

Estimating Cycle Time Percentile Curves for Manufacturing Systems via Simulation

Feng Yang

Industrial and Management Systems Engineering Department
West Virginia University

Bruce E. Ankenman

Barry L. Nelson

Department of Industrial Engineering and Management Sciences
Northwestern University
Evanston, IL 60208-3119, U.S.A.

January 8, 2008

Abstract

Cycle time-throughput (CT-TH) percentile curves quantify the relationship between percentiles of cycle time and factory throughput, and they can play an important role in strategic planning for manufacturing systems. In this paper, a highly flexible distribution, the generalized gamma, is used to represent the underlying distribution of cycle time. To obtain CT-TH percentile curves, we use a factory simulation to fit metamodels for the first three CT-TH moment curves throughout the throughput range of interest, determine the parameters of the generalized gamma by matching moments, and obtain any percentile of interest by inverting the distribution. To insure efficiency and control estimation error, simulation experiments are built up sequentially using a multistage procedure. Numerical results are presented to demonstrate the effectiveness of the approach.

Key words: discrete event simulation; response surface modeling; design of experiments; semiconductor manufacturing; queueing

1 Introduction

Planning for manufacturing, either at the factory or enterprise level, requires answering “what if” questions involving (perhaps a very large number of) different scenarios for product mix, production targets, and capital expansion. Computer simulation is an essential tool

for the design and analysis of complex manufacturing systems. Often, before a new system is deployed or changes are made to an existing system, a simulation model will be created to predict the system's performance. Even when no substantial changes are envisioned, simulation is used to allocate capacity among production facilities. In either case, simulation is faster and much more cost effective than experimenting with the physical system (when that is even possible). This is especially true in the semiconductor industry, which is the motivating application for this research (see, for instance Schömig and Fowler 2000).

Simulation is popular because it can incorporate any details that are important, and the now common practice of animating simulations means that they have a face validity that a system of equations can never hope to achieve. However, simulation can be a clumsy tool for planning: Simulation can only evaluate one scenario at a time and, depending on the complexity of the system and the details of the simulation model, it may take several minutes or even hours to complete a simulation run for each scenario. This can lead to fewer questions being asked and fewer options being considered especially when scenarios are discussed and debated in real time.

To make simulation an effective tool for planning, our approach is to use simulation to parameterize sophisticated response surface models (RSMs) that are easily explored or optimized with respect to the controllable decision variables. Although it might take significant simulation time to build the RSM, once it has been generated, the RSM is able to instantly answer what-if questions in real time (e.g., at a quarterly production planning meeting that only lasts an hour). Analytically tractable queueing models or approximations can also produce such surfaces, but they invariably require significant simplification of the actual manufacturing system. Our RSMs, generated by simulation, excel in the sense that they are relatively simple formulas like those provided by queueing models, but are fitted to a high-fidelity simulation.

In this paper, we propose a simulation-based methodology to quantify the relationship between percentiles of steady-state cycle time (CT) and throughput (TH). Cycle time is defined as a random variable representing the time required for a job or lot to traverse a given routing in a production system (Hopp and Spearman 2001). A planner can control cycle time by controlling the rate at which lots are started in the factory (lot-start rate or equivalently, throughput rate). A CT-TH percentile curve is simply a given percentile of the cycle time distribution as a function of the throughput, and it can be used to discuss the tradeoffs of lead time vs. throughput. For instance, if the 95th percentile CT-TH curve is used to set the throughput level, then 95% of the time the actual lead time of any given product will meet the promised delivery time. Hence, such CT-TH percentile curves can

play an important role in strategic planning for manufacturing. They may be used to answer questions like what throughput would this system be able to sustain if the lead times were quoted to be four weeks for a particular product? How much additional throughput could be generated if the lead time was quoted at six weeks? And do we have sufficient production capacity to satisfy customer demands and how should we distribute the production among facilities? The curves are not, however, designed for making detailed order-release decisions.

We perform a sequence of simulation experiments to characterize the cycle-time distribution as a function of the throughput. The goal is to provide a methodology that requires nothing of the user beyond 1) the simulation model; 2) a throughput range of interest, say $[x_L, x_U]$ (the throughput has been rescaled so that $0 < x_L < x_U < 1$, where 1 is the factory capacity); 3) a percentile range of interest, say $[\alpha_L, \alpha_U]$ ($0 < \alpha_L < \alpha_U < 1$, where 1 corresponds to 100%) and 4) a measure of the required precision for the estimated curves. The result is a complete response profile that quantifies the relationship of percentiles of cycle time to throughput rate.

In our procedure, the precision of the percentile estimates is selected by the user and is expressed as a relative error (e.g., 5% or 10%). Here “precision” only refers to the estimated percentiles of the simulated cycle time. The validity of the simulation model itself, although of great importance, is beyond the scope of this research. We assume that the company is satisfied that the simulation model is sufficiently detailed to provide useful information about the behavior of the manufacturing system in question. Once the CT-TH percentile curves are constructed, they allow strategists to instantly see the limits imposed on throughput rate with decreasing lead-time requirements.

The remainder of this paper is organized as follows. Section 2 provides an overview of our approach. Section 3 describes how we simultaneously estimate the first three moment curves of cycle time. In Section 4, the properties of the Generalized Gamma distribution are provided in detail. Section 5 discusses the estimation of percentiles and the statistical inference made on the estimators. Section 6 describes the experiment design used to carry out the sequential simulation experiments and gives a comprehensive presentation of the multistage procedure we have developed. The procedure is evaluated in Section 7 based on some queueing models and a full factory simulation.

2 Overview

In this section, we provide an overview of the methodology we propose to generate CT-TH percentile curves.

2.1 Distribution of Cycle Time

Our focus is on *cycle time*, in the sense used by Hopp and Spearman (2001, p. 321), “as a random variable that gives the time an individual job takes to traverse a routing.” However, our objective in this paper is to go beyond the standard summary measure of average cycle time (which we addressed in Yang et al. 2007) and consider additional summary measures of the distribution of cycle time, in particular percentiles of cycle time, as a function of throughput. Throughput is the rate (e.g., number of jobs per week) that jobs are completed, which is the same as the release rate of new jobs into the system over the long run if new jobs are released at a constant rate and the system itself is unchanging. Thus, we consider the throughput to be an independent variable that can be controlled by setting the release rate.

To be precise, let CT_h be the cycle time, as defined above, of the h th product or job completed. We assume that as $h \rightarrow \infty$, CT_h converges weakly to a random variable CT whose distribution is independent of h (see Billingsley 1999 for a definition of weak convergence) and has finite first four moments. The distribution of CT clearly depends on the throughput x , and we assume convergence of $CT_h(x)$ to $CT(x)$ for all $x \in (0, 1)$, where we have normalized throughput so that 1 corresponds to the capacity of the system. In fact we actually require a bit more: We also assume that the sample estimate of the ν^{th} ($\nu = 1, 2, 3$) moment $H(x)^{-1} \sum_{h=1}^{H(x)} (CT_h(x))^\nu$ (where $H(x)$ is the selected number of jobs simulated in steady state for simulations at x) is strongly consistent as $H(x) \rightarrow \infty$, which requires certain mild regularity conditions on the dependence structure of the cycle-time output process to insure that it satisfies a strong law of large numbers (e.g., Glynn and Iglehart 1986 give conditions for regenerative processes, Chien, Goldsman and Melamed 1997 provide conditions based on mixing, and Meyn and Tweedie 1993, Chapter 17 provide conditions for general state-space Markov chains). We are interested in percentiles of $CT(x)$ as a function of x .

Of course, such convergence never occurs in a physical manufacturing system, but for planning and analysis purposes we often approximate the finite-time behavior of a stochastic system by the limiting behavior of a stationary stochastic model of it (e.g., Buzacott and Shanthikumar 1993). When the model is a mathematically tractable (or readily approximated) queueing network model, then the conditions that insure the existence of a “steady state” can often be verified. However, if the model is a discrete-event stochastic simulation, as it is in this paper, then we can argue for the existence of a steady state by analogy to more precisely specified stochastic processes (see, for instance, Henderson 2001), but rarely can we formally prove it. At a practical level, we are assuming that if the driving stochastic processes

(job arrivals, machine processing times, failure and repair processes) are stationary, and the system logic (job priorities, queue disciplines, workcenter capacities) is unchanging, then a conceptually infinitely long simulation run will yield cycle times that satisfy our assumptions. See, for instance, Law and Kelton (2000) for more on the “steady-state simulation problem”. From here on when we refer to “cycle time” we are referring to the random variable $CT(x)$.

2.2 Overview of the Method

Simulation is often used to provide percentile estimates, and substantial research effort has been devoted to the estimation of cycle time percentiles via simulation. However, efficiently generating cycle time percentile estimates remains a challenging topic for at least two reasons: Standard estimators based on order statistics may require excessive data storage unless all of the percentiles of interest are known in advance, and even then it is difficult to do sequential estimation until a fixed precision is reached (Chen and Kelton 1999). On the other hand, approximations based on only the first two moments of cycle time and assuming a normal distribution can be grossly inaccurate (McNeill et al. 2003). A technique based on the Cornish-Fisher expansion has been proposed by McNeill et al. (2003) to estimate percentiles of cycle time; it takes into account the first four moments of the cycle time distribution and allows accurate and precise percentile estimates to be generated for moderately non-normal distributions. However, this method can only give percentiles at fixed, prespecified throughputs where simulation experiments have been performed. The methodology proposed in this paper aims at providing a more comprehensive profile of the system by generating CT-TH percentile curves throughout a throughput range.

Our approach to approximating percentiles of $CT(x)$ is to fit curves to the first three moments (equivalently mean, variance and skewness) of $CT(x)$ as a function of throughput x , match a highly flexible distribution, the Generalized Gamma distribution (GGD), to these moments, and then to invert the fitted distribution to obtain percentiles. More specifically, the strategy we propose for estimating $\mathcal{C}_\alpha(x)$, the 100α ($\alpha \in [\alpha_L, \alpha_U]$) percentile of cycle time at throughput rate $x \in [x_L, x_U]$, is outlined as follows:

1. Use an extended version of the methodology of Yang et al. (2006) to estimate not only the CT-TH mean (1^{st} moment) curve, but also the CT-TH 2^{nd} and 3^{rd} moment curves over the throughput range of interest. This allows for the prediction of the first three moments of cycle time at any throughput x , say $\mu_1(x)$, $\mu_2(x)$, and $\mu_3(x)$.
2. Use the method of moments to fit a GGD distribution $G(t; a(x), b(x), k(x))$ as an

approximation for the cycle time distribution ($a(x)$, $b(x)$, and $k(x)$ are distribution parameters that depend on x). We write the resulting fitted GGD as $G(t; \hat{a}(x), \hat{b}(x), \hat{k}(x))$.

3. Estimate the percentile $\mathcal{C}_\alpha(x)$ by taking the inverse of the c.d.f. (cumulative distribution function) of the cycle time: $\hat{\mathcal{C}}_\alpha(x) = G^{-1}(\alpha; \hat{a}(x), \hat{b}(x), \hat{k}(x))$.

The functional form we have chosen as a metamodel for the moment curves (see Section 3) was motivated by a combination of the known moment curves of some simple, single-queue models (e.g., M/M/1), and heavy traffic results for more general models (including networks of queues, where a single bottleneck queue dominates in heavy traffic; see for instance Whitt 1989). Of course, a complex manufacturing system behaves neither like a single queue nor (typically) like a queueing network in extremely heavy traffic. Therefore, our metamodel has more free parameters than these simple models, providing greater flexibility. When considering the first moment, Yang et al. (2006) showed that this model worked remarkably well. In fact, there have been a number of papers in which queueing systems have been well represented by metamodels, including Cheng and Kleijnen (1999), Fowler et al. (2001) and Park et al. (2002). However, Allen (2003) and Johnson et al. (2004) demonstrated that the models used in these papers can be inadequate for complex manufacturing systems, which motivated our more flexible formulation.

The ubiquitous use of the normal distribution in statistics might tempt one to try to get by with a two-moment approximation for the distribution of CT . However, even very simple models (such as the M/M/1 queue, where the steady-state cycle time is exponentially distributed) demonstrate that this will be woefully inadequate. On the other hand, we might consider using four or more moments as in McNeill et al. (2003), based on the premise that more moments provide a better characterization of the distribution. Our choice of three moments is a compromise between the obvious need for more than two moments, and the practical difficulty of precisely estimating curves for higher moments.

As already described, we adopt an indirect method to derive cycle-time percentiles from the moment estimates. But why not just run simulations at a very fine grid of x values and save the results? Because this approach not only requires the storage of a very large amount of data, it could also require many hours of simulation to cover a fine grid, while our approach will simulate no more than five values of throughput x , yet still deliver any percentile at any x nearly instantly, once the simulations are complete.

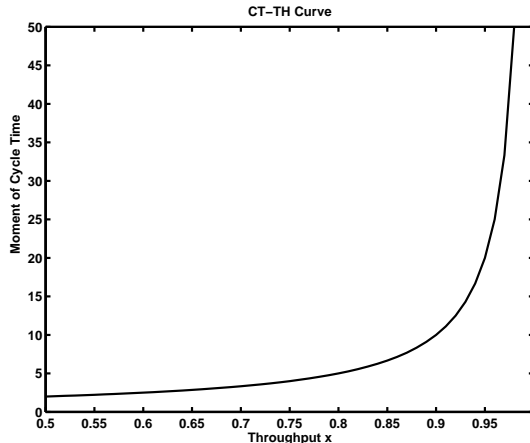


Figure 1: A generic moment CT-TH curve

3 Estimation of CT-TH Moment Curves

As indicated in the previous section, providing the first three moment CT-TH curves over the throughput range of interest is the primary step in the estimation of $\mathcal{C}_\alpha(x)$.

