I_{hi}(\cdot)$  is given by Equation (18). As  $T \rightarrow \infty$ ,  $AEOC-B$  is minimized if

$$\frac{\rho_i^*}{\rho_j^*} = \frac{\eta_j}{\eta_i}, \quad \forall i, j \in \Theta, \quad (37)$$

and thus, the asymptotic optimal sample allocation strategy is given by

$$N_i^* = \frac{1/\eta_i}{\sum_{k=1}^K 1/\eta_k} T, \quad \forall i \in \Theta. \quad (38)$$

*Proof.* Since  $\rho_i \doteq \lim_{T \rightarrow \infty} \frac{N_i}{T}$  exists from Assumption 2 and  $\sum_{i=1}^K \rho_i = 1$  with  $\rho_i \geq 0$ , we can treat  $\rho_i$  as continuous variables between 0 and 1. Therefore, problem  $\mathcal{OCBA-P-B}$  becomes a linear program with continuous variables. As such, the optimal solution is obtained when the right-hand sides of all constraints (31)–(33) are equal. To see why,  $z$  can only be decreased further from this solution if all  $\rho_i$  for  $i \in \Theta$  are decreased simultaneously so as to decrease all right-hand sides in Equations (31)–(33). However, decreasing one  $\rho_i$  will increase another since  $\rho_b + \sum_{i \in \Theta_O} \rho_i + \sum_{i \in \Theta_F} \rho_i = 1$  from Equation (34). Hence, by the definition in Equation (36), we

have

$$\rho_i^* \eta_i = \rho_j^* \eta_j, \quad \forall i, j \in \Theta. \quad (39)$$

Note that  $I_{hi}(\psi_{hi})$  in Equation (18) is a convex function with its minimum value of 0 achieved at  $\psi_{hi} = p_{hi}$ . By the definition of problem  $\mathcal{OCBA-P}$ ,  $p_{hi} \neq \alpha_h$  for all  $h$  and  $i$ , and  $p_{0b} < c_i < p_{0i}$  by choice, we have  $\eta_i > 0$ . Thus, the optimality condition in Equation (37) readily follows, so as the asymptotically optimal sample allocation strategy in Equation (38) considering the constraint (34).  $\square$

Intuitively speaking, the simulation replications allocated to solution  $\mathbf{x}_i$  should be proportional to the inverse of its rate function  $\eta_i$  given by Equation (38). Observe that the rate function in Equation (18) is a convex function with its minimum achieved at  $\psi_{hi} = p_{hi}$ . The closer  $\psi_{hi}$  is to  $p_{hi}$ , the smaller  $I_{hi}(\psi_{hi})$  is. Therefore, the intuition is that we should increase the sample size to a solution if either optimality or feasibility is harder to detect. More specifically, the sample size allocated to solution  $\mathbf{x}_i$  should become larger as  $c_i$  gets closer to  $p_{0i}$  for any  $i \in \Theta_O$ , and as the constraint violation becomes less substantial for any  $i \in \Theta_F$ ; for the best solution  $\mathbf{x}_b$ , its sample size should be increased if one of the attendant service levels  $p_{hb}$  is close to the required service level  $\alpha_h$ , or if one of the  $c_i$  gets close to  $p_{0b}$ . This intuition is consistent with the standard OCBA allocation strategy for normally distributed random variables in Chen et al. (2000), where more samples are allocated to solutions whose optimality is harder to decide, or whose variance is larger.

Theorem 2 provides a simple closed-form formula for allocating the total computing budget to all solutions, such that the  $AEBC-B$  is asymptotically optimized. That is, such an allocation strategy optimizes the Bayesian EOC of selecting the estimated best solution  $\hat{b}$  given a sufficiently large computing budget. The allocation strategy in Theorem 2 may not be optimal when  $T$  is relatively small. However, the performance shall get better as  $T$  becomes larger.

**Remark 4.** It may be of interest to consider two degenerate cases of problem  $\mathcal{P}$ . The first case is an unconstrained version of problem  $\mathcal{P}$ , where stochastic constraints in Equation (7) no longer exist. In this case, Theorem 2 still holds, except that  $H = 0$ ,  $\Theta_F = \emptyset$  (no infeasible solution for the unconstrained problem), and  $\Theta_O = \Theta \setminus \{b\}$ , resulting in a degenerate case of Equation (36) given by

$$\eta_i = \begin{cases} \min_{k \in \Theta \setminus \{b\}} I_{0b}(c_k), & i = b \\ I_{0i}(c_i), & i \in \Theta \setminus \{b\} \end{cases}. \quad (40)$$

The second case is a feasibility detection variant of problem  $\mathcal{P}$ , where the objective function in Equation (6) no longer exists. The goal is to find all feasible solutions meeting the service targets, that is,

$$\begin{aligned} (\mathcal{P}_f) \quad & \text{all } \mathbf{x} \in \mathcal{X} \\ \text{s.t.} \quad & P\{F_h(\mathbf{x}, \omega) > d_h\} \leq \alpha_h \quad h = 1, \dots, H. \end{aligned} \quad (41)$$

Let  $\mathcal{S}_f = \{i \in \Theta : \mathbf{x}_i \in \mathcal{X} \text{ and } P\{F_h(\mathbf{x}_i, \omega) > d_h\} \leq \alpha_h \text{ for } h = 1, \dots, H\}$  represent the set of true feasible solutions and let  $\hat{\mathcal{S}}_f = \{i \in \Theta : \mathbf{x}_i \in \mathcal{X} \text{ and } \mathbb{E}_T[\Psi_{hi}] \leq \alpha_h \text{ for } h = 1, \dots, H\}$  represent the set of feasible solutions estimated from samples. The Bayesian EOC for problem  $\mathcal{P}_f$  can be defined as

$$\begin{aligned} EOC_{f:Bayes} = & \int \cdots \int_{[0,1]^{K \times H}} \sum_{i \in \mathcal{S}_f} \sum_{h=1}^H OC_{hi}^{f_1} \pi_1^T(\psi_1) \\ & \cdots \pi_H^T(\psi_H) d\psi_1 \cdots d\psi_H \\ & + \int_{\{\psi_1: \psi_{1i} \leq \alpha_i\}} \cdots \int_{\{\psi_H: \psi_{Hi} \leq \alpha_H\}} \sum_{i \notin \mathcal{S}_f} \sum_{h=1}^H OC_{hi}^{f_2} \pi_1^T(\psi_1) \\ & \cdots \pi_H^T(\psi_H) d\psi_1 \cdots d\psi_H \end{aligned} \quad (42)$$

where  $OC_{hi}^{f_1} = \varpi_h \max\{\psi_{hi} - \alpha_h, 0\}$  and  $OC_{hi}^{f_2} = \varpi_h \max\{\alpha_h - \psi_{hi}, 0\}$ ,  $h = 1, \dots, H$ . Using analysis similar to Theorem 1, we can show that the rate of decay of  $EOC_{f:Bayes}$  is bounded by

$$\min \left\{ \min_{i \in \mathcal{S}_f} \rho_i \min_{h=1, \dots, H} I_{hi}(\alpha_h), \min_{i \notin \mathcal{S}_f} \rho_i I_{hi}^*(\alpha_{hi}^*) \right\}. \quad (43)$$

Maximizing Equation (43), the asymptotic optimal sample allocation strategy has the same form as Equation (38), where  $\eta_i$  in Equation (36) becomes

$$\eta_i = \begin{cases} \min_{h=1, \dots, H} I_{hi}(\alpha_h), & i \in \mathcal{S}_f \\ I_{hi}^*(\alpha_{hi}^*), & i \notin \mathcal{S}_f \end{cases}. \quad (44)$$

### 4.3 | Iterative algorithm and convergence analysis

Recall that  $\eta_i$  in Theorem 2 is computed using the true success probability of Bernoulli variables,  $p_{hi}$ . However,  $p_{hi}$  is practically unknown, but can be estimated by the posterior mean. Therefore, we propose an iterative algorithm in Algorithm 1 to sequentially allocate simulation replications to the solutions in consideration, by implementing Theorem 2.