In Yang et al. (2006), a metamodeling-based methodology was developed for estimating mean (1^{st} moment $\mu_1(x)$) CT-TH curves, which quantifies the relationship of long-run average cycle time to throughput rate. A nonlinear regression metamodel is developed to represent the underlying CT-TH curve, and simulation experiments are sequentially built up in an efficient manner to collect data for fitting the curve. In this section, we will generalize the method of metamodeling to simultaneously estimate the first three moment CT-TH curves.

In manufacturing simulations, moment CT-TH curves typically follow the shape in Figure 1 (see, for instance, Fowler et al. (2001), Park et al. (2002), Allen (2003) and Johnson et al. (2004)). The methodology used to fit the mean CT-TH curve (Yang et al. 2006) can be extended to the estimation of higher moment curves, so in this subsection we will restate the estimation method in Yang et al. (2006) in a general way to cover estimating the first three moment curves simultaneously over a given throughput range $[x_L, x_U]$ based on a single set of simulation experiments.

We suppose that the experiment is made up of a number of independent simulation replications performed at m distinct levels of throughput $\mathbf{x} = (x_1, x_2, \dots, x_m)$ with $x_i \in [x_L, x_U]$ for $i = 1, 2, \dots, m$. From the j^{th} replication performed at throughput x , an output response $\{Y_j^{(\nu)}(x), \nu = 1, 2, 3\}$, can be obtained for the purpose of estimating the ν^{th} moment

curve:

$$Y_j^{(\nu)}(x) = \frac{1}{H(x)} \sum_{h=1}^{H(x)} (CT_{jh}(x))^\nu \quad j = 1, 2, \dots, n(x). \quad (1)$$

Here $n(x)$ is the number of replications placed at throughput x ; $CT_{jh}(x)$ represents the individual cycle time of the h^{th} job completed in the j^{th} replication at x ; and $H(x)$ is the selected number of products simulated in steady state for simulations at x . The lower bound on the value of $H(x)$ could be determined following the guidelines given in Law and Kelton (2000). As explained in Yang et al. (2007), for simplicity $H(x)$ could be set as $H(x) = H$ for all values of x (if $H(x)$ varies with x , then a simple additional step needs to be taken in the moment curve fitting). Obviously, for a given experiment consisting of a number of simulation replications carried out at m design points, the data sets

$$\mathbf{Y}^{(\nu)} = (Y_1^{(\nu)}(x_1), \dots, Y_{n(x_1)}^{(\nu)}(x_1), \dots, Y_1^{(\nu)}(x_m), \dots, Y_{n(x_m)}^{(\nu)}(x_m)) \quad (2)$$

can be extracted for $\nu = 1, 2, 3$. The integer vector $\mathbf{n} = (n(x_1), n(x_2), \dots, n(x_m))$ represents the allocation of replications to the m design points.

To these data sets $\{\mathbf{Y}^{(\nu)}, \nu = 1, 2, 3\}$, the three moment curves can be fitted. The curve fitting is based on the regression models that will be introduced below, and justification for the specific form of the models is given in Appendix A.1. The metamodeling methodology applies to fitting moment curves of any order, and for the sake of clarity, we omit the superscript ν representing the ν^{th} moment in the regression models that appear in the remainder of this section.

The CT-TH relationship for the ν^{th} moment curve can be represented by the following metamodel:

$$Y_j(x) = \mu(x, \mathbf{c}, p) + \varepsilon_j(x) \quad j = 1, 2, \dots, n(x) \quad (3)$$

where

$$\mu(x, \mathbf{c}, p) = \frac{\sum_{\ell=0}^t c_\ell x^\ell}{(1-x)^p}. \quad (4)$$

Extensive experiments have shown this model to be flexible enough to model cycle-time moments of realistic manufacturing simulations (Allen 2003 and Johnson et al. 2004). The exponent p , the polynomial order t , and the coefficient vector $\mathbf{c} = (c_0, c_1, \dots, c_t)$ are unknown parameters in the model.

As explained in Yang et al. (2007), the error term $\varepsilon_j(x)$ has expectation 0 and variance $\sigma^2(x)$ which depends on x through a ‘‘variance model’’ of the form:

$$\sigma^2(x) = \frac{\sigma^2}{(1-x)^{2q}}. \quad (5)$$

Both σ^2 and q are unknown parameters. Substituting sample variance as the response in (5), we can estimate the variance model. With the estimated \hat{q} , a simple transformation of Model (3) will yield a standard nonlinear regression model with approximately constant error variance:

$$Z_j(x) = Y_j(x) \times (1 - x)^q \quad (6)$$

$$= \eta(x, \mathbf{c}, r) + \delta_j = \sum_{\ell=0}^t c_\ell x^\ell (1 - x)^r + \delta_j, \quad (7)$$

where $r = q - p$ is an unknown parameter and the error δ_j is assumed to have a constant variance σ^2 . Since the error term in (7) has constant variance, we estimate Model (7) directly. The parameters of the original moment model (3) will be obtained indirectly by noting that the coefficients \mathbf{c} in (3) coincide with those in (7), and p is estimated by the difference of the q and r estimates. The polynomial order t in the moment model is determined via extra sum of squares analysis in a forward-selection manner.

To conclude this subsection, we summarize the method described above for estimating the ν^{th} ($\nu = 1, 2, 3$) moment CT-TH curve. First, based on the data set $\mathbf{Y}^{(\nu)}$ as defined in (2), the variance model (5) is fitted and the estimated parameter \hat{q}_ν is obtained; with \hat{q}_ν , the data transformation is performed on $\mathbf{Y}^{(\nu)}$ as shown in (6), and the resulting transformed data set with stabilized variance can be represented by the following vector:

$$\mathbf{Z}^{(\nu)} = (Z_1^{(\nu)}(x_1), \dots, Z_{n(x_1)}^{(\nu)}(x_1), \dots, Z_1^{(\nu)}(x_m), \dots, Z_{n(x_m)}^{(\nu)}(x_m)). \quad (8)$$

Finally, Model (7) is fitted to $\mathbf{Z}^{(\nu)}$ and using the moment model (4), $\mu_\nu(x)$ is estimated for $\nu = 1, 2, 3$ over the throughput range $[x_L, x_U]$.

4 The Generalized Gamma Distribution

The distribution family chosen to fit the individual cycle times for manufacturing settings should be able to provide a good fit for a variety of cycle-time distributions. As noted by Rose (1999), for complicated manufacturing systems cycle times tend to be close to normally distributed. However, as the system is loaded with heavier traffic, even for complicated systems, cycle-time distributions become more and more skewed (McNeill et al. 2003). In our method, the generalized gamma distribution is adopted because, to the best of our knowledge, it is the most flexible three-parameter distribution in terms of the coverage in the skewness-kurtosis plane.

The three-parameter generalized gamma distribution (GGD3), first presented in Stacy (1962), has the following p.d.f. (probability density function)

$$g(t, a, b, k) = \frac{|k|}{\Gamma(a)} \cdot \frac{t^{ak-1}}{b^{ak}} \cdot \exp[-(t/b)^k], \quad t > 0 \quad (9)$$

$$a > 0, b > 0, k \neq 0,$$

where a and k are the shape parameters, and b the scale parameter. As illustrated in Ashkar et al. (1988), the GGD can cover a wide range of skewness as well as kurtosis. Also, the GGD includes a variety of distributions as special cases, such as exponential ($a = k = 1$), gamma ($k = 1$), and Weibull ($a = 1$) distributions. The lognormal and normal distributions also arise as limiting cases.

In addition to its shape flexibility, the reason why we adopt GGD (as opposed to the Cornish-Fisher expansion proposed by McNeill et al. (2003) or other flexible distributions such as the Johnson family) is because that it only involves three parameters, which means only the first three moment curves $\{\mu_\nu(x), \nu = 1, 2, 3\}$ need to be estimated to provide a fit of the cycle-time distribution. In our experience, precisely estimating higher moment curves could be very difficult. As explained in Section 3, the ν^{th} moment curve is estimated based on the data set (2). When $\nu \geq 4$, the steepness of the moment curve $\mu_\nu(x)$ and the heteroscedasticity of variance in the data (2) become so pronounced that it requires substantially more simulation data to obtain the well-estimated moment curves. This trend with the increasing of ν is illustrated through the M/M/1 example in Appendix A.1.

A location parameter t_0 can be added to GGD3, and the resulting 4-parameter distribution (GGD4) is obtained by shifting the lower bound of GGD3 from $t = 0$ to $t = t_0$. The properties of any variable T following a GGD4 can be derived from those of $T - t_0 \sim \text{GGD3}$. In GGD-based fitting of a cycle-time distribution, t_0 signifies the lower bound of individual cycle times, which might be known in advance. In cases where t_0 is difficult to specify, we can set $t_0 = 0$ because GGD3 is flexible enough to give an adequate fit even if the origin of the underlying distribution deviates from 0. In light of these features, we will focus our attention on GGD3 and present some of its properties that will be used in our research.

Noncentral moments of GGD3 are given by:

$$m_\nu = \frac{b^\nu \Gamma(a + \nu/k)}{\Gamma(a)} \quad \nu = 1, 2, 3, \dots, \quad (10)$$

where ν is the order of the moment. The moments are defined only if $a + \nu/k > 0$. In the remainder of this paper, we assume that the moments exist, i.e., $a + \nu/k > 0$ for $\nu = 1, 2, 3$. Choosing any three distinct values for ν will provide the equations required by the method of moments to obtain the three distribution parameters, a , b , and k .

GGD3 can be regarded as a generalization of the 2-parameter gamma distribution (GD2) by supplying a positive parameter k as an exponent for the exponential factor of GD2(a, b) whose p.d.f is:

$$f(t) = \frac{t^{a-1}e^{-t/b}}{b^a\Gamma(a)}, \quad t > 0 \quad (11)$$

$$a > 0, b > 0.$$

Suppose $T \sim \text{GGD3}(a, b, k)$, then $T' = (T/b)^k \sim \text{GD2}(a, 1)$ (standard gamma distribution). We have implemented numerical methods to calculate the α percentile of T' , $\mathcal{C}'(\alpha; a, 1)$. The corresponding percentile of variable T , say $\mathcal{C}(\alpha; a, b, k)$ can be obtained by the straightforward transformation:

$$\mathcal{C}(\alpha; a, b, k) = b(\mathcal{C}'(\alpha; a, 1))^{1/k}. \quad (12)$$

The partial derivatives of the percentile $\mathcal{C}(\alpha; a, b, k)$ with respect to the GGD parameters are as follows:

$$\begin{aligned} \frac{\partial \mathcal{C}(\alpha; a, b, k)}{\partial a} &= \frac{b}{k} (\mathcal{C}'(\alpha; a, 1))^{1/k-1} \frac{\partial \mathcal{C}'(\alpha; a, 1)}{\partial a}, \\ \frac{\partial \mathcal{C}(\alpha; a, b, k)}{\partial b} &= (\mathcal{C}'(\alpha; a, 1))^{1/k}, \\ \frac{\partial \mathcal{C}(\alpha; a, b, k)}{\partial k} &= \frac{-b}{k^2} (\mathcal{C}'(\alpha; a, 1))^{1/k} \log \mathcal{C}'(\alpha; a, 1). \end{aligned} \quad (13)$$

Note that the partial derivative $\partial \mathcal{C}'(\alpha; a, 1)/\partial a$ can be approximately calculated using the finite difference:

$$\frac{\partial \mathcal{C}'(\alpha; a, 1)}{\partial a} \approx \frac{\mathcal{C}'(\alpha; a + \Delta a/2, 1) - \mathcal{C}'(\alpha; a - \Delta a/2, 1)}{\Delta a} \quad (14)$$

Therefore the first derivatives (13) can all be obtained numerically.

5 Estimation of Percentiles

In this section, we describe in detail how we fit the generalized gamma distribution at a given throughput level $x \in [x_L, x_U]$ based on the first three moment estimates, how the percentiles are estimated once $G(t, \hat{a}(x), \hat{b}(x), \hat{k}(x))$ is obtained, and how an approximate standard error can be provided for the percentile estimators.

5.1 Point Estimation

As explained in Section 3, the first three moment curves can be fitted simultaneously based on a number of simulation experiments performed at different levels of throughput. Therefore,

for any $x \in [x_L, x_U]$, the first three moments can be predicted by $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$. Substituting the moment estimates into (10) results in the following equations:

$$\begin{aligned}\widehat{\mu}_1(x) &= \frac{\widehat{b}(x)\Gamma(\widehat{a}(x) + 1/\widehat{k}(x))}{\Gamma(\widehat{a}(x))}, \\ \widehat{\mu}_2(x) &= \frac{\widehat{b}(x)^2\Gamma(\widehat{a}(x) + 2/\widehat{k}(x))}{\Gamma(\widehat{a}(x))}, \\ \widehat{\mu}_3(x) &= \frac{\widehat{b}(x)^3\Gamma(\widehat{a}(x) + 3/\widehat{k}(x))}{\Gamma(\widehat{a}(x))}.\end{aligned}\tag{15}$$

Solving the three Equations (15) numerically gives the three estimated distribution parameters $(\widehat{a}(x), \widehat{b}(x), \widehat{k}(x))$ for the fitted GGD distribution at throughput x . With the estimated distribution of cycle time at throughput rate x , $G(t; \widehat{a}(x), \widehat{b}(x), \widehat{k}(x))$, the percentile $\mathcal{C}_x(\alpha)$ can be estimated for any $\alpha \in [\alpha_L, \alpha_U]$ utilizing the relationship (12).

5.2 Statistical Inference for the Percentile Estimator

Drawing inference about a parameter obtained indirectly is in general difficult. In this paper, the delta method (Lehmann 1999) is applied to make inferences concerning the estimated percentiles.

The percentile $\mathcal{C}_\alpha(x)$ is estimated based on the fitted GGD distribution, and is obviously a function of the distribution parameters, a , b and k . The first order approximation using the delta method provides the following estimation for the variance of percentile estimators, where for convenience we suppress the dependence of a , b , and k on x :

$$\begin{aligned}\text{Var}[\widehat{\mathcal{C}}_\alpha(x)] &\doteq \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial a}\right)^2 \text{Var}[\widehat{a}] + \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial b}\right)^2 \text{Var}[\widehat{b}] + \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial k}\right)^2 \text{Var}[\widehat{k}] \\ &+ 2 \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial a}\right) \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial b}\right) \text{Cov}[\widehat{a}, \widehat{b}] + 2 \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial b}\right) \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial k}\right) \text{Cov}[\widehat{b}, \widehat{k}] \\ &+ 2 \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial k}\right) \left(\frac{\partial \mathcal{C}_\alpha(x)}{\partial a}\right) \text{Cov}[\widehat{k}, \widehat{a}].\end{aligned}\tag{16}$$

In (16), the partial derivatives of the percentile $\mathcal{C}_\alpha(x)$ with respect to the GGD parameters can be approximately calculated from (13) by substituting the estimates, \widehat{a} , \widehat{b} , \widehat{k} and $\widehat{\mathcal{C}}_\alpha(x)$.

Since the GGD parameters are estimated by matching the first three moments of the GGD distribution to the moment estimates $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$, the variances and covariances in (16) are functions of the variances and covariances of $\widehat{\mu}_1(x)$, $\widehat{\mu}_2(x)$ and $\widehat{\mu}_3(x)$. This is where the delta method (first-order approximation) is applied for a second time. Using

matrix notation, we have the following relationship as derived in Ashkar et al. (1988):

$$\begin{pmatrix} \text{Var}[\widehat{a}] \\ \text{Var}[\widehat{b}] \\ \text{Var}[\widehat{k}] \\ \text{Cov}[\widehat{a}, \widehat{b}] \\ \text{Cov}[\widehat{a}, \widehat{k}] \\ \text{Cov}[\widehat{b}, \widehat{k}] \end{pmatrix} \doteq \begin{pmatrix} C_{11} & C_{12} & \cdots & C_{16} \\ C_{21} & C_{22} & \cdots & C_{26} \\ C_{31} & C_{32} & \cdots & C_{36} \\ C_{41} & C_{42} & \cdots & C_{46} \\ C_{51} & C_{52} & \cdots & C_{56} \\ C_{61} & C_{62} & \cdots & C_{66} \end{pmatrix}^{-1} \begin{pmatrix} \text{Var}[\widehat{\mu}_1(x)] \\ \text{Var}[\widehat{\mu}_2(x)] \\ \text{Var}[\widehat{\mu}_3(x)] \\ \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_2(x)] \\ \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_3(x)] \\ \text{Cov}[\widehat{\mu}_2(x), \widehat{\mu}_3(x)] \end{pmatrix}, \quad (17)$$

where the matrix \mathbf{C} is given by:

$$\begin{pmatrix} D_{11}^2 & D_{12}^2 & D_{13}^2 & 2D_{11}D_{12} & 2D_{11}D_{13} & 2D_{12}D_{13} \\ D_{21}^2 & D_{22}^2 & D_{23}^2 & 2D_{21}D_{22} & 2D_{21}D_{23} & 2D_{22}D_{23} \\ D_{31}^2 & D_{32}^2 & D_{33}^2 & 2D_{31}D_{32} & 2D_{31}D_{33} & 2D_{32}D_{33} \\ D_{11}D_{21} & D_{12}D_{22} & D_{13}D_{23} & D_{11}D_{22} + D_{12}D_{21} & D_{11}D_{23} + D_{13}D_{21} & D_{12}D_{23} + D_{13}D_{22} \\ D_{11}D_{31} & D_{12}D_{32} & D_{13}D_{33} & D_{11}D_{32} + D_{12}D_{31} & D_{11}D_{33} + D_{13}D_{31} & D_{12}D_{33} + D_{13}D_{32} \\ D_{21}D_{31} & D_{22}D_{32} & D_{23}D_{33} & D_{21}D_{32} + D_{22}D_{31} & D_{21}D_{33} + D_{23}D_{31} & D_{22}D_{33} + D_{23}D_{32} \end{pmatrix} \quad (18)$$

with $D_{rj} = \partial\mu_r(x)/\partial\zeta_j$ ($r, j = 1, 2, 3$) and $\zeta_1 = a$, $\zeta_2 = b$, $\zeta_3 = k$.

From (10), the partial derivatives can be approximately calculated by substituting the estimated GGD parameters into

$$\begin{aligned} D_{11} &= \frac{-b}{k^2\Gamma(a)}[\Gamma(a+1/k)\Psi(a+1/k)], & D_{12} &= \frac{\Gamma(a+1/k)}{\Gamma(a)}, \\ D_{13} &= b\frac{\Gamma(a+1/k)}{\Gamma(a)}[\Psi(a+1/k) - \Psi(a)], \\ D_{21} &= \frac{-2b^2}{k^2\Gamma(a)}[\Gamma(a+2/k)\Psi(a+2/k)], & D_{22} &= 2k\frac{\Gamma(a+2/k)}{\Gamma(a)}, \\ D_{23} &= b^2\frac{\Gamma(a+2/k)}{\Gamma(a)}[\Psi(a+2/k) - \Psi(a)], \\ D_{31} &= \frac{-3b^3}{k^2\Gamma(a)}[\Gamma(a+3/k)\Psi(a+3/k)], & D_{32} &= 3k^2\frac{\Gamma(a+3/k)}{\Gamma(a)}, \\ D_{33} &= b^3\frac{\Gamma(a+3/k)}{\Gamma(a)}[\Psi(a+3/k) - \Psi(a)], \end{aligned} \quad (19)$$

where

$$\Psi(t) = \frac{1}{\Gamma(t)} \frac{d\Gamma(t)}{dt}.$$

Clearly, from the derivation above, estimating $\text{Var}[\widehat{\mathcal{C}}_\alpha(x)]$ requires obtaining the moment estimators $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$ and their variances and covariances. The estimators $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$ can be obtained by following the methodology explained in Section 3, while estimating their variances and covariances is discussed in Appendix A.3.

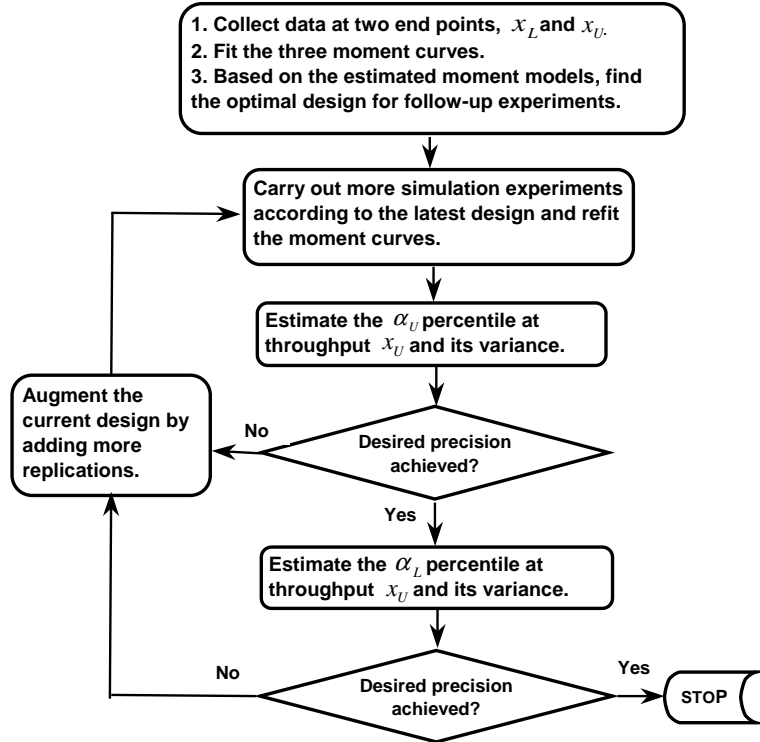


Figure 2: Flow chart for the multistage procedure.

6 Procedure for Estimating Percentiles of Cycle Time

In this section, we discuss issues related to experiment design, and give a description of the proposed procedure for estimating percentiles. To provide context, a high-level description of the procedure is provided in Figure 2.

In summary, simulation experiments are carried out sequentially until the prespecified stopping criterion is satisfied. The experimentation is initiated with a starting design which allocates an equal number (chosen by the user) of replications to the two end points of the throughput range $[x_L, x_U]$. As the procedure progresses, new design points are included and additional replications are added in batches. Each batch of replications is allocated to the design points to minimize PM (defined in equation (22)), an experiment design criterion that is related to the variance of the percentile estimators. Since the design criterion depends on unknown parameters of the moment curves, the current best estimates of the parameters are used in the allocation of each batch of replications. As more simulation data are collected, increasingly precise estimators are obtained until the precision of the estimators matches the stopping criterion.

6.1 Experiment Design

As already noted, the experiments are started by a design that allocates an equal number of replications to the upper and lower end of the throughput range. This design will then be augmented by including more design points and adding more replications as the procedure progresses. To determine the follow-up design, we need to answer three questions based on the estimates obtained from the current data set. (i) How many additional replications, say ΔN , will be added? (ii) At what design points (throughput levels) will the simulations be executed? (iii) How many of the ΔN replications should be allocated to each design point? We use the vector $\mathbf{x} = (x_1, x_2, \dots, x_m)$ to represent the set of design points included in the design, and $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_m)$ the fractions for the total replications assigned to each design point. At each step, the values of \mathbf{x} and $\boldsymbol{\pi}$ will be determined conditional on the fact that some replications have already been allocated to certain design points.

6.1.1 Design Criterion

Our goal is to develop a method to estimate the percentile $\mathcal{C}_\alpha(x)$ for $x \in [x_L, x_U]$ and $\alpha \in [\alpha_L, \alpha_U]$. Therefore, the experiment design will seek to minimize some measure of the variance of $\widehat{\mathcal{C}}_\alpha(x)$. Suppose that N is the number of replications available for allocation. A natural performance measure, which is inherited from Cheng and Kleijnen (1999), is the weighted average variance over the throughput range of interest normalized by N :

$$PM_0 = N \frac{\int_{x_L}^{x_U} w(x) \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x)] dx}{\int_{x_L}^{x_U} w(x) dx} \quad (20)$$

where $w(x)$ is the weight function which in the simplest case is 1, and $N \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x)]$ is the normalized variance which is independent of N . As explained in Section 5.2, from the first three fitted moment curves $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$, the variance $\text{Var}[\widehat{\mathcal{C}}_\alpha(x)]$ can be estimated for any percentile estimator $\widehat{\mathcal{C}}_\alpha(x)$ ($x \in [x_L, x_U]$, $\alpha \in [\alpha_L, \alpha_U]$). We chose to base (20) on the variance of the largest percentile α_U because $\widehat{\mathcal{C}}_{\alpha_U}(x)$ is typically much more variable than other percentile estimators. Unfortunately, it is not practical to determine $[\mathbf{x}, \boldsymbol{\pi}]$ by minimizing PM_0 , because $\text{Var}[\widehat{\mathcal{C}}_\alpha(x)]$ can only be numerically estimated for given values of x and α (Section 5.2). Hence, we use the simple finite difference approximation of (20):

$$PM_0 \doteq N \sum_{\kappa \in \Omega_x} \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)] \cdot \Delta\kappa \quad (21)$$

where Ω_x is a chosen set of evenly spaced grid points in the range $[x_L, x_U]$, and $\Delta\kappa$ is the interval between two neighboring points. Obviously, $\Delta\kappa$ is a constant which can be dropped

from (21), so we define our design criterion as

$$PM = N \sum_{\kappa \in \Omega_x} \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)]. \quad (22)$$

Measure PM is a function of the design $[\mathbf{x}, \boldsymbol{\pi}]$ as illustrated in Appendix A.5. Evaluating PM for given $[\mathbf{x}, \boldsymbol{\pi}]$ comes down to providing an estimate for $\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)]$, which can be obtained at any later stage of experimentation where simulation data are available for the estimation of the first three moment curves (for details, see Appendix A.5). Hence, at a point where further experiments are to be carried out, PM can be approximately calculated for given $[\mathbf{x}, \boldsymbol{\pi}]$, which enables us to apply a numerical search method to the problem of optimizing PM . Note that $N \times \boldsymbol{\pi}$ is not restricted to be integer in the search for the optimal solution of $[\mathbf{x}, \boldsymbol{\pi}]$.