In Step 5 of Algorithm 1, for  $i \in \hat{\Theta}_O$ ,  $\hat{c}_i$  can be chosen between  $\hat{p}_{0b}$  and  $\hat{p}_{0i}$ , and we will show in the Supporting Information EC.2 that  $\hat{c}_i = \frac{\hat{p}_{0b} + \hat{p}_{0i}}{2}$  is a recommended choice. Further, parameter  $\eta_i$  is strictly positive by definition, and thus  $\hat{\eta}_i^r$  is bounded by a sufficiently small positive number in Step 5. Finally, in Step 6, it may be necessary to round the values computed by Equation (47) to integers, for example, by preserving the total sum while minimizing the maximum relative deviation of such rounding.

While the asymptotic performance of Algorithm 1 is guaranteed by Theorem 2, a more important practical concern is its finite-time performance. Specifically, two metrics are critical: (1) How close is the posterior mean  $\hat{p}_{hi}$  to the true success probability  $p_{hi}$ ? and (2) How close is the empirical sample

**ALGORITHM 1** Iterative Sample Allocation Algorithm for Solving  $OCBA-P-B$ 

0. **Initialize:** Set iteration counter  $r = 0$  and the incremental budget  $\Delta$ . Set initial sample size  $N_i^r = n_0$  and incremental sample size  $\delta_i^r = N_i^r$  for all  $i \in \Theta$ . Specify the prior  $\pi_h^0(\psi_h)$  for  $h = 0, 1, \dots, H$ .
1. **Simulate:** Run simulation for each solution  $\mathbf{x}_i$ ,  $i \in \Theta$ , for  $\delta_i^r$  replications. For the  $j$ -th simulation replication of  $\mathbf{x}_i$ , record the simulation output  $F_h(\mathbf{x}_i, \omega_{ij})$ , and compute  $X_{hij}$  using Equation (9) for all  $h$ .
2. **Estimate:** Update the posterior distributions  $\pi_h^t(\psi_h)$  using Equation (10) for all  $h$ , where  $t = \sum_{i=1}^K N_i^r$ . Estimate the success probability of service level  $h$  for solution  $\mathbf{x}_i$  as  $\hat{p}_{hi} = \mathbb{E}_t[\Psi_{hi}]$  using Equation (11) for all  $h$  and  $i \in \Theta$ .
3. **Divide:** Compute the set of violated constraints for any  $i \in \Theta$  as  $\hat{\Pi}_i = \{h \in \{1, \dots, H\} : \hat{p}_{hi} > \alpha_h\}$ . Then, find the sample best solution:

$$\hat{b} = \arg \min_{i \in \Theta: \hat{\Pi}_i = \emptyset} \hat{p}_{0i}, \quad (45)$$

and divide the non-best solutions into two sets:

$$\hat{\Theta}_O = \{i | i \neq \hat{b}, \hat{\Pi}_i = \emptyset\}, \quad \hat{\Theta}_F = \{i | \hat{\Pi}_i \neq \emptyset\}. \quad (46)$$

4. **Stop:** If  $t \geq T$ , terminate and output  $\hat{b}$  as the best solution; otherwise, continue.
5. **Update:** For each  $i \in \hat{\Theta}_O$ , choose  $\hat{c}_i = w_i \hat{p}_{0b} + (1 - w_i) \hat{p}_{0i}$  where  $0 < w_i < 1$  with  $w_i = \frac{1}{2}$  recommended; and for each  $i \in \hat{\Theta}_F$ , choose  $\hat{h}_i^* = \arg \max_{h \in \Pi_i} (\hat{p}_{hi} - \alpha_h)$ . Then, for each  $i \in \Theta$ , compute  $\hat{\eta}_i^r$  using Equation (36), where  $\Theta_F$  and  $\Theta_O$  are respectively approximated by  $\hat{\Theta}_F$  and  $\hat{\Theta}_O$ , and  $I_{hi}(\cdot)$  is computed using Equation (18) with  $p_{hi}$  approximated by  $\hat{p}_{hi}$ . If  $\hat{\eta}_i^r = 0$ , bound  $\hat{\eta}_i^r$  by a sufficiently small positive number  $\kappa$ :  $\hat{\eta}_i^r \leftarrow \kappa$ .
6. **Allocate:** Increase the computing budget by  $\Delta$ , and determine the new allocation using

$$\hat{N}_i^{(*,r)} = \frac{1/\hat{\eta}_i^r}{\sum_{k=1}^K 1/\hat{\eta}_k^r} (n_0 K + (r+1)\Delta), \quad \forall i \in \Theta. \quad (47)$$

Allocate additional  $\delta_i^{r+1} = \max\{\hat{N}_i^{(*,r)} - N_i^r, 0\}$  simulation replications to solution  $\mathbf{x}_i$ . Update  $N_i^{r+1} \leftarrow N_i^r + \delta_i^{r+1}$ , and  $r \leftarrow r + 1$ . Go to Step 1.

allocation ratio in Equation (47) to the asymptotically optimal one in Equation (38)? For the latter, we define  $\rho_i^* = \frac{1/\eta_i}{\sum_{k=1}^K 1/\eta_k}$  as the asymptotic optimal ratio of samples allocated to solution  $\mathbf{x}_i$ , and  $\hat{\rho}_i^r = \frac{1/\hat{\eta}_i^r}{\sum_{k=1}^K 1/\hat{\eta}_k^r}$  as the corresponding empirical ratio at iteration  $r$ . To this end, Theorem 3 characterizes the finite-time performance of Algorithm 1, with the proof supported by Lemma 4.

**Lemma 4.** Suppose the prior joint distributions  $\pi_h^0(\psi_h)$  ( $h = 0, 1, \dots, H$ ) satisfy Assumption 1. Further, suppose at iteration  $r$  of Algorithm 1,  $N_i^r$  samples have been allocated to  $\mathbf{x}_i$  out of a total of  $t$  samples, and define the estimated success probability  $\hat{p}_{hi} \doteq \mathbb{E}_t[\Psi_{hi}]$  from Equation (11). Then, for any  $\gamma > 0$ ,

$\mathbb{E}[e^{\gamma W}] < \infty$ , where

$$W = \max_{r \in \mathbb{N} \cup \{0\}} \max_{i \in \Theta} \max_{h=0,1,\dots,H} \sqrt{\frac{N_i^r + 1}{\log(e + N_i^r)}} |\hat{p}_{hi} - p_{hi}|. \quad (48)$$

*Proof.* The proof is given in the Supporting Information EC.1.3 to this paper.  $\square$

**Theorem 3.** Suppose Assumption 1 holds. For any sufficiently small  $\epsilon > 0$ , let

$$R^\epsilon = \inf \{R \in \mathbb{N} : |\hat{p}_{hi} - p_{hi}| \leq \epsilon \text{ and } |\hat{\rho}_i^r - \rho_i^*| \leq \epsilon \text{ for all } i \in \Theta, h = 0, 1, \dots, H, \text{ and } r \geq R\} \quad (49)$$

be the earliest iteration such that for all iterations  $r \geq R^\epsilon$  in Algorithm 1, the posterior mean  $\hat{p}_{hi}$  and the empirical allocation ratio  $\hat{\rho}_i^r$  are sufficiently close to their asymptotic counterparts,  $p_{hi}$  and  $\rho_i^*$ , respectively. Further, as the computing budget  $T \rightarrow \infty$ ,  $\mathbb{E}[R^\epsilon] < \infty$  and  $\text{Var}[R^\epsilon] < \infty$ .