Next, we will give the details regarding how the optimization problem is constructed and solved to augment the current experiments at each stage.

6.1.2 Optimal Experiment Design

We propose solving the following constrained nonlinear optimization problem to guide further simulation experiments given that some replications have already been allocated.

$$\min_{\mathbf{x}, \boldsymbol{\pi}} PM(\mathbf{x}, \boldsymbol{\pi}) \quad (23)$$

$$s.t. \quad \{x_1, x_2, \dots, x_m\} \supseteq \{\widehat{x}_1, \widehat{x}_2, \dots, \widehat{x}_{m_c}\} \quad (24)$$

$$x_L \leq x_1 < x_2 < \dots < x_m \leq x_U$$

$$\sum_{i=1}^m \pi(x_i) = 1$$

$$\pi(x_i) \geq lb(x_i) \quad \text{for } i = 1, 2, \dots, m$$

The input parameters, decision variables, and constraints of (23) are given as follows.

Input parameters:

- The range of throughput $[x_L, x_U]$.
- m_c and m ($m \geq m_c$), the number of design points before and after augmenting the design, respectively.
- The old design points $\{\widehat{x}_1, \widehat{x}_2, \dots, \widehat{x}_{m_c}\}$, and the allocation of simulation replications already made at those points $\{n_c(\widehat{x}_1), n_c(\widehat{x}_2), \dots, n_c(\widehat{x}_{m_c})\}$. Note that $n_c(x) = 0$ for $x \notin \{\widehat{x}_1, \widehat{x}_2, \dots, \widehat{x}_{m_c}\}$.

- The total number of replications already allocated N_c , and the increment of replications to be added to the current design ΔN . Therefore we have $N = N_c + \Delta N$. For the same reason as explained in Yang et al. (2007), guiding the choice of ΔN at each stage is important, and the way to determine ΔN is detailed in Appendix A.6. Both N_c and ΔN are used to calculate the lower bounds $lb(x_i) = \max\{n_c(x_i), 2\}/(N_c + \Delta N)$ for $i = 1, 2, \dots, m$. We set $lb_i \geq 2/(N_c + \Delta N)$ to insure that at least 2 replications are assigned to any point x_i included in the design.

Decision variables:

- The new set of design points $\mathbf{x} = \{x_1, x_2, \dots, x_m\}$, whose values are forced to be increasing in the subscript.
- The updated allocation of simulation effort $\boldsymbol{\pi} = \{\pi_1, \pi_2, \dots, \pi_m\}$, which has to be mapped into integer numbers of replications, say $\{n(x_i), i = 1, 2, \dots, m\}$, assigned to the design points. We use a simple rounding in our method by setting $n(x_i) = \lceil N\pi_i \rceil$. By rounding up we insure that each design point gets at least as many replications as called for by the optimal solution, and since our goal is to achieve a fixed precision (rather than optimize a fixed budget) this can do no harm. The resulting integer solution could be suboptimal, but seems to work well in our numerical experiments.

Constraints:

- The constraint (24) forces the new set of design points to include the old points.
- The meanings of the other constraints are obvious.

To solve the optimization problem, we need to choose starting values for each of the decision variables. In all the experiments considered in this paper, the starting values of x are chosen to be evenly spaced throughout the interval of throughput, and the fraction of replications at each design point x_i initiates from $N^{-1}(n_c(x_i) + \Delta N/m)$ ($i = 1, 2, \dots, m$).

In the procedure, (23) is solved to augment the current design when an assessment of the chosen percentile estimates shows that subsequent experimental effort is necessary. The design may be augmented in two different ways: 1) adding design points and replications, and 2) adding replications only. Augmentation of type 1 only occurs once in the procedure when we expand the starting design which only consists of experiments performed at the two end points, x_L and x_U , to a m -point design. (Guidelines for determining the number

of design points m is provided in Section 6.3.) Afterwards, the location of design points are fixed, and only the allocation of simulation effort can be modified by assigning more replications to the current design points. In our experiments, the optimization problem (23) is coded in Matlab, and the Matlab optimization function “fmincon” is used to solve the nonlinear constrained problem (with $m = 5$ as will be explained in Section 6.3), which takes about 150 seconds on a computer with processor speed of 3GHz.

6.2 Stopping Rule

The proposed procedure collects simulation data to allow for estimation of $\mathcal{C}_\alpha(x)$ for $x \in [x_L, x_U]$ and $\alpha \in [\alpha_L, \alpha_U]$ with both ranges of interest being specified by the user. Moreover, our methods provides an error estimate for any percentile estimator $\widehat{\mathcal{C}}_\alpha(x)$ ($x \in [x_L, x_U]$, $\alpha \in [\alpha_L, \alpha_U]$)(Section 5.2). Obviously, the upper end of throughput is where the variability of cycle time is most pronounced, and it is known that estimators of larger percentiles are more variable than their lower counterparts. Consequently, $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$ is considered to possess the highest variability among all the estimable percentiles, which motivates us to use the relative error of $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$ as the stopping criterion for our procedure. By controlling the precision of the most variable estimator $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$, we hope that other percentiles will also be well estimated.

Specifically, we let the user specify a precision level, say $100\gamma\%$, and the procedure terminates only when the condition

$$\frac{2\text{SE}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)]}{\widehat{\mathcal{C}}_{\alpha_U}(x_U)} \leq 100\gamma\%$$

is satisfied. We define $\text{SE}[\cdot] = \sqrt{\text{Var}[\cdot]}$. Moreover, a safe fall-back strategy is adopted. As illustrated in Figure 2, a check is also performed on the precision of $\widehat{\mathcal{C}}_{\alpha_L}(x_U)$, and simulation data will be collected until

$$\frac{2\text{SE}[\widehat{\mathcal{C}}_{\alpha_L}(x_U)]}{\widehat{\mathcal{C}}_{\alpha_L}(x_U)} \leq 100\gamma\%.$$

Constraining the precision of any percentile estimator $\mathcal{C}_\alpha(x)$ within a certain prespecified level is difficult, but by controlling the relative precision of the two estimators $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$ and $\widehat{\mathcal{C}}_{\alpha_L}(x_U)$, we hope to impose precision control on percentile estimators throughout the range we consider.

6.3 The Multistage Procedure

This subsection is devoted to an overall description of the multistage procedure, which is diagrammed in Figure 2.

The procedure is divided into two stages: In the initial stage, pilot simulation runs are performed at the two end points of the throughput range to provide the preliminary data for model estimation; in the second stage, the experiments are augmented to include, say, m design points, and simulation runs are added in an efficient manner until the desired precision level on the chosen percentile estimators is achieved.

The number of design points m is a user-specified parameter, which has to be set to allow for the good estimation of the first three moment curves; that is, m should be sufficiently large to allow the moment model (4) to include enough polynomial terms to generate a good fit for the ν^{th} ($\nu = 1, 2, 3$) moment curve. As pointed out by Yang et al. (2007), the value of m must be determined through consideration of the system being investigated. In our extensive empirical experiments with both simple queueing models and realistic manufacturing systems, we have never encountered a situation where 5 design points cannot provide an adequate fit for the moment curves over $[x_L, x_U] = [0.5, 0.95]$ (a throughput range much wider than the range within which a real manufacturing system is typically operated). Hence, we recommend setting $m = 5$ if no reliable information is available to suggest the use of less points.

Inputs: Simulation model of the system being investigated, precision level $100\gamma\%$ which is defined as the relative error on the chosen percentiles, throughput range $[x_L, x_U]$, percentage range $[\alpha_L, \alpha_U]$, number of design points m , initial number of replications N_0 .

Outputs: Fitted moment curves $\{\hat{\mu}_\nu(x), \nu = 1, 2, 3\}$ and the inferred variance-covariance information, from which the percentile estimate $\hat{\mathcal{C}}_\alpha(x)$ and an approximate standard error $\widehat{\text{SE}}[\hat{\mathcal{C}}_\alpha(x)] = \sqrt{\widehat{\text{Var}}[\hat{\mathcal{C}}_\alpha(x)]}$ can be provided for $x \in [x_L, x_U]$ and $\alpha \in [\alpha_L, \alpha_U]$.

Stage 0. Initially, N_0 replications are allocated evenly to the two end points x_L and x_U . The three moment curves $\{\hat{\mu}_\nu(x), \nu = 1, 2, 3\}$ are then estimated by fitting Models (5) and (7) as described in Section 3. At this point, the polynomial order of $\{\hat{\mu}_\nu(x), \nu = 1, 2, 3\}$ is equal to 0 because of the constraint imposed by the number of design points, 2. With the estimated moment models and the inferred variance information, we

1. determine ΔN , the number of replications to be added to the initial design (see Appendix A.6 for the determination of the value of ΔN);

2. find the optimal design $(\mathbf{x}, \boldsymbol{\pi})$ consisting of m points by solving the nonlinear optimization problem (23).

Stage 1. In this stage, we fix the m design points, and keep allocating more replications to those points until the desired precision is achieved. Three tasks are to be completed in the following steps.

Step 1: Run more simulation experiments. Assign ΔN additional runs to the design points found in the previous stage according to the latest updated loadings $\boldsymbol{\pi}$. Refit the three moment curves, and search for appropriate polynomial order of each fitted model; we follow the forward-selection method suggested in Yang et al. (2007). Obtain the estimate of the distribution of cycle time at x_U , $G(t; \hat{a}(x_U), \hat{b}(x_U), \hat{k}(x_U))$, and then estimate the percentile $\mathcal{C}_{\alpha_U}(x_U)$ by inverting the c.d.f.

Step 2: Evaluate the precision of $\hat{\mathcal{C}}_{\alpha_U}(x_U)$. Estimate the standard error of $\hat{\mathcal{C}}_{\alpha_U}(x_U)$. If the desired precision is achieved ($2\widehat{\text{SE}}[\hat{\mathcal{C}}_{\alpha_U}(x_U)]$ is less than $\gamma\%$ of $\hat{\mathcal{C}}_{\alpha_U}(x_U)$), then move to Step 3. Otherwise, conditional on the current design points, find the value of ΔN at the current point and solve (23) to adjust the loadings $\boldsymbol{\pi}$ of the design according to the latest estimated moment curves. Go back to Step 1.

Step 3: Evaluate the precision of $\hat{\mathcal{C}}_{\alpha_L}(x_U)$. Estimate the standard error of $\hat{\mathcal{C}}_{\alpha_L}(x_U)$. If the desired precision is achieved ($2\widehat{\text{SE}}[\hat{\mathcal{C}}_{\alpha_L}(x_U)]$ is less than $\gamma\%$ of $\hat{\mathcal{C}}_{\alpha_L}(x_U)$), then stop. Otherwise, conditional on the current design points, solve (23) to adjust the loadings $\boldsymbol{\pi}$ of the design according to the latest estimated moment curves. Go back to Step 1.

7 Numerical Evaluation

In this section, we evaluate the performance of the proposed procedure based on queueing models. In our experiments, we have considered the following G/G/1 queueing systems: M/M/1, M/E₂/1, D/E₂/1 and D/M/1. These models cover deterministic, Erlang, and exponential (representing no, moderate and high variability) distributions for the interarrival and processing times, and they represent a range of cycle-time distributions while still being analytically tractable. We use these simple models to allow control of factors that might affect procedure performance; a realistic full factory simulation is studied in the next section.

Not surprisingly, our procedure performs best on M/M/1, where the assumptions concerning the form of moment models and the distribution of cycle times are known to be true. Among these four systems, our procedure has the worst performance on the D/M/1 system. Due to space constraints, we only present the results for M/M/1 and D/M/1.

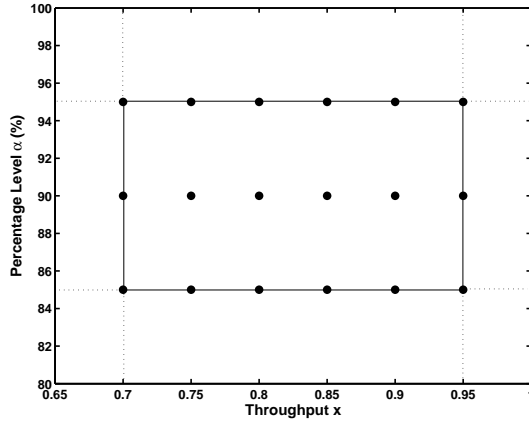


Figure 3: Check Points Selected in the Feasible Region

7.1 Results for Queueing Systems

For both M/M/1 and D/M/1, the true percentiles of cycle time at different throughputs can be analytically computed, and hence the quality of percentile estimation can be easily evaluated. For each model, the proposed procedure was applied 100 times, and from each of the 100 macro-replications, selected percentile estimates were compared to their true values.

In our experiments, the throughput range of interest was chosen to be $[x_L, x_U] = [0.7, 0.95]$, and the percentile range $[\alpha_L, \alpha_U] = [0.85, 0.95]$, where we have normalized the throughput so that the maximum system capacity is 1. The precision level of the relative error used as the stopping criterion was set at $100\gamma\% = 5\%$ (see Section 6.2). For all the queueing models considered, the location parameter t_0 (see Section 4) was set at 0 throughout the throughput range. As already noted, our procedure is able to give percentile estimates $\mathcal{C}_\alpha(x)$ for any point in the two-dimensional region defined by the percentile $\alpha \in [\alpha_L, \alpha_U]$ and throughput $x \in [x_L, x_U]$. We call this region the feasible region. To evaluate the accuracy and precision of the percentile estimation, check points were selected inside this feasible region, as shown in Figure 3. At each of these points, the estimates were compared to the true percentiles of the queueing system.