*Proof.* The proof is given in the Supporting Information EC.1.4 to this paper.  $\square$

**Remark 5.** By the Chebyshev inequality, the finite expectation and variance indicate that by a very high probability,  $EOC_{Bayes}$  will converge at the almost optimal convergence rate after a finite number of iterations. For example, let  $r^* = \mathbb{E}[R^\epsilon] + \sqrt{10\text{Var}[R^\epsilon]}$ , the Chebyshev inequality tells us that  $|\hat{p}_{hi} - p_{hi}| \leq \epsilon$  and  $|\rho_i^* - \hat{\rho}_i^r| \leq \epsilon$  for all  $i \in \Theta$ ,  $h = 0, 1, \dots, H$ , and  $r > r^*$  with a probability greater than 90%.

Note that by Equation (19), the empirical value of  $AEOC-B$  provided by Algorithm 1 converges to its theoretically optimal counterpart as the sample allocation converges. That is,  $|\rho_i^* - \hat{\rho}_i^r| \leq \epsilon$  means

$$\begin{aligned} & \left| \min \left( \min_{h=1,2,\dots,H} \rho_b^* I_{hb}(\alpha_h), \min_{i \in \Theta_O} \rho_b^* I_{0b}(c_i), \right. \right. \\ & \quad \left. \min_{i \in \Theta_F} \rho_i^* I_{h_i^* i}(\alpha_{h_i^*}), \min_{i \in \Theta_O} \rho_i^* I_{0i}(c_i) \right) \\ & \quad - \min \left( \min_{h=1,2,\dots,H} \hat{\rho}_b^r I_{hb}(\alpha_h), \min_{i \in \Theta_O} \hat{\rho}_b^r I_{0b}(c_i), \right. \\ & \quad \left. \min_{i \in \Theta_F} \hat{\rho}_i^r I_{h_i^* i}(\alpha_{h_i^*}), \min_{i \in \Theta_O} \hat{\rho}_i^r I_{0i}(c_i) \right) \Big| \leq \epsilon, \end{aligned} \quad (50)$$

where  $\epsilon$  decreases to zeros as  $\epsilon \rightarrow 0$ . In other words, Algorithm 1 achieves the almost optimal value of  $AEOC-B$  when  $r \geq r^*$ .

## 5 | NUMERICAL EXPERIMENTS

In this section, we provide numerical results that show the effectiveness of the proposed algorithm. To this end, we test

randomly generated test cases and a case study inspired by the real-world staff allocation problem introduced in Section 3.1.

## 5.1 | Random test cases

We first test the performance of Algorithm 1 (short for PB in the following context) using randomly generated test cases. Three benchmark algorithms are selected in comparison, namely, the OCBA for constrained optimization in Lee et al. (2012) (short for CO), the stochastically constrained R&S via SCORE in Pasupathy et al. (2015) (short for SC), and the heuristic procedure based on statistical hypothesis tests in Choi et al. (2021) (short for HE). Although the CO, SC, and HE algorithms do not consider the performance correlations in solutions as PB does, they are the most relevant algorithms developed to solve the constrained R&S problems under study in this paper. As such, we test the PB algorithm with an uncorrelated uniform prior (short for PB<sub>u</sub>) as a direct comparison to the three benchmarks. Further, we test the PB algorithm with a correlated prior (short for PB<sub>c</sub>) in order to quantify the benefit of using a correlated prior in the PB implementation. To this end, we generate random test cases with different parameter settings, and compute the optimal solutions using the aforesaid five algorithms by varying computing budgets. Note that the true optimum for each random test case can be analytically computed (to become apparent once we introduce the generation procedure for random cases below) and will be used as ground truths for performance comparison. Note further that our implementations of all algorithms strictly use sample approximations without leveraging the knowledge of the true optimum.

We compare the algorithm performance under different values of  $|\mathcal{X}|$ ,  $H$ , and  $T$ . Recalling that the CO and HE algorithms were developed based on the assumption that samples follow a normal underlying distribution, we will further compare the performance with different distributions, including normal, uniform, and exponential. Given each parameter setting ( $|\mathcal{X}|$ ,  $H$ ,  $T$ , and an underlying distribution), 1,000 random test cases were generated based on the following procedure:

1. Generate the value of  $d_h$  for  $h = 0, 1, \dots, H$ , each following a discrete uniform distribution  $\mathcal{DU}(1, 200)$ .
2. Generate the value of  $\alpha_h$  for  $h = 1, \dots, H$ , each following a uniform distribution  $\mathcal{U}(0.01, 0.1)$ .
3. For each solution  $\mathbf{x}_i \in \mathcal{X}$  for  $i = 1, \dots, K$ , generate the value of  $F_h(\mathbf{x}_i, \omega)$  for  $h = 0, 1, \dots, H$ , each following the chosen underlying distribution.
  - For a normal underlying distribution,  $F_h(\mathbf{x}_i, \omega) \sim \mathcal{N}(\mu_i, \sigma^2)$ , where  $\mu_i \sim \mathcal{U}((0.75 + 0.0015i) d_h, (0.85 + 0.0015i) d_h)$  and  $\sigma = 0.12 d_h$ .
  - For a uniform underlying distribution,  $F_h(\mathbf{x}_i, \omega) \sim \mathcal{U}(a_i, b_i)$ , where  $b_i - a_i \sim \mathcal{U}(0, d_h)$  and  $\frac{b_i - d_h}{b_i - a_i} \sim \mathcal{U}(0.0015i, 0.0015i + 0.1)$ .

- For an exponential underlying distribution,  $F_h(\mathbf{x}_i, \omega) \sim \exp(-\frac{\log \lambda_i}{d_h})$ , where  $\lambda_i \sim \mathcal{U}(0.0015i, 0.0015i + 0.1)$ .

For each random test case, the value of  $\mathbb{E}[F_h(\mathbf{x}_i, \omega)]$  depends on the randomly generated distribution parameters, while  $\mathbb{E}_s[\mathbb{E}[F_h(\mathbf{x}_i, \omega)]]$  increases with  $i$  where  $\mathbb{E}_s[\cdot]$  denotes the expectation over the randomness caused by the distribution parameters. As seen, the test cases were generated to mimic the correlation in solution performances in a way that the performance of a solution may be more correlated with its neighboring ones.

In the PB<sub>c</sub> implementation, we set the prior distribution of  $\log(\Psi_{hi}/(1 - \Psi_{hi}))$ ,  $i \in \Theta$  as a multivariate normal (Howard, 1998). The covariance function is squared exponential (Rasmussen & Williams, 2006); that is,  $\Sigma(\mathbf{x}_i, \mathbf{x}_{i'}) = \tau^2 e^{-\phi \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2^2}$  with  $\tau > 0$  and  $\phi > 0$ . Let  $Cov$  be the covariance matrix whose  $(i, i')$  entry is  $(Cov)_{ii'} = \Sigma(\mathbf{x}_i, \mathbf{x}_{i'})$ . Then, the prior belief has a density  $\pi_h^0(\boldsymbol{\psi}_h)$  proportional to

$$\exp\left(-\frac{1}{2}\left(\log\left(\frac{\boldsymbol{\psi}_h}{1 - \boldsymbol{\psi}_h}\right) - \mathbf{y}\right) \cdot (Cov)^{-1} \times \left(\log\left(\frac{\boldsymbol{\psi}_h}{1 - \boldsymbol{\psi}_h}\right) - \mathbf{y}\right)^\top\right) \prod_{i \in \Theta} \frac{d \log\left(\frac{\psi_{hi}}{1 - \psi_{hi}}\right)}{d\psi_{hi}}, \quad (51)$$