7.1.1 Point Estimators

All the point estimators for percentiles performed similarly well in terms of deviation from the true value for both M/M/1 and D/M/1. Two types of plots were made to display graphically the 100 realizations of each percentile estimator made at the check points: (i)

relative error plots, where the y -axis is defined as

$$\frac{\text{Percentile Estimate} - \text{True Percentile}}{\text{True Percentile}} \times 100\%, \quad (25)$$

and (ii) absolute error plots, in which percentile estimates are plotted around their true values.

Figure 4 shows the percentile estimation results for M/M/1. Figure 4a–c are relative error plots with the percentile α being 85%, 90%, and 95%, respectively. For these graphs, the x -axis represents throughput rate x , and every point in the graph represents the relative deviation at corresponding check point (α, x) calculated by (25) from one of the 100 macro-replications. Notice that a very high proportion of the relative deviations of the percentile estimates at the selected check points are within 5% (the precision level $100\gamma\%$ imposed prior to experimentation). Figures 4a'–c' are the absolute error plots, in which the solid curve represents a piecewise linear version of the true percentile curve across the throughput range and the percentile estimates are plotted in absolute units. From these plots, it is evident that the variability of the percentile estimators at the highest throughput $x_U = 0.95$ is the most pronounced, and as explained in Section 6.2, it has been well controlled in our procedure.

Figure 5 shows an analogous plot for the D/M/1 system, and similar conclusions can be drawn, although the performance not as good as the M/M/1 especially when the throughput is at $x = 0.95$.

7.1.2 Standard Error

An estimator of the standard error $\text{SE}[\widehat{\mathcal{C}}_\alpha(x)] = \sqrt{\text{Var}[\widehat{\mathcal{C}}_\alpha(x)]}$ is provided for each percentile estimator $\widehat{\mathcal{C}}_\alpha(x)$ by the procedure described in Section 6.3. Our goal in this section is to evaluate the quality of the SE estimator. Tables 1 and 2 show the results for M/M/1 and D/M/1, respectively. The column labeled “Sample Stdev” is the sample standard deviation of the percentile point estimators calculated from the 100 realizations of the percentile estimator; therefore, it is an unbiased estimator of the true standard error. The “Average SE” column is the average of the 100 standard error estimators $\widehat{\text{SE}}[\widehat{\mathcal{C}}_\alpha(x)]$, each one of which is estimated from within a single macro-replication.

Table 1 shows that for M/M/1, the mean of the standard error estimate in the “Average SE” column is close to, but consistently less than, the unbiased external estimate of the standard deviation found in the “Sample Stdev” column. The underestimation trend is more apparent for the D/M/1. Nevertheless, the estimated standard error $\widehat{\text{SE}}[\widehat{\mathcal{C}}_\alpha(x)]$ provided by

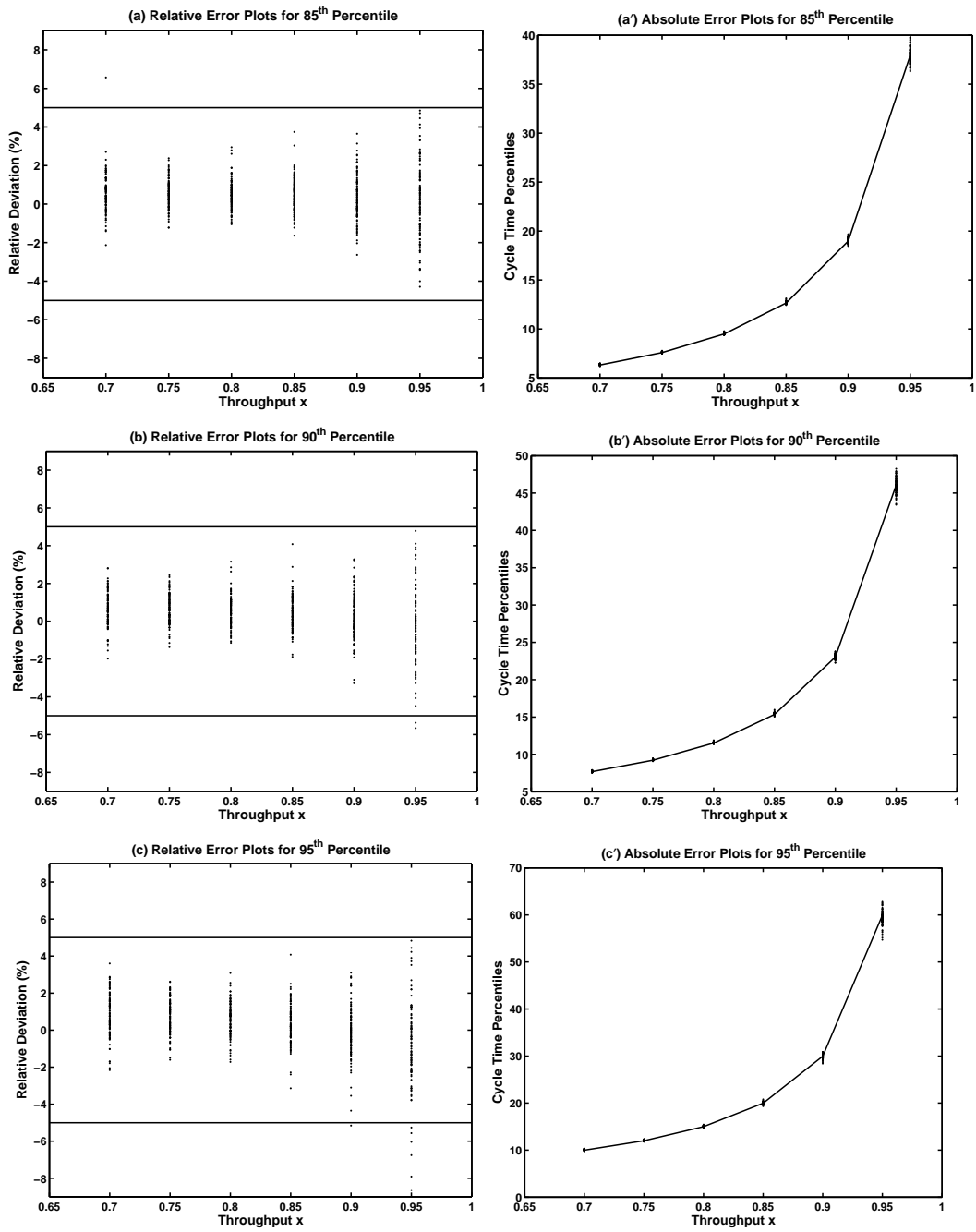


Figure 4: Plots of the Percentile Estimates for M/M/1 (100 Macro-replications)

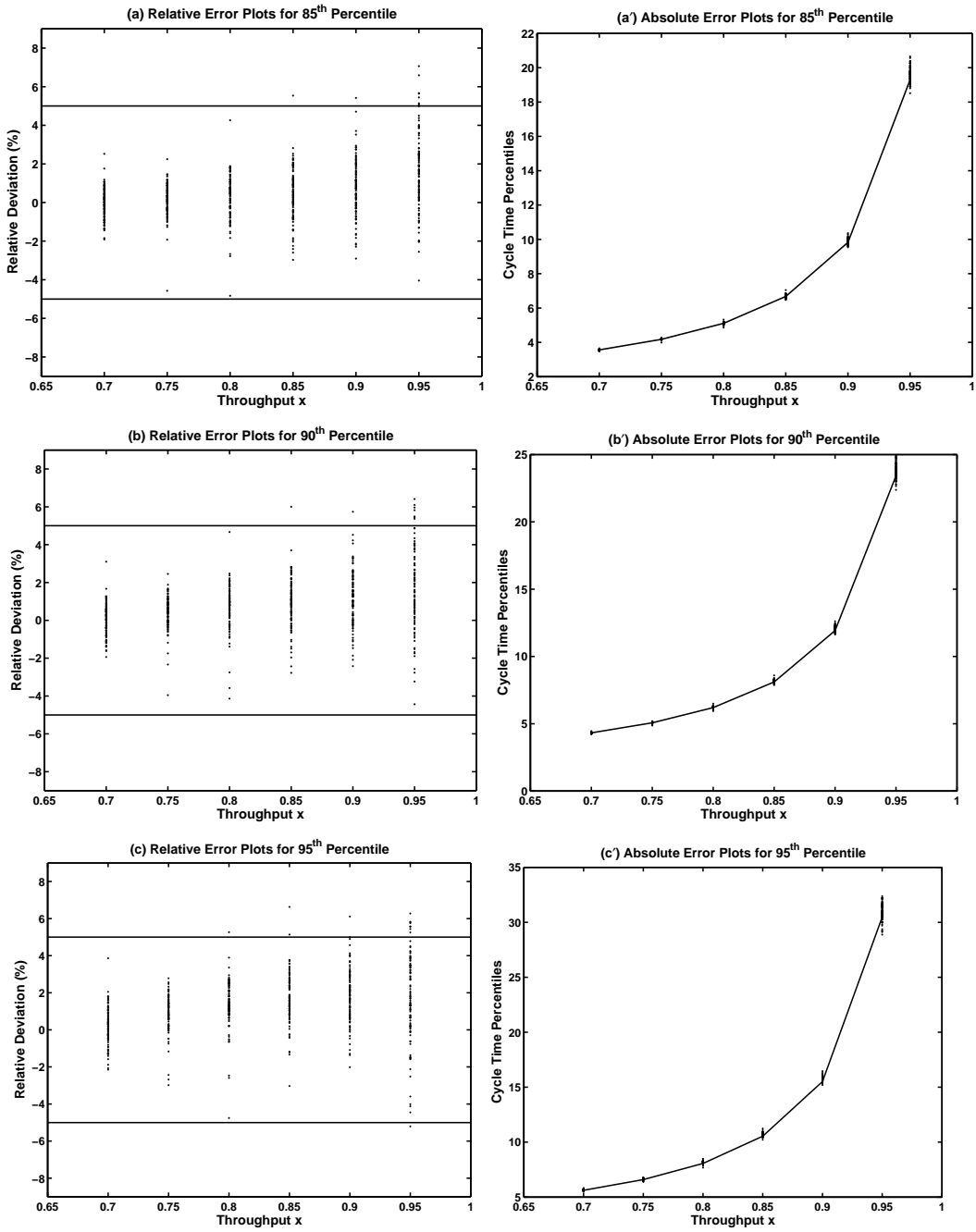


Figure 5: Plots of the Percentile Estimates for D/M/1 (100 Macro-replications)

Table 1: Estimated Standard Errors of Percentile Estimates for M/M/1

TH x	85 th Percentile		90 th Percentile		95 th Percentile	
	Sample Stdev	Average SE	Sample Stdev	Average SE	Sample Stdev	Average SE
0.70	0.055	0.053	0.072	0.067	0.117	0.105
0.75	0.056	0.053	.071	0.067	0.112	0.106
0.80	0.070	0.064	0.089	0.080	0.137	0.125
0.85	0.110	0.100	0.139	0.126	0.213	0.190
0.90	0.225	0.213	0.282	0.270	0.430	0.399
0.95	0.759	0.734	0.950	0.926	1.439	1.361

Table 2: Estimated Standard Errors of Percentile Estimates for D/M/1

TH x	85 th Percentile		90 th Percentile		95 th Percentile	
	Sample Stdev	Average SE	Sample Stdev	Average SE	Sample Stdev	Average SE
0.70	0.028	0.024	0.035	0.029	0.056	0.046
0.75	0.035	0.026	0.044	0.032	0.065	0.048
0.80	0.060	0.042	0.074	0.052	0.103	0.079
0.85	0.090	0.065	0.106	0.082	0.142	0.125
0.90	0.144	0.105	0.171	0.132	0.224	0.195
0.95	0.399	0.373	0.503	0.471	0.733	1.696

the procedure can still give the user a rough idea about how variable the percentile estimator is.

In the absence of any knowledge about the distribution of the percentile estimators, it would be natural to attempt to form a 95% confidence interval for the percentile by using:

$$\widehat{\mathcal{C}}_\alpha(x) \pm 1.96 \times \widehat{\text{SE}}[\widehat{\mathcal{C}}_\alpha(x)]. \quad (26)$$

For M/M/1, (26) works well in terms of coverage and gives a conservative CI. However, for D/M/1, the coverage probability was lower than the nominal level. This can be explained by underestimation of the standard error, and non-normality of the percentile of cycle time estimator. In the case with D/M/1, it appears that non-normality is the dominant factor.

7.2 Summary of Results

Through experimentation with queueing models, it has been shown that the proposed procedure has the potential to be effective in providing accurate and precise percentile estimators. By controlling the relative standard error of the percentile estimators at the upper end of the throughput range, high precision has been achieved for estimators of percentiles throughout the feasible region.

For each percentile estimator, an estimate of the standard error is also provided which gives the user a sense of its variability. However, in the scope of our work, there is not sufficient information to draw any conclusion regarding the distribution (or limiting distribution) of the percentile estimators. Thus, no reliable confidence interval can be created based on the standard error estimation.

Of course, real manufacturing systems are networks of queues. To stress our procedure in a realistic setting, we next consider a semiconductor fabrication model.

8 An Example of Manufacturing Systems

In this section, we apply the proposed procedure to a semiconductor wafer fab simulation model representing a full manufacturing factory. The model (Testbed Data Set # 1 created in Factory Explorer) was taken from the website of the Modeling and Analysis for Semiconductor Manufacturing Lab at Arizona State University (www.eas.asu.edu/~masmlab/). The model is designed to process two types of jobs, Prod1 and Prod2, with each type being released into the system at a constant (deterministic) rate. Jobs of different types follow different process steps, and thus have different cycle-time distributions. The primary sources of variability are machine failures and repairs.

In our experiments, the product mix (expressed as a percent of production dedicated to each product type) is set as 66.7% Prod1 and 33.3% Prod2. We investigate the CT-TH relationships for the two types of products separately. For the percentile of cycle time curves to be generated, the independent variable, throughput, was defined as the overall production rate (as a percentage of the capacity) of both types of jobs that are mixed with a constant ratio. Note that the cycle-time distribution for a particular type of product also depends on the product mix, and in this paper, we restrict ourselves to situations where the jobs are released with fixed product mix. The construction of CT-TH-PM (product mix) surfaces is the subject of ongoing research.

As already indicated, our objective is to estimate the CT-TH percentile curves for both Prod1 and Prod2 based on a single set of simulation runs. In our experiments, we chose to drive the simulation by the precision of Prod2. After accumulating sufficient data for the estimation of Prod2, we estimate the percentile curves for both products. For the implementation of our procedure, the range of throughput was chosen to be $[0.7, 0.95]$, where “1” corresponds to system capacity, and the percentile range $[85\%, 95\%]$. The precision level was set at $100\gamma\% = 1\%$, and the number of design points $m = 5$. In the remainder of this section, we will discuss the results for Prod2 in detail; similar conclusions can be drawn for Prod1.

As explained in Section 4, we allow the user to introduce a fourth parameter t_0 , which represents the lower bound of the GGD representing cycle time. There are at least three different ways to set the location parameter $t_0(x)$ for the GGD distribution fitted at throughput rate x . 1) In the absence of any knowledge about the lower bound of cycle time, 0 can always be used as the default value of $t_0(x)$ for any x . 2) The pure minimum processing time of the product being considered, which is usually available to the user, can be safely used for the location parameter throughout the range of throughput. These two simple settings can provide good percentile estimates as will be explained later. However, to achieve better precision, we recommend using a third method which imposes a much tighter lower bound on the cycle time. As already noted, the distribution of steady-state cycle time varies with the range of throughput as illustrated in Figure 6 which gives histograms of 50,000 individual cycle times for Prod2 at two throughput levels, 0.7 and 0.95. Although the theoretical pure processing time is not a function of throughput, the impact of queueing is to make the effective minimum cycle time much larger, from about 450 hours at $x = 0.7$ to about 680 hours at $x = 0.95$. As can be seen from the graphs, at high throughput levels, the steady-state cycle times are bounded well away from their pure processing time (223 hours), the theoretical minimum. Our experiment results have shown that using the empirical minimum brings

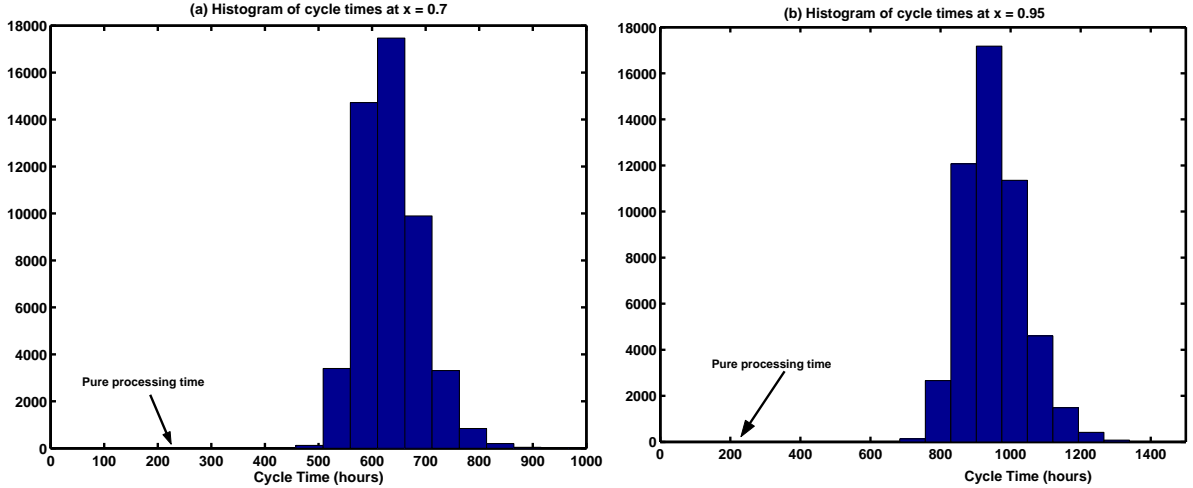


Figure 6: Histograms of cycle times (50,000 cycle times in each graph)

significant improvement to the percentile estimates. To obtain the empirical minimum of cycle time at any $x \in [x_L, x_U]$, we propose the following method: at the 5 design points where simulation experiments are performed, the empirical minimal cycle time can be easily obtained; for other points, we use linear interpolation based on the 5 minimums at those design points.

To evaluate the percentile estimates, the check points as shown in Figure 3 are used again. Since the true percentiles at those points are unknown, substantial additional data were collected at the check points to obtain the “nearly true” estimates for percentiles of cycle times. We present in Table 3 the numerical results for Prod2 with the throughput-sensitive location parameter being the empirical minimum at each throughput level. In Table 3, absolute deviations (defined as percentile estimates minus “true percentiles”) and relative deviations (defined by (25)) of the estimated percentiles are given, and an estimate of the standard error of the percentile estimator is also provided.

The point estimates of percentiles are good in terms of the relative error; with the exception of two check points at which the absolute value of the relative error is slightly above 1%, the deviations are well within the desired precision. From the sign of the deviations, it is obvious that at the lowest throughput rate the percentiles are overestimated, and that at the higher throughput levels the percentiles are underestimated. We conjecture that the bias in percentile estimates was inherited from the moment estimates. We compared the estimated moments obtained from the procedure with (very precisely estimated) “true” moments, and detected the same pattern: the first three moments are all slightly overestimated at the lower end of the throughput range while being slightly underestimated at the other throughput

Table 3: Results for the semiconductor manufacturing model

TH x	85 th Percentile			90 th Percentile			95 th Percentile		
	Abs.Dev.	Rel.Dev.	Est.SE	Abs.Dev.	Rel.Dev.	Est.SE	Abs.Dev.	Rel.Dev.	Est.SE
0.70	4.62	0.66%	0.84	5.49	0.76%	1.21	5.67	0.76%	1.96
0.75	-2.53	-0.35%	0.63	-3.27	-0.44%	0.87	-1.98	-0.26%	1.39
0.80	-6.88	-0.91%	0.47	-7.77	-1.00%	0.66	-7.34	-0.91%	1.07
0.85	-6.93	-0.86%	0.63	-7.12	-0.86%	0.90	-6.87	-0.80%	1.48
0.90	-3.98	-0.45%	0.92	-3.38	-0.37%	1.17	-1.95	-0.21%	1.71
0.95	-9.60	-0.91%	1.64	-10.04	-0.94%	2.31	-11.78	-1.06%	3.73

levels. As explained in Yang et al. (2007), this consistent pattern in the estimation on moment estimates is what we expected. Due to the form of the moment model (4), the fitted moment curve is likely to increase smoothly and intersect with the underlying true moment curve at some point within the throughput range. In this case, for all three moment curves, the intersection point is somewhere close to the lower end: at throughput levels lower than the intersection we overestimate the moments, and at throughput levels higher than the intersection we underestimate the moments.

Based on the same data set, percentile estimates were also obtained using different settings of the location parameter: (i) $t_0(x) = 0$; and (ii) $t_0(x)$ equal to the pure processing time. The percentile estimates obtained from these two settings are still fairly good in terms of relative error: for case (i), the relative error at all the check points was within 3.5%; and the precision achieved in case (ii) is slightly better.

9 Summary

Estimating percentiles of cycle time via simulation is difficult due to the high variability of percentile estimators and the diversity of cycle-time distributions. This paper proposes a new methodology for estimating multiple cycle-time percentiles throughout the throughput range of interest based on a single set of simulation runs. It has been shown through experiments on queueing models, such as M/M/1 and D/M/1, and a real semiconductor manufacturing simulation, that the multistage procedure provides good point estimators for percentiles of cycle time.

As a by-product of our research, our moment curves can also be used to obtain other summary statistics, such as the standard deviation of cycle time. Our fitting techniques could

also be employed to estimate simulation-generated “clearing functions” (a type of throughput vs. work-in-process inventory curve; see Asmundsson et al. 2006). Clearing functions are used in production planning optimization models to allow the model to assess the impact on work in process and throughput of altering the production plan. This is the subject of ongoing research.

Acknowledgments

This research was supported by National Science Grant DMI-0140385 and SRC grant 2004-0J-1225. Additional thanks go to Professors John Fowler and Gerald Mackulak from Arizona State University. The authors thank the Area Editor, Associate Editor and referees for help in clarifying the presentation.

References

- Allen, C. (2003) The impact of network topology on rational-function models of the cycle time-throughput curve. Honors Thesis, Department of Industrial Engineering & Management Sciences, Northwestern University. Available online via users.iems.northwestern.edu/~nelsonb/Publications/CarlAllenThesis.pdf
- Asmundsson, J., Rardin, R. L., and Uzsoy, R. 2006. Tractable nonlinear production planning models for semiconductor wafer fabrication facilities. *IEEE Transactions on Semiconductor Manufacturing*, **19**, 95–111,
- Ashkar, F., Bobée, B., Leroux, D. and Morissette, D. (1988) The generalized method of moments as applied to the generalized gamma distribution. *Stochastic Hydrology and Hydraulics*, **2**, 161–174.
- Bates, D. M. and Watts, D. G. (1988) *Nonlinear Regression Analysis and Its Applications*, New York, John Wiley & Sons Inc.
- Bickel, P. J. and K. A. Doksum (2001) *Mathematical Statistics: Basic Ideas and Selected Topics*, Volume I, 2nd edition, NJ: Prentice Hall.
- Billingsley, P. (1999) *Convergence of Probability Measures*, 2nd edition. New York: John Wiley.
- Buzacott, J. A., and J. G. Shanthikumar. (1993) *Stochastic Models of Manufacturing Systems*, Upper Saddle River, NJ: Prentice-Hall.
- Chen, E. J. and Kelton, W. D. (1999) Simulation-based estimation of quantiles. In *Proceedings of the 1999 Winter Simulation Conference*, ed. P. A. Farrington, H. B. Nembhard,

- D. T. Sturrock, and G. W. Evans, 428–434. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers.
- Available via (<http://www.informs-cs.org/wsc99papers/059.PDF>)
- Cheng, R. C. H. and Kleijnen, J. P. C. (1999) Improved design of queueing simulation experiments with highly heteroscedastic responses. *Operations Research*, **47**, 762–777.
- Chien, C., D. Goldsman and B. Melamed. (1997) Large-sample results for batch means. *Management Science*, **43** (9): 1288–1295.
- Fowler, J. W., S. Park, G. T. Mackulak, and D. L. Shunk. (2001) Efficient Cycle Time-Throughput Curve Generation Using a Fixed Sample Size Procedure. *International Journal of Production Research* **39**: 2595–2613.
- Glynn, P. W and D. L. Iglehart. (1986) Estimation of steady-state central moments by the regenerative method. *Operations Research Letters* **5**: 271–276.
- Henderson, S. G. (2001) Mathematics for simulation. In *Proceedings of the 2001 Winter Simulation Conference*, B. A. Peters, J. S. Smith, D. J. Medeiros, and M. W. Rohrer, eds., IEEE, 83–94.
- Hopp, W. J., and M. L. Spearman. (2001) *Factory Physics: Foundations of Manufacturing Management*. 2nd edition. Chicago: Irwin.
- Johnson, R., F. Yang, B. E. Ankenman, and B. L. Nelson. (2004) Nonlinear regression fits for simulated cycle time vs. throughput curves for semiconductor manufacturing. *Proceedings of the 2004 Winter Simulation Conference*, 1951–1955. Available online via <<http://www.informs-cs.org/wsc04papers/260.pdf>>
- Lehmann, E. L. (1999) *Elements of Large-Sample Theory*. Springer-Verlag New York, Inc.
- Law, A. M. and W. D. Kelton (2000) *Simulation Modeling and Analysis*. 3rd edition. NY: McGraw-Hill.
- Meyn, S. P. and R. L. Tweedie. (1993) *Markov Chains and Stochastic Stability*. NY: Springer.
- McNeill, J. E., Mackulak, G. T. and Fowler, J. W. (2003) Indirect estimation of cycle time quantiles from discrete event simulation models using the Cornish-Fisher expansion. *Proceeding of the 2003 Winter Simulation Conference*, ed. S. Chick, P. J. Sánchez, D. Ferrin, and D. J. Morrice, 1377–1382. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers. Available via (<http://www.informs-cs.org/wsc03papers/173.pdf>)
- Park, S., J. W. Fowler, G. T. Mackulak, J. B. Keats, and W. M. Carlyle. (2002) D-optimal sequential experiments for generating a simulation-based cycle time-throughput curve. *Operations Research* **50**: 981–990.
- Rose, O. (1999) Estimation of the cycle time distribution of a wafer fab by a simple sim-

- ulation model. In *Proceedings of the 1999 International Conference on Semiconductor Manufacturing Operational Modeling and Simulation*, 133–138
- Schömig, A. and J. W. Fowler. (2000). Modelling semiconductor manufacturing operations. *Proceedings of the 9th ASIM Simulation in Production and Logistics Conference*, 55–64. Berlin, Germany.
- Stacy, E. W. (1962) A generalization of the gamma distribution. *Ann. Math. Stat*, **33**, 1187–1192.
- Whitt, W. (1989) Planning queueing simulations. *Management Sciences* **35**: 1341–1366.
- Whitt, W. (2002) Personal communication.
- Yang, F., Ankenman, B. E. and Nelson, B. L. (2007) Efficient generation of cycle time-throughput curves through simulation and metamodeling. *Naval Research Logistics*, **54**, 78–93.