where  $\mathbf{y}$  is the expectation of the random vector  $\log(\boldsymbol{\Psi}_h/(1 - \boldsymbol{\Psi}_h)) = (\log(\Psi_{h1}/(1 - \Psi_{h1})), \dots, \log(\Psi_{hK}/(1 - \Psi_{hK})))$ . In the experiments, we set  $\tau = 1$ ,  $\phi = 1$ , and  $\mathbf{y} = (y_1, y_2, \dots, y_K)^\top$  to be  $y_i = \log(\frac{\mathbb{E}_s[p_{hi}]}{1 - \mathbb{E}_s[p_{hi}]})$ , where  $p_{hi} = P\{F_h(\mathbf{x}_i, \omega) > d_h\}$ . Here, the choice of  $\mathbf{y}$  mimics practical situations when decision makers may have some prior knowledge on  $p_{hi}$ , but the prior is inaccurate if  $p_{hi}$  is predicted by  $\mathbb{E}_s[p_{hi}]$ . We mention that the posterior mean does not have a closed form and is approximated using the Markov chain Monte Carlo approach.

For each random case, its true (theoretical) optimal solution for the attendant problem  $\mathcal{P}$  can be analytically computed based on the distribution information. Then, the same problem is solved using the five algorithms given a computing budget  $T$ . The solution obtained is compared to the true optimal solution. Three metrics are computed based on the 1,000 cases:

- Probability of correct selection (PCS): the percentage of cases where a true optimal solution is identified using the chosen algorithm;
- Probability of selecting an infeasible solution (PIF): the percentage of cases where an infeasible solution is identified as an optimal solution using the chosen algorithm;
- Expect opportunity cost (EOC): the expected performance gap between the selected solution and the true optimum. Specifically, the frequentist opportunity cost of selecting  $\mathbf{x}_i \neq \mathbf{x}_b$  is defined as  $OC_i^* = \sum_{h=0}^H OC_{hi}^*$ , where  $OC_{hi}^* = \max\{p_{hi} - \alpha_h, 0\}$ ,  $h = 1, 2, \dots, H$ , and

**TABLE 2** Performance comparison with normal underlying distributions ( $\alpha_h \sim \mathcal{U}(0.01, 0.1)$ )

Setting		PCS (%)					PIF (%)					EOC (×1,000)				
(K, H)	T	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>
(200, 1)	8,000	14.9	15.5	15.0	28.6	37.2	38.3	37.3	36.7	12.6	17.0	27.1	26.1	26.3	16.9	15.1
	16,000	35.2	29.1	29.7	54.2	62.7	22.6	24.0	28.2	4.2	5.8	13.7	13.8	14.9	7.8	6.2
	24,000	55.1	52.0	48.5	64.7	73.3	13.7	11.8	15.1	3.5	3.7	8.5	6.2	7.5	4.8	3.4
	32,000	64.4	68.0	66.7	74.6	77.7	11.2	7.1	6.3	2.2	2.8	6.7	3.3	4.2	3.6	2.9
	40,000	69.5	77.0	75.2	78.1	82.3	10.0	4.3	4.7	2.0	2.3	6.3	2.0	3.4	2.7	2.2
(500, 3)	20,000	27.9	19.9	20.5	38.2	42.7	46.2	56.2	57.7	14.2	20.6	30.6	36.4	34.3	28.8	26.1
	40,000	59.1	28.1	28.1	69.6	74.3	17.8	44.5	45.2	5.6	7.3	15.1	26.7	25.8	10.8	8.5
	60,000	73.0	32.2	35.0	81.3	85.0	10.6	38.8	36.5	3.3	4.1	11.8	22.2	20.9	6.4	4.0
	80,000	77.6	37.4	38.8	85.1	88.7	8.2	34.4	33.6	2.5	2.6	10.2	20.7	19.2	5.3	2.9
	100,000	80.3	42.5	42.5	86.8	90.3	7.8	29.0	29.3	2.6	2.2	9.9	18.0	17.8	4.2	2.3
(500, 1)	20,000	14.4	15.4	17.2	34.9	44.1	39.3	38.1	34.7	10.5	13.5	27.2	27.2	25.1	14.5	12.3
	40,000	40.2	23.2	24.8	63.3	74.1	20.1	30.2	27.1	5.6	3.6	11.9	16.5	15.3	5.5	3.3
	60,000	61.4	29.7	30.4	77.5	83.3	11.9	24.2	24.5	2.6	1.5	7.2	11.9	11.8	2.6	1.8
	80,000	71.9	34.5	33.3	82.9	86.7	9.9	20.2	21.6	1.9	1.4	5.9	10.0	9.9	2.1	1.4
	100,000	74.9	37.6	38.0	86.9	90.6	9.3	19.4	19.9	1.2	1.0	5.7	8.8	8.0	1.2	0.8

$OC_{0i}^* = \max\{p_{0i} - p_{0b}, 0\}$ . Let  $OC_{no}^* = 0.1$  when  $\hat{b}$  does not exist. Then,  $EOC = OC_{no}^*P(\hat{b} \text{ does not exist}) + \sum_{i \neq b} OC_i^*P(\hat{b} = i)$ , where the simulation results are used to estimate  $P(\hat{b} \text{ does not exist})$  and  $P(\hat{b} = i)$ .

Here, PCS measures an algorithm's expected performance (i.e., how good it is in selecting the optimal solution); PIF measures its worst-case performance (i.e., how bad it is in selecting an infeasible solution); and EOC measures the expected performance gap between the selected solution and the true optimum (i.e., the overall performance on the solution quality).

Table 2 exhibits the performance comparison results for the normal underlying distribution. Specifically, the table shows the PCS, PIF, and EOC metrics for three groups of settings, ( $K = 200, H = 1$ ), ( $K = 500, H = 3$ ), and ( $K = 500, H = 1$ ), where the difficulties of problems in each group increase in the aforesaid order. It is straightforward to see that group ( $K = 500, H = 1$ ) is much more challenging than group ( $K = 200, H = 1$ ), while it may be counterintuitive to see that group ( $K = 500, H = 3$ ) is also easier than ( $K = 500, H = 1$ ). The latter is due to the fact that each problem in group ( $K = 500, H = 1$ ) is expected to have more feasible solutions than that in group ( $K = 500, H = 3$ ). Since identifying an infeasible solution is typically less computationally expensive than verifying a feasible solution, group ( $K = 500, H = 1$ ) is much harder than group ( $K = 500, H = 3$ ). In summary, we characterize three groups of problems as follows: ( $K = 200, H = 1$ ) has a small size of candidate solutions; ( $K = 500, H = 3$ ) has a small size of feasible solutions and a large size of infeasible ones; and ( $K = 500, H = 1$ ) has large sizes of feasible and infeasible solutions. Note that, for all algorithms, we set  $n_0 = 20$  and  $\Delta = 4,000$  for ( $K = 200, H = 1$ ), and  $\Delta = 10,000$

for both ( $K = 500, H = 3$ ) and ( $K = 500, H = 1$ ). For each group, we further vary the values of computing budget  $T$  to show the convergence trend.