A Appendix

The appendix is organized as follows: In Appendix A.1, we justify the specific forms of the regression models chosen to represent the CT-TH moment curves. In Appendix A.2, we define the notation that will be used in Appendix A.3–Appendix A.6. Appendix A.3 provides the statistical inference on the moment estimators $\{\hat{\mu}_\nu(x) = \mu(x, \hat{\theta}_\nu), \nu = 1, 2, 3\}$. Appendix A.4 derives the normalized variances and covariances of the moment estimators in terms of the design parameters $[\mathbf{x}, \boldsymbol{\pi}]$. The numerical evaluation of the design criterion PM is explained in Appendix A.5, and the method for determining the number of additional replications ΔN is given in Appendix A.6.

A.1 Justification for the Metamodels

This subsection is devoted to justifying the specific form of the metamodels (4) and (5) adopted in the estimation of moment curves. Both models are motivated by some elementary queueing models as well as heavy-traffic analysis of queueing systems.

In Section 3, moment CT-TH curves are represented by Model (4) which consists of two parts: a polynomial function in the numerator, and $(1 - x)^p$ in the denominator. Analytical results for some simple queueing systems have shown that the moment curves for these systems are of the form of (4). For example, the first three moment curves of the M/M/1

queue in steady state are

$$\begin{aligned} \mathbb{E}[CT(x)] &= 1/(1-x) \\ \mathbb{E}[CT(x)^2] &= 2/(1-x)^2 \\ \mathbb{E}[CT(x)^3] &= 6/(1-x)^3 \end{aligned}$$

In Whitt (1989), it has been shown that, at least as $x \rightarrow 1$, moment curves of many queueing systems will follow the form of (4). As explained in Yang et al. (2007), introducing parameter p in Model (4) is necessary because it has been demonstrated that the p value for the underlying CT-TH first-moment (mean) curves could deviate substantially from 1 for real manufacturing systems (Johnson et al. 2004).

Before we discuss the variance model (5), we first define the asymptotic variance of the sample mean $\sum_{h=1}^{H(x)} (CT_h(x))^\nu / H(x)$ as follows:

$$[\sigma_A^{(\nu)}(x)]^2 = \lim_{H(x) \rightarrow \infty} H(x) \text{Var} \left[H(x)^{-1} \sum_{h=1}^{H(x)} (CT_h(x))^\nu \right].$$

We assume that $\{CT_h(x), h = 1, 2, 3, \dots\}$ is a stationary discrete-time stochastic process. Whitt (1989) showed that for the M/M/1 queue, the asymptotic variance of $\sum_{h=1}^{H(x)} CT_h(x) / H(x)$ is given by

$$[\sigma_A^{(1)}(x)]^2 = \frac{x(2 + 5x - 4x^2 + x^3)}{(1-x)^4} \approx \frac{4}{(1-x)^4} \quad (27)$$

as x approaches 1. Under heavy traffic, (28) is consistent with Model (5). Furthermore, it has been shown by Whitt (2002) that the asymptotic variance of higher moments will also have form (5) in heavy traffic:

$$[\sigma_A^{(\nu)}(x)]^2 \approx \frac{\sigma^2}{(1-x)^{2\nu+2}}, \quad \nu = 1, 2, 3, \dots \quad (28)$$

A.2 Notation

We define the notation that will be used for the remainder of this appendix.

$\mathbf{n} = (n(x_1), n(x_2), \dots, n(x_m))$ is the integer vector that represents the allocation of simulation replications to the m design points.

$\mathbf{Z}^{(\nu)} = (Z_1^{(\nu)}(x_1), \dots, Z_{n(x_1)}^{(\nu)}(x_1), \dots, Z_1^{(\nu)}(x_m), \dots, Z_{n(x_m)}^{(\nu)}(x_m))'$ for $\nu = 1, 2, 3$; the transformed data vector consists of independent observations collected at m different throughput levels.

$\mathbf{Z}^{(\nu)}(x) = (Z_1^{(\nu)}(x), Z_2^{(\nu)}(x), \dots, Z_{n(x)}^{(\nu)}(x))'$ is the transformed data vector consists of i.i.d observations collected at throughput x .

$\text{Cov}[\mathbf{Z}^{(k)}(x), \mathbf{Z}^{(\ell)}(x)]$ ($k, \ell = 1, 2, 3$) is a diagonal matrix with identical diagonal elements, $\sigma_{k\ell}(x) \cdot \mathbf{I}_{n(x) \times n(x)}$.

$\text{Cov}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}]$ ($k, \ell = 1, 2, 3, k \neq \ell$) is a diagonal matrix given by:

$$\begin{pmatrix} \sigma_{k\ell}(x_1) \cdot \mathbf{I}_{n(x_1) \times n(x_1)} & & & & \\ & \sigma_{k\ell}(x_2) \cdot \mathbf{I}_{n(x_2) \times n(x_2)} & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \sigma_{k\ell}(x_m) \cdot \mathbf{I}_{n(x_m) \times n(x_m)} \end{pmatrix}. \quad (29)$$

N : the total number of replications.

$\boldsymbol{\theta}_\nu = (\theta_{\nu 1}, \theta_{\nu 2}, \dots, \theta_{\nu T_\nu})$: the unknown parameter vector of the ν^{th} ($\nu = 1, 2, 3$) transformed moment model (7) including the coefficients (c_0, c_1, \dots, c_t) and the exponent parameter r . T_ν is the number of unknown parameters.

$\{\widehat{\mu}_\nu(x) = \mu(x, \widehat{\boldsymbol{\theta}}_\nu); \nu = 1, 2, 3\}$: the estimated moment model (4).

$\{\eta(x, \widehat{\boldsymbol{\theta}}_\nu); \nu = 1, 2, 3\}$: the estimated transformed moment model (7).

$\mathbf{u}(x, \boldsymbol{\theta}_\nu)$, the $T_\nu \times 1$ derivative vector with $\nu = 1, 2, 3$:

$$\begin{aligned} \mathbf{u}(x, \widehat{\boldsymbol{\theta}}_\nu) &= (u_1(x, \widehat{\boldsymbol{\theta}}_\nu), u_2(x, \widehat{\boldsymbol{\theta}}_\nu), \dots, u_{T_\nu}(x, \widehat{\boldsymbol{\theta}}_\nu)) \\ &= \left[\frac{\partial \mu(x, \boldsymbol{\theta}_\nu)}{\partial \theta_{\nu 1}}, \frac{\partial \mu(x, \boldsymbol{\theta}_\nu)}{\partial \theta_{\nu 2}}, \dots, \frac{\partial \mu(x, \boldsymbol{\theta}_\nu)}{\partial \theta_{\nu T_\nu}} \right]' \end{aligned} \quad (30)$$

$\widehat{\mathbf{V}}_\nu$: the $N \times T_\nu$ derivative matrix with elements $\{v_{kj}\}$ defined as

$$v_{kj}^{(\nu)} = v_j^{(\nu)}(x_k) = \frac{\partial \eta(x_k, \boldsymbol{\theta}_\nu)}{\partial \theta_{\nu j}} \Big|_{\boldsymbol{\theta}_\nu = \widehat{\boldsymbol{\theta}}_\nu} \quad k = 1, 2, \dots, N; j = 1, 2, \dots, T_\nu. \quad (31)$$

$\{s_\nu^2; \nu = 1, 2, 3\}$: the residual mean square error resulting from fitting the transformed model (7) for the ν^{th} moment.

A.3 Variance-Covariance Matrix of Moment Estimators

As can be seen from Section 5.2, the variance of the percentile estimator $\widehat{\mathcal{C}}_\alpha(x)$ is a function of the variances and covariances of the moment estimators $\{\widehat{\mu}_\nu(x) = \mu(x, \widehat{\boldsymbol{\theta}}_\nu), \nu = 1, 2, 3\}$.

Consider the random vector of moment estimators $\widehat{\boldsymbol{\mu}}(x) = (\widehat{\mu}_1(x), \widehat{\mu}_2(x), \widehat{\mu}_3(x))'$ as a whole, and its variance-covariance matrix is:

$$\text{Var}[\widehat{\boldsymbol{\mu}}(x)] = \begin{pmatrix} \text{Var}[\widehat{\mu}_1(x)] & \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_2(x)] & \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_3(x)] \\ \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_2(x)] & \text{Var}[\widehat{\mu}_2(x)] & \text{Cov}[\widehat{\mu}_2(x), \widehat{\mu}_3(x)] \\ \text{Cov}[\widehat{\mu}_1(x), \widehat{\mu}_3(x)] & \text{Cov}[\widehat{\mu}_2(x), \widehat{\mu}_3(x)] & \text{Var}[\widehat{\mu}_3(x)] \end{pmatrix}. \quad (32)$$

We provide an estimator for each element in the matrix.

A.3.1 Estimation of the Variances

As illustrated in Section 3, the moment estimators $\{\mu(x, \widehat{\boldsymbol{\theta}}_\nu), \nu = 1, 2, 3\}$ are obtained indirectly from the estimated transformed model $\eta(x, \widehat{\boldsymbol{\theta}}_\nu)$:

$$\mu(x, \widehat{\boldsymbol{\theta}}_\nu) = \frac{\eta(x, \widehat{\boldsymbol{\theta}}_\nu)}{(1-x)^{q_\nu}}$$

The fitting of $\eta(x, \widehat{\boldsymbol{\theta}}_\nu)$ is based on the transformed data set $\{\mathbf{Z}^{(\nu)}, \nu = 1, 2, 3\}$. From nonlinear regression analysis (Bates and Watts 1988), $\text{Var}[\widehat{\boldsymbol{\theta}}_\nu]$ is estimated as

$$\begin{aligned} \widehat{\text{Var}}[\widehat{\boldsymbol{\theta}}_\nu] &= \widehat{\text{Var}}[(\widehat{\mathbf{V}}'_\nu \widehat{\mathbf{V}}_\nu)^{-1} \widehat{\mathbf{V}}'_\nu \mathbf{Z}_\nu] \\ &= (\widehat{\mathbf{V}}'_\nu \widehat{\mathbf{V}}_\nu)^{-1} \widehat{\mathbf{V}}'_\nu \text{Var}[\mathbf{Z}_\nu] ((\widehat{\mathbf{V}}'_\nu \widehat{\mathbf{V}}_\nu)^{-1} \widehat{\mathbf{V}}'_\nu)' \\ &= s_\nu^2 (\widehat{\mathbf{V}}'_\nu \widehat{\mathbf{V}}_\nu)^{-1}. \end{aligned} \quad (33)$$

For reasons explained in Yang et al. (2007), we assume that $q_\nu = \widehat{q}_\nu$ is known when making statistical inference on $\mu(x, \widehat{\boldsymbol{\theta}}_\nu)$. Since $\mu(x, \boldsymbol{\theta}_\nu)$ is nonlinear in the parameters $\boldsymbol{\theta}_\nu$, the first-order approximation by the delta method (Lehmann 1999) is used in the estimation of $\text{Var}[\mu(x, \widehat{\boldsymbol{\theta}}_\nu)]$:

$$\begin{aligned} \text{Var}[\mu(x, \widehat{\boldsymbol{\theta}}_\nu)] &\doteq \text{Var}[\mu(x, \boldsymbol{\theta}_\nu) + \mathbf{u}'(x, \boldsymbol{\theta}_\nu)(\widehat{\boldsymbol{\theta}}_\nu - \boldsymbol{\theta}_\nu)] \\ &= \text{Var}[\mathbf{u}'(x, \boldsymbol{\theta}_\nu) \widehat{\boldsymbol{\theta}}_\nu] \\ &= \mathbf{u}'(x, \boldsymbol{\theta}_\nu) \text{Var}[\widehat{\boldsymbol{\theta}}_\nu] \mathbf{u}(x, \boldsymbol{\theta}_\nu). \end{aligned} \quad (34)$$