From Table 2, it is seen that, for the PCS metric,  $PB_u$  clearly outperforms three benchmarks for all three groups of problems; particularly, for the hardest (third) group,  $PB_u$  reaches a PCS of 77.5% when  $T = 60,000$ , a value that the benchmark algorithms fail to reach even when  $T = 100,000$ . For the PIF metric,  $PB_u$  shows good controls of not choosing an infeasible solution especially as  $T$  becomes sufficiently large, while the performance of the benchmarks are not as satisfactory. For the EOC metric,  $PB_u$  again outperforms benchmarks by a substantial margin for the harder problems (second and third groups), indicating a strong overall performance by taking into account the quality of the selected solution. Further,  $PB_c$  outperforms  $PB_u$  for all test cases in terms of PCS and EOC, showing the benefit of considering a correlated prior when performance correlations indeed exist. In terms of PIF,  $PB_c$  performs better in harder problems, while  $PB_u$  is superior in easy problems. Note that the advantage of  $PB_c$  over  $PB_u$  tends to diminish as  $T$  increases.

To check the robustness of PB across different underlying distributions, Table 3 and Table 4 exhibit performance comparisons with uniform and exponential underlying distributions, respectively. It is seen that PB is superior to all three benchmarks robustly across various underlying distributions. We also point out that the advantage of PB is more significant for the hardest group of cases, indicating that PB is a better choice when the size of problem becomes large.

In the Supporting Information EC.3, we modified the procedure of generating random test cases by setting the values of  $\alpha_h$  ( $h = 1, \dots, H$ ) to follow a uniform distribution  $\mathcal{U}(0.01, 0.3)$  instead of  $\mathcal{U}(0.01, 0.1)$ , capturing a wider range

**TABLE 3** Performance comparison with uniform underlying distributions ( $\alpha_h \sim \mathcal{U}(0.01, 0.1)$ )

Setting		PCS (%)					PIF (%)					EOC ( $\times 1,000$ )				
(K, H)	T	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>
(200,1)	8,000	21.6	23.1	23.9	39.9	45.7	39.8	39.0	38.8	11.0	14.9	29.8	29.5	28.7	18.3	16.4
	16,000	44.9	42.3	36.7	61.7	73.3	23.5	24.3	30.2	5.7	6.0	14.6	14.0	15.8	8.0	4.6
	24,000	60.4	61.8	55.3	73.5	82.2	13.5	13.0	17.2	2.8	2.3	10.3	5.6	7.8	5.1	2.5
	32,000	70.0	74.2	68.2	79.9	84.9	10.8	8.3	11.2	2.6	2.2	8.8	3.2	4.2	3.7	1.8
	40,000	73.2	81.0	78.0	83.5	87.5	10.3	5.5	6.4	1.6	2.2	8.4	1.8	2.7	2.8	1.5
(500,3)	20,000	31.7	24.5	25.9	45.8	50.6	45.4	53.4	52.5	12.5	17.0	36.6	38.9	36.7	31.9	26.0
	40,000	59.0	57.2	33.7	76.5	78.4	18.0	21.9	42.6	4.8	5.0	20.8	17.8	26.8	9.8	8.8
	60,000	69.9	70.6	43.5	85.3	87.6	11.1	10.4	33.1	2.7	3.3	16.7	14.1	22.8	5.6	4.8
	80,000	75.3	77.1	55.3	88.3	92.0	8.6	6.2	20.8	1.7	2.6	14.6	12.7	16.9	4.6	2.8
	100,000	77.7	80.8	64.6	90.6	93.3	7.9	5.1	14.8	1.3	1.7	13.2	11.8	14.4	3.5	2.6
(500,1)	20,000	21.5	23.9	24.2	47.8	56.6	39.3	39.5	36.8	7.2	10.7	31.1	28.8	29.8	16.1	11.8
	40,000	47.7	43.9	38.2	77.7	82.1	21.4	23.0	27.4	3.1	2.6	13.3	13.4	15.4	4.1	2.6
	60,000	68.0	65.1	60.4	86.1	90.4	12.2	9.9	15.5	1.7	0.9	7.4	5.4	7.1	2.2	0.9
	80,000	77.2	79.9	74.6	91.0	92.7	9.8	4.8	8.4	1.1	0.9	6.3	2.6	4.0	1.5	0.7
	100,000	80.0	87.5	84.8	92.6	94.7	9.8	2.7	4.0	0.4	0.3	6.3	1.8	2.5	1.3	0.6

**TABLE 4** Performance comparison with exponential underlying distributions ( $\alpha_h \sim \mathcal{U}(0.01, 0.1)$ )

Setting		PCS (%)					PIF (%)					EOC ( $\times 1,000$ )				
(K, H)	T	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>	CO	SC	HE	PB <sub>u</sub>	PB <sub>c</sub>
(200,1)	8,000	22.2	22.6	20.9	42.3	44.8	39.3	41.2	40.7	10.6	13.0	31.7	30.6	30.3	17.4	17.6
	16,000	44.6	43.2	37.7	66.5	70.8	21.9	23.3	27.9	5.2	4.0	15.5	14.0	15.1	7.1	6.1
	24,000	60.3	62.3	56.1	76.4	82.3	13.5	11.7	16.1	2.1	2.2	9.9	6.7	8.5	4.4	3.4
	32,000	69.1	74.5	66.8	82.0	85.7	10.7	6.6	9.8	1.4	2.7	8.5	3.6	6.3	3.0	2.7
	40,000	73.1	82.7	76.2	85.6	88.4	10.2	4.7	6.3	1.2	1.4	7.9	2.3	4.3	2.5	2.3
(500,3)	20,000	31.8	25.0	22.2	43.9	50.2	48.9	54.6	54.2	15.0	17.8	34.0	36.7	38.9	31.2	24.0
	40,000	62.0	55.1	34.1	76.3	79.7	17.5	22.2	40.8	4.3	5.9	16.9	17.7	26.5	8.6	7.2
	60,000	74.4	72.1	44.0	84.2	88.1	11.4	10.5	31.9	2.9	3.3	12.8	11.8	20.2	5.6	4.1
	80,000	80.4	79.0	57.2	88.0	91.7	6.8	6.3	20.4	2.3	1.9	11.3	10.5	16.0	4.2	2.8
	100,000	81.5	82.8	66.9	90.5	93.4	7.3	4.5	14.5	1.9	1.9	10.9	9.5	13.3	2.9	2.3
(500,1)	20,000	21.7	21.6	20.4	45.0	55.0	39.6	42.6	40.7	8.6	11.7	29.8	29.8	30.0	18.6	12.7
	40,000	48.3	45.1	37.3	75.1	80.6	21.9	22.4	26.9	2.8	3.6	14.6	13.1	14.5	4.9	3.1
	60,000	66.9	66.3	55.1	85.0	90.3	14.7	10.0	18.3	2.2	1.3	8.9	5.6	8.4	2.2	1.7
	80,000	73.2	80.5	72.7	88.8	91.6	12.2	4.7	8.2	1.5	1.3	8.0	3.0	5.1	1.7	1.4
	100,000	75.7	87.1	81.7	91.4	93.4	11.4	3.3	5.5	1.1	0.7	7.6	1.9	3.4	1.2	1.1

of service levels seen in applications. With this modification, the random test cases are expected to have more feasible solutions in general. The performance comparisons are shown in Tables EC.1, EC.2, and EC.3 in the Supporting Information, and the observations are quite similar to those in Tables 2, 3, and 4, showing the robustness of PB across different settings of problem  $\mathcal{P}$  by varying  $\alpha_h$ . Further, we varied the incremental budget  $\Delta$  in the PB algorithm implementation. The impact of  $\Delta$  on performance and suggestions on the

choice of  $\Delta$  are also discussed in the Supporting Information EC.3.

In summary, these numerical results show that, compared to CO, SC, and HE, PB is more likely to correctly select the optimal solution, and less likely to select an infeasible one; even if an optimal solution is not selected, PB is much more likely to select a near-optimal one. All these are desired characteristics for an R&S algorithm. Further, the performance improvement is quite significant even when the computing

budget is relatively small. This will benefit many applications as each simulation replication can be expensive and time consuming to run and the decision-making time frame may be short (e.g., within hours or half a day). Finally, when prior knowledge is available for a problem, it would be beneficial to incorporate a correlated prior.

## 5.2 | An ED case study

In this section, we will demonstrate how the proposed model and algorithm can be used to improve services in the ED example introduced in Section 3.1. A discrete-event simulation model is used to evaluate the performance of the highly complex and dynamic processes in ED. The data used in this case study were collected from July 1, 2009, to June 30, 2010, and the service process and simulation model used have been verified in Guo et al. (2017) and Chen et al. (2020).

From the historical data obtained, the ED received various patients with the volume varying from 350 to 430 per day, where the patients fall into five categories with the following percentages on average (from categories I to V): 0.87%, 1.45%, 42.09%, 48.76%, and 6.83%. In the simulation model, the patient arrivals follow a nonhomogeneous Poisson process, whose hourly arrival rates were calculated empirically based on the historical data and are shown in the Supporting Information EC.4.

The discrete event simulation model follows the actual service processes (Chen et al., 2020), where the detailed process map and the distributions used to generate the service times at each process are displayed in the Supporting Information EC.4. Specifically, an incoming patient, depending on the situation, may go through different service areas, including registration and triage, resuscitation (for critical patients), consultation, laboratory test and imaging investigation, and observation. Main medical resources required include different types of nurses, doctors with different seniority levels and expertise, different medical devices (e.g., electrocardiography [ECG] and ultrasonography [USG]), as well as different spaces (offices, rooms, and labs). Medical staff worked three shifts around the clock: morning shift (from 8 a.m. to 4 p.m.), evening shift (from 4 p.m. to 12 p.m.), and night shift (from 0 a.m. to 8 a.m.).

The simulation model was implemented in Arena, and was validated using the 1-year data collected. For the model validation, we set the length of the simulation to be 1 month with a 7-day warm-up period and ran the simulation model for 500 replications. We used three indicators, waiting time at the triage (WTT), waiting time for physicians (WTP), and length of stay (LOS), to assess the validity of outputs, as shown in Table 5. Note that the data of WTT and WTP for category I patients are not available from the cooperative hospital. It is seen that the simulation can closely mimic the real data.

Next, we experiment with the simulation model to identify the optimal solution  $\mathbf{x}^*$  that optimizes the resource allocation problem defined in Equations (1)–(5). Recall that each candidate solution  $\mathbf{x}$  consists of the medical staff assignment for

TABLE 5 Simulation model validation results

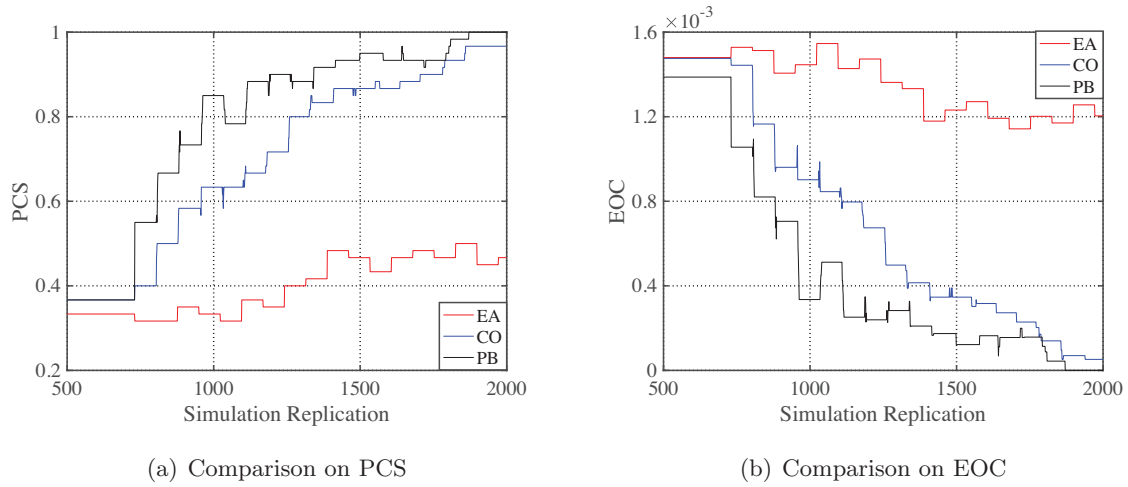
		CAT I	CAT II	CAT III	CAT IV	CAT V
WTT (h)	Real data	–	0.04	0.09	0.11	0.11
	Simulated result	–	0.03	0.11	0.11	0.11
WTP (h)	Real data	–	0.11	0.31	1.65	2.07
	Simulated result	–	0.07	0.27	1.64	2.17
LOS (h)	Real data	1.30	1.23	3.53	4.15	3.13
	Simulated result	1.56	1.51	3.49	4.37	3.20

three 8-h shifts of a day. The medical staff is composed of the admission nurse (AN), triage nurse (TN), junior nurse (JN), senior nurse (SN), junior doctor (JD), and senior doctor (SD). To protect the data privacy, we set the labor cost for each type of staff using the cost unit (CU) as follows: 1 CU per AN, 1 CU per TN, 1.5 CUs per JN, 2 CUs per SN, 3 CUs per JD, and 4 CUs per SD. For the actual solution implemented in the ED, the total staffing cost is 85.5 CUs. By applying the following filters to all possible solutions, we identified 73 promising candidate solutions for the R&S problem ( $|\mathcal{X}| = 73$ ).

- Solutions were filtered by the deterministic constraints specified in Equation (5), including the personnel and budget constraints. Specifically, the personnel constraint imposes practical bounds on the number of medical staff of each type and shift. Any solution resulting in a staffing cost of more than 85.5 CUs was also filtered, since the ED preferred maintaining the current budget.
- Solutions dominated by others were filtered. For example, consider solution A with 2 SDs for the night shift and solution B with 1 SD for the night shift, while all other staff assignments are equal. We expect solution A to dominate solution B without evaluating them since the objective is to find a solution that maximizes the service level of type IV patients, thereby excluding solution B from the solution space.
- Solutions were further filtered by consulting domain experts, who would remove nonpromising solutions based on their prior experience on staff schedules or preliminary simulation runs.

The full set of the 73 candidate solutions is included in the Supporting Information EC.4. We mention that the above filtering process substantially reduces the solution space, making it suitable to construct an R&S problem studied in this paper, while the original problem without filtering is more suitable to be solved using a search-based algorithm, such as the one developed in Chen et al. (2020).

We compare the performance of PB to two benchmarks: One is the CO in Lee et al. (2012) and the other is the equal allocation (EA). The EA approach is a simple strategy of allocating the computing budget equally among all candidate solutions, and is often used in the literature as a baseline. The computations were conducted in Arena 14.0 via VBA codes. The prior belief was set to have a density



**FIGURE 1** Performance comparison on the ED example [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**TABLE 6** Comparison between the actual and optimal solutions

	Shift	Staff allocation						Service level				
		AN	TN	JN	SN	JD	SD	$P_{WT_2}$	$P_{WT_3}$	$P_{LOS}$	$P_{WT_4}$	Cost
Actual solution	Morning	1	2	2	4	2	3	93.79%	89.10%	88.31%	67.59%	85.5
	Evening	1	2	2	4	2	3					
	Night	1	1	1	2	2	2					
Optimal solution	Morning	2	2	3	2	3	3	97.08%	91.41%	98.46%	99.39%	84.5
	Evening	1	2	2	2	2	3					
	Night	1	1	2	2	2	2					

$\pi_h^0(\psi_h)$  proportional to Equation (51), where  $\tau = 100$ ,  $\phi = 1$ , and  $\mathbf{y}$  is a vector of zeros. We set the total computing budget  $T = 2,117$ , and therefore each solution was simulated for 29 replications using EA. For PB and CO, we ran  $n_0 = 10$  initial replications for each solution, and set the incremental budget  $\Delta = 73$ . Note that among the set of solution candidates, only one of them has been implemented in practice and thus no ground-truth values are known for these solutions. To obtain approximated ground-truth values, we simulated each solution for 500 replications (sufficiently large) and selected the best solution as the estimated optimal solution. Then, we ran each algorithm to select the best solution and compared it to this estimated optimal solution. For each algorithm, we ran the same experiment for 60 times, and Figure 1 plots the average values of PCS and EOC evolving as simulation replications increase.

It is seen from Figure 1 that the PCS converges to 1 and EOC converges to 0 faster when using PB than using CO, while the PCS and EOC do not converge within the tested range of the computing budget when using EA. EA is apparently not a good strategy, and needs a much larger computing budget to conclude a reliable solution. To further illustrate that the optimal solution obtained by PB is indeed better than the actual staff allocation schedule, we simulated both solu-

tions for 500 replications. Table 6 shows the performance comparison between the actual solution implemented by the hospital and the optimal solution obtained by PB, where  $P_{WT_2}$ ,  $P_{WT_3}$ ,  $P_{LOS}$ , and  $P_{WT_4}$  are the probabilistic measures defined on the left-hand side of Equations (2), (3), (4), and (1), respectively. It is seen that the optimized solution meets all probabilistic service levels specified in Table 1, while the actual solution fails to do so (with a small margin). Further, the optimal solution performs much better in terms of the objective (1) (i.e.,  $P_{WT_4}$ ) compared to the actual solution, with a 31.80% of improvement on this measure. The total staff costs for both solutions are very close and within the specified budget. The comparison shows the superiority of the optimized solution, and confirms the necessity of optimizing the resource allocation for this ED.

## 6 | CONCLUSION

The purpose of this research is to formulate a model to study resource allocation problems in service systems, where the system performance is evaluated via simulation. Particularly, the objective and constraints in the class of problems under study are in the form of probabilistic measures.

Simulation models are known to be capable of capturing complex processes and system dynamics, which are hard to model in mathematical programs but are often present in service systems. However, a major challenge in optimizing systems using simulation models, known as the research field of simulation optimization or optimization via simulation, is the trade-off between simulation accuracy and the affordable computing budget. To this end, we formulate an OCBA model to efficiently balance the effort of exploring the solution space and exploiting given solutions with more simulation replications. From the methodological point of view, the OCBA problem formulated in this paper possesses a couple of unique features compared to the literature: (1) Instead of maximizing the PCS, the problem under study aims to minimize the EOC, considering that the service provider is typically more concerned with the quality of the selected solution than the statistical significance of selecting the true optimal one; and (2) the proposed model considers the prior knowledge accumulated on solutions, as well as the performance correlations on solutions, and hence uses a Bayesian modeling approach. To the best of our knowledge, such a Bayesian OCBA model has not been studied in the literature. Next, based on the property of the probabilistic measures, we derive the asymptotic optimality conditions of the OCBA formulation. Then, an efficient iterative algorithm is developed to sequentially allocate the simulation budget to the most necessary solutions, guided by the optimality conditions. The efficiency of the algorithm is shown via the theoretical finite-time convergence analysis, as well as the numerical experiments.

From a practical vintage point, the generalized resource allocation problem is motivated by a real-life resource allocation problem stemming from a hospital ED, and is widely applicable to service systems, such as hospitals, call centers, and public service offices. The proposed algorithm can expedite the identification of an optimal resource allocation scenario using a simulation model, with the presence of probabilistic constraints and objective. Extensive numerical examples with known true optimum conducted in this paper have demonstrated the superiority of the proposed algorithm compared to benchmarks, evaluated by the PCS, PIF, and EOC metrics. Further, we apply the algorithms to the aforesaid staff allocation problem in the ED in Hong Kong, whose simulation model is driven by real data. This case study has further validated the practicality of the proposed model and algorithm.

Several future research directions are available to explore. Sometimes, the service system may want to guarantee the PCS (e.g., 90%) and minimize the required computing budget (Chen et al., 2014). Thus, a future research direction is to explore alternative simulation optimization methods, such as the IZ approaches. Further, this research considers problems with a manageable size of solution space. In the proposed algorithm, every solution is simulated for a minimum number of replications and evaluated to obtain an initial assessment. For a large-size problem, it may not

be computationally feasible to provide an initial evaluation for every solution. As such, a metamodel-based or random search-based method needs to be developed to balance the exploration of unknown solutions and the exploitation of known solutions with more simulation replications. Finally, when the resource constraints in Equations (1)–(5) or Equations (6)–(8) can be expressed as mathematical constraints, it would be interesting to explore a combined simulation optimization and mathematical programming approach.

## ACKNOWLEDGMENTS

We gratefully thank the department editor, the senior editor, and the two reviewers for their valuable comments that have greatly improved the paper. Siyang Gao thanks the support by the City University of Hong Kong (Grants 7005269 and 7005568). Jianzhong Du thanks the support by the National Natural Science Foundation of China (Grant 72091211).

## ORCID

Weiwei Chen  <https://orcid.org/0000-0002-7736-3411>

## REFERENCES

- Andradóttir, S., & Kim, S.-H. (2010). Fully sequential procedures for comparing constrained systems via simulation. *Naval Research Logistics*, 57(5), 403–421.
- Atlason, J., Epelman, M. A., & Henderson, S. G. (2008). Optimizing call center staffing using simulation and analytic center cutting-plane methods. *Management Science*, 54(2), 295–309.
- Batur, D., & Kim, S.-H. (2010). Finding feasible systems in the presence of constraints on multiple performance measures. *ACM Transactions on Modeling and Computer Simulation*, 20(3), 1–26.
- Beraldi, P., Bruni, M. E., & Conforti, D. (2004). Designing robust emergency medical service via stochastic programming. *European Journal of Operational Research*, 158(1), 183–193.
- Best, T. J., Sandıkçı, B., Eisenstein, D. D., & Meltzer, D. O. (2015). Managing hospital inpatient bed capacity through partitioning care into focused wings. *Manufacturing & Service Operations Management*, 17(2), 157–176.
- Bodur, M., & Luedtke, J. R. (2017). Mixed-integer rounding enhanced benders decomposition for multiclass service-system staffing and scheduling with arrival rate uncertainty. *Management Science*, 63(7), 2073–2091.
- Buckley, P., & Majumdar, R. (2018). *The services powerhouse: Increasingly vital to world economic growth*. Tech. rep., Deloitte Insights.
- Cezik, M. T., & L'Ecuyer, P. (2008). Staffing multiskill call centers via linear programming and simulation. *Management Science*, 54(2), 310–323.
- Chen, C.-H., Lin, J., Yücesan, E., & Chick, S. E. (2000). Simulation budget allocation for further enhancing the efficiency of ordinal optimization. *Discrete Event Dynamic Systems*, 10, 251–270.
- Chen, W., Gao, S., Chen, C.-H., & Shi, L. (2014). An optimal sample allocation strategy for partition-based random search. *IEEE Transactions on Automation Science and Engineering*, 11(1), 177–186.
- Chen, W., Guo, H., & Tsui, K.-L. (2020). A new medical staff allocation via simulation optimisation for an emergency department in Hong Kong. *International Journal of Production Research*, 58(19), 6004–6023.
- Chick, S. E., Branke, J., & Schmidt, C. (2010). Sequential sampling to myopically maximize the expected value of information. *INFORMS Journal on Computing*, 22(1), 71–80.
- Chick, S. E., & Wu, Y. (2005). Selection procedures with frequentist expected opportunity cost bounds. *Operations Research*, 53(5), 867–878.
- Choi, S. H., Kang, B. G., & Kim, T. G. (2021). A heuristic procedure for selecting the best feasible design in the presence of stochastic constraints. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 51(2), 1016–1026.

- Dmitry, I., Sokolov, B., Chen, W., Dolgui, A., Werner, F., & Potryasaev, S. (2021). A control approach to scheduling flexibly configurable jobs with dynamic structural-logical constraints. *IIE Transactions*, 53(1), 21–38.
- Frazier, P., Powell, W., & Dayanik, S. (2009). The knowledge-gradient policy for correlated normal beliefs. *INFORMS Journal on Computing*, 21(4), 599–613.
- Ganesh, A. J., O'Connell, N., & Wischik, D. J. (2004). *Big queues*, vol. 1838 of *lecture notes in mathematics*. Springer.
- Gans, N., Shen, H., Zhou, Y.-P., Korolev, N., McCord, A., & Ristock, H. (2015). Parametric forecasting and stochastic programming models for call-center workforce scheduling. *Manufacturing & Service Operations Management*, 17(4), 571–588.
- Gao, F., Gao, S., Xiao, H., & Shi, Z. (2019). Advancing constrained ranking and selection with regression in partitioned domains. *IEEE Transactions on Automation Science and Engineering*, 16(1), 382–391.
- Gao, S., & Chen, W. (2016). A partition-based random search for stochastic constrained optimization via simulation. *IEEE Transactions on Automatic Control*, 62(2), 740–752.
- Gao, S., Chen, W., & Shi, L. (2017). A new budget allocation framework for the expected opportunity cost. *Operations Research*, 65, 787–803.
- Gardner, J., Kusner, M., Xu, Z., Weinberger, K., & Cunningham, J. (2014). Bayesian optimization with inequality constraints. In Eric P. Xing & Tony Jebara (Eds.), *Proceedings of the 31st international conference on machine learning (ICML)* (Vol. 2014, pp. 937–945).
- Guo, H., Gao, S., Tsui, K.-L., & Niu, T. (2017). Simulation optimization for medical staff configuration at emergency department in Hong Kong. *IEEE Transactions on Automation Science and Engineering*, 14(4), 1655–1665.
- Healey, C., Andradóttir, S., & Kim, S.-H. (2014). Selection procedures for simulations with multiple constraints under independent and correlated sampling. *ACM Transactions on Modeling and Computer Simulation*, 24(3), 1–25.
- Hong, L. J., Luo, J., & Nelson, B. L. (2015). Chance constrained selection of the best. *INFORMS Journal on Computing*, 27(5), 317–334.
- Howard, J. V. (1998). The 2x2 table: A discussion from a Bayesian viewpoint. *Statistical Science*, 13(4), 351–367.
- Huh, W. T., Liu, N., & Truong, V.-A. (2013). Multiresource allocation scheduling in dynamic environments. *Manufacturing & Service Operations Management*, 15(2), 280–291.
- Hunter, S. R., & Pasupathy, R. (2013). Optimal sampling laws for stochastically constrained simulation optimization on finite sets. *INFORMS Journal on Computing*, 25(3), 527–542.
- Izady, N., & Worthington, D. (2012). Setting staffing requirements for time dependent queueing networks: The case of accident and emergency departments. *European Journal of Operational Research*, 219(3), 531–540.
- Kim, K., & Mehrotra, S. (2015). A two-stage stochastic integer programming approach to integrated staffing and scheduling with application to nurse management. *Operations Research*, 63(6), 1431–1451.
- Lee, L. H., Pujowidianto, N. A., Li, L.-W., Chen, C.-H., & Yap, C. M. (2012). Approximate simulation budget allocation for selecting the best design in the presence of stochastic constraints. *IEEE Transactions on Automatic Control*, 57(11), 2940–2945.
- Letham, B., Karrer, B., Ottoni, G., & Bakshy, E. (2019). Constrained Bayesian optimization with noisy experiments. *Bayesian Analysis*, 14(2), 495–519.
- Luo, Y., & Lim, E. (2013). Simulation-based optimization over discrete sets with noisy constraints. *IIE Transactions*, 45(7), 699–715.
- Mason, A. J., Ryan, D. M., & Panton, D. M. (1998). Integrated simulation, heuristic and optimisation approaches to staff scheduling. *Operations Research*, 46(2), 161–175.
- Mattia, S., Rossi, F., Servilio, M., & Smriglio, S. (2017). Staffing and scheduling flexible call centers by two-stage robust optimization. *Omega*, 72, 25–37.
- Milner, J. M., & Olsen, T. L. (2008). Service-level agreements in call centers: Perils and prescriptions. *Management Science*, 54(2), 238–252.
- Nagaraj, K., & Pasupathy, R. (2013). R-spline for local integer-ordered simulation optimization problems with stochastic constraints. In R. Pasupathy, S.-H. Kim, A. Tolc, R. Hill, & M. E. Kuhl (Eds.), *Proceedings of the 2013 winter simulation conference* (pp. 846–855).
- Örsdemir, A., Deshpande, V., & Parlaktürk, A. K. (2019). Is servicization a win-win strategy? Profitability and environmental implications of servicization. *Manufacturing & Service Operations Management*, 21(3), 674–691.
- Park, C., & Kim, S.-H. (2015). Penalty function with memory for discrete optimization via simulation with stochastic constraints. *Operations Research*, 63(5), 1195–1212.
- Pasupathy, R., Hunter, S. R., Pujowidianto, N. A., Lee, L. H., & Chen, C.-H. (2015). Stochastically constrained ranking and selection via SCORE. *ACM Transactions on Modeling and Computer Simulation*, 25(1), 1–26.
- Pujowidianto, N. A., Lee, L. H., & Chen, C.-H. (2013). Minimizing opportunity cost in selecting the best feasible design. In R. Pasupathy, S.-H. Kim, A. Tolc, R. Hill, & M. E. Kuhl (Eds.), *Proceedings of the 2013 winter simulations conference* (pp. 898–907).
- Rasmussen, C. E., & Williams, C. K. (2006). *Gaussian process for machine learning*. MIT Press.
- Robbins, T. R., & Harrison, T. P. (2010). A stochastic programming model for scheduling call centers with global service level agreements. *European Journal of Operational Research*, 207(3), 1608–1619.
- Russo, D. (2020). Simple Bayesian algorithms for best-arm identification. *Operations Research*, 68(6), 1625–1647.
- Taigel, F., Meller, J., & Rothkopf, A. (2018). Data-driven capacity management with machine learning: A novel approach and a case-study for a public service office. In Hui Yang & Robin Qiu (Eds.), *Advances in service science: Proceedings of the 2018 INFORMS international conference on service science* (pp. 105–115).
- Ton, Z., & Raman, A. (2010). The effect of product variety and inventory levels on retail store sales: A longitudinal study. *Production and Operations Management*, 19(5), 546–560.
- Tsai, S. C., & Zheng, Y.-X. (2013). A simulation optimization approach for a two-echelon inventory system with service level constraints. *European Journal of Operational Research*, 229(2), 364–374.
- Ungredda, J., & Branke, J. (2021). Bayesian optimisation for constrained problems. arXiv:2105.13245v1 [cs.LG].

## SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Chen, W., Gao, S., Chen, W., & Du, J. (2023). Optimizing resource allocation in service systems via simulation: A Bayesian formulation. *Production and Operations Management*, 32, 65–81.  
<https://doi.org/10.1111/poms.13825>