Plugging (33) into (34) and replace $\boldsymbol{\theta}_\nu$ by $\widehat{\boldsymbol{\theta}}_\nu$, we have:

$$\widehat{\text{Var}}[\mu(x, \widehat{\boldsymbol{\theta}}_\nu)] = s_\nu^2 \mathbf{u}'(x, \widehat{\boldsymbol{\theta}}_\nu) (\widehat{\mathbf{V}}'_\nu \widehat{\mathbf{V}}_\nu)^{-1} \mathbf{u}(x, \widehat{\boldsymbol{\theta}}_\nu). \quad (35)$$

A.3.2 Estimation of the Covariances

A similar development applies to estimating $\{\text{Cov}[\mu(x, \hat{\boldsymbol{\theta}}_k), \mu(x, \hat{\boldsymbol{\theta}}_\ell)]; k, \ell = 1, 2, 3, k \neq \ell\}$. By the first order approximation of the delta method,

$$\begin{aligned} \text{Cov}[\hat{\mu}(x, \hat{\boldsymbol{\theta}}_k), \hat{\mu}(x, \hat{\boldsymbol{\theta}}_\ell)] &\doteq \text{Cov}[\mu(x, \boldsymbol{\theta}_k) + \mathbf{u}'(x, \boldsymbol{\theta}_k)(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k), \mu(x, \boldsymbol{\theta}_\ell) + \mathbf{u}'(x, \boldsymbol{\theta}_\ell)(\hat{\boldsymbol{\theta}}_\ell - \boldsymbol{\theta}_\ell)] \\ &= \text{Cov}[\mathbf{u}'(x, \boldsymbol{\theta}_k)\hat{\boldsymbol{\theta}}_k, \mathbf{u}'(x, \boldsymbol{\theta}_\ell)\hat{\boldsymbol{\theta}}_\ell] \\ &= \mathbf{u}'(x, \boldsymbol{\theta}_k)\text{Cov}[\hat{\boldsymbol{\theta}}_k, \hat{\boldsymbol{\theta}}_\ell]\mathbf{u}(x, \boldsymbol{\theta}_\ell) \end{aligned} \quad (36)$$

And from nonlinear regression analysis of $\eta(x, \hat{\boldsymbol{\theta}}_\nu)$, we have

$$\begin{aligned} \widehat{\text{Cov}}[\hat{\boldsymbol{\theta}}_k, \hat{\boldsymbol{\theta}}_\ell] &= \widehat{\text{Cov}}[(\widehat{\mathbf{V}}'_k \widehat{\mathbf{V}}_k)^{-1} \widehat{\mathbf{V}}'_k \mathbf{Z}^{(k)}, (\widehat{\mathbf{V}}'_\ell \widehat{\mathbf{V}}_\ell)^{-1} \widehat{\mathbf{V}}'_\ell \mathbf{Z}^{(\ell)}] \\ &= (\widehat{\mathbf{V}}'_k \widehat{\mathbf{V}}_k)^{-1} \widehat{\mathbf{V}}'_k \widehat{\text{Cov}}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}] ((\widehat{\mathbf{V}}'_\ell \widehat{\mathbf{V}}_\ell)^{-1} \widehat{\mathbf{V}}'_\ell)'. \end{aligned} \quad (37)$$

Combining (36) and (37), and substituting the parameter estimates $\{\hat{\boldsymbol{\theta}}_\nu, \nu = 1, 2, 3\}$, we have:

$$\begin{aligned} \widehat{\text{Cov}}[\mu(x, \hat{\boldsymbol{\theta}}_k), \mu(x, \hat{\boldsymbol{\theta}}_\ell)] &= \mathbf{u}'(x, \hat{\boldsymbol{\theta}}_k) (\widehat{\mathbf{V}}'_k \widehat{\mathbf{V}}_k)^{-1} \widehat{\mathbf{V}}'_k \widehat{\text{Cov}}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}] (\widehat{\mathbf{V}}_\ell (\widehat{\mathbf{V}}'_\ell \widehat{\mathbf{V}}_\ell)^{-1}) \mathbf{u}(x, \hat{\boldsymbol{\theta}}_\ell) \\ &k, \ell = 1, 2, 3; \quad k \neq \ell. \end{aligned} \quad (38)$$

The only part left unsolved in (38) is the estimation of $\text{Cov}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}]$, which is a diagonal matrix given by (29). The natural estimator of $\sigma_{k\ell}(x)$, the $k\ell$ element of the covariance matrix $\text{Cov}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}]$, is the sample covariance estimated at x :

$$\begin{aligned} \tilde{\sigma}_{k\ell}(x) &= \frac{1}{n(x) - 1} \sum_{j=1}^{n(x)} (Z_j^{(k)}(x) - \bar{Z}^{(k)}(x))(Z_j^{(\ell)}(x) - \bar{Z}^{(\ell)}(x)) \\ &k, \ell = 1, 2, 3; \quad k \neq \ell \end{aligned} \quad (39)$$

where $\bar{Z}^{(\nu)}(x) = n(x)^{-1} \sum_{j=1}^{n(x)} Z_j^{(\nu)}(x)$ is the sample mean.

However, using the sample covariance (39) to estimate $\text{Cov}[\mu(x, \hat{\boldsymbol{\theta}}_k), \mu(x, \hat{\boldsymbol{\theta}}_\ell)]$ may lead to a covariance matrix that is not positive semi-definite, which is a fundamental requirement on the variance-covariance matrix of any random vector. Therefore, we propose estimating $\sigma_{k\ell}(x)$ in the following way, which provides some consistency in how we estimate the variances and covariances, and preserves the positive semi-definite property of the estimated matrix (32). We estimate the sample correlation

$$\tilde{\varrho}_{k\ell}(x) = \frac{\tilde{\sigma}_{k\ell}(x)}{\sqrt{\tilde{\sigma}_k^2(x) \tilde{\sigma}_\ell^2(x)}} \quad (40)$$

where $\tilde{\sigma}_{k\ell}(x)$ is the sample covariance (39), and $\hat{\sigma}_\nu^2(x) = (n(x)-1)^{-1} \sum_{j=1}^{n(x)} (Z_j^{(\nu)}(x) - \bar{Z}^{(\nu)}(x))^2$ is the sample variance. Then an estimate of $\sigma_{ij}(x)$ can be obtained utilizing the residual error s_ν^2 which is also used in (33) for estimating $\text{Var}[\mu(x, \hat{\boldsymbol{\theta}}_\nu)]$:

$$\hat{\sigma}_{ij}(x) = \tilde{\varrho}_{ij}(x) \times s_i s_j. \quad (41)$$

Notice that the estimation discussed in the subsection is used either for making statistical inference on the percentile estimator or numerically evaluating the design criterion in search of the optimal design. The case worthy of mention is when the current design is to be expanded by including more design points. At those additional points, no simulation data have yet been collected, and we cannot directly obtain the corresponding covariance estimates $\hat{\sigma}_{ij}(x)$ for the new design points. A simple approximation has been adopted in our procedure: we assign $\hat{\sigma}_{ij}(x) = (\hat{\sigma}_{ij}(x_L) + \hat{\sigma}_{ij}(x_U))/2$ for those new design points to be included.

A.4 Normalized Variance and Covariance Derivation for Moment Estimators

In this section, we derive the normalized variances ($N\text{Var}[\cdot]$) and covariances ($N\text{Cov}[\cdot]$) of the moment estimators at a particular throughput rate $x \in [x_L, x_U]$ in terms of \mathbf{x} and $\boldsymbol{\pi}$. This provides the basis for evaluating the normalized variances of percentile estimators, which is critical to the design of simulation experiments as shown later in Appendices A.5 and A.6.

The estimation is based on N replications allocated to m design points $\mathbf{x} = (x_1, x_2, \dots, x_m)$ according to the loadings $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_m)$. In the following derivations, we assume that the allocation of replications $n_i = N\pi_i$ ($i = 1, 2, \dots, m$) is an integer.

The estimated variance of $\hat{\mu}_\nu(x)$ is given by (35) as

$$\widehat{\text{Var}}[\hat{\mu}_\nu(x)] = s_\nu^2 \mathbf{u}'(x, \hat{\boldsymbol{\theta}}_\nu) (\widehat{\mathbf{V}}_\nu' \widehat{\mathbf{V}}_\nu)^{-1} \mathbf{u}(x, \hat{\boldsymbol{\theta}}_\nu). \quad (42)$$

Let the matrix $\mathbf{A}^{(\nu)} = \widehat{\mathbf{V}}_\nu' \widehat{\mathbf{V}}_\nu$, with element $a_{st}^{(\nu)} = \sum_{i=1}^m n_i v_s^{(\nu)}(x_i) v_t^{(\nu)}(x_i)$. We can write the determinant of $\mathbf{A}^{(\nu)}$ as

$$|\mathbf{A}^{(\nu)}| = \sum_{i_1=1}^m n_{i_1} v_1^{(\nu)}(x_{i_1}) \sum_{i_2=1}^m n_{i_2} v_2^{(\nu)}(x_{i_2}) \cdots \sum_{i_T=1}^m n_{i_T} v_T^{(\nu)}(x_{i_T}) \begin{vmatrix} v_1^{(\nu)}(x_{i_1}) & v_1^{(\nu)}(x_{i_2}) & \cdots & v_1^{(\nu)}(x_{i_T}) \\ v_2^{(\nu)}(x_{i_1}) & v_2^{(\nu)}(x_{i_2}) & \cdots & v_2^{(\nu)}(x_{i_T}) \\ \vdots & \vdots & \ddots & \vdots \\ v_T^{(\nu)}(x_{i_1}) & v_T^{(\nu)}(x_{i_2}) & \cdots & v_T^{(\nu)}(x_{i_T}) \end{vmatrix}.$$

Then $|\mathbf{A}^{(\nu)}|$, the determinant of $\mathbf{A}^{(\nu)}$, and $|\mathbf{A}_{st}^{(\nu)}|$, the cofactor of $\mathbf{A}^{(\nu)}$ can be written as

follows:

$$\begin{aligned}
|\mathbf{A}^{(\nu)}| &= \sum_{\substack{i_s \neq i_t \\ 1 \leq i_1, \dots, i_T \leq m}} n_{i_1} n_{i_2} \cdots n_{i_T} v_1^{(\nu)}(x_{i_1}) v_2^{(\nu)}(x_{i_2}) \cdots v_T^{(\nu)}(x_{i_T}) \\
&\times \begin{vmatrix} v_1^{(\nu)}(x_{i_1}) & v_1^{(\nu)}(x_{i_2}) & \cdots & v_1^{(\nu)}(x_{i_T}) \\ v_2^{(\nu)}(x_{i_1}) & v_2^{(\nu)}(x_{i_2}) & \cdots & v_2^{(\nu)}(x_{i_T}) \\ \vdots & \vdots & \ddots & \vdots \\ v_T^{(\nu)}(x_{i_1}) & v_T^{(\nu)}(x_{i_2}) & \cdots & v_T^{(\nu)}(x_{i_T}) \end{vmatrix} \\
&= N^{T\nu} b^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})
\end{aligned} \tag{43}$$

$$\begin{aligned}
|\mathbf{A}_{st}^{(\nu)}| &= \sum_{\substack{i_a \neq i_b \\ i_a \neq i_t \\ 1 \leq i_1, \dots, i_T \leq m}} n_{i_1} \cdots n_{i_{t-1}} n_{i_{t+1}} \cdots n_{i_T} v_1^{(\nu)}(x_{i_1}) \cdots v_{t-1}^{(\nu)}(x_{i_{t-1}}) v_{t+1}^{(\nu)}(x_{i_{t+1}}) \cdots v_T^{(\nu)}(x_{i_T}) \\
&\times \begin{vmatrix} v_1^{(\nu)}(x_{i_1}) & \cdots & v_1^{(\nu)}(x_{i_{t-1}}) & v_1^{(\nu)}(x_{i_{t+1}}) & \cdots & v_1^{(\nu)}(x_{i_T}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{s-1}^{(\nu)}(x_{i_{s-1}}) & \cdots & v_{s-1}^{(\nu)}(x_{i_{t-1}}) & v_{s-1}^{(\nu)}(x_{i_{t+1}}) & \cdots & v_{s-1}^{(\nu)}(x_{i_T}) \\ v_{s+1}^{(\nu)}(x_{i_{s+1}}) & \cdots & v_{s+1}^{(\nu)}(x_{i_{t-1}}) & v_{s+1}^{(\nu)}(x_{i_{t+1}}) & \cdots & v_{s+1}^{(\nu)}(x_{i_T}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ v_T^{(\nu)}(x_{i_1}) & \cdots & v_T^{(\nu)}(x_{i_{t-1}}) & v_T^{(\nu)}(x_{i_{t+1}}) & \cdots & v_T^{(\nu)}(x_{i_T}) \end{vmatrix} \\
&= N^{T-1} b_{st}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})
\end{aligned} \tag{44}$$

Equations (43) and (44) follow because $n_i = N\pi_i$. The functions $b^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})$ and $b_{st}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})$ can be evaluated given \mathbf{x} and $\boldsymbol{\pi}$. Then the inverse of $\mathbf{A}_{st}^{(\nu)}$ can be obtained as

$$[(\mathbf{A}^{(\nu)})^{-1}]_{st} = \frac{|\mathbf{A}_{st}^{(\nu)}|}{|\mathbf{A}^{(\nu)}|} = N^{-1} \frac{b_{st}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})}{b^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})}. \tag{45}$$

Thus, $(\mathbf{A}^{(\nu)})^{-1} = N^{-1} \mathbf{B}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})$ where the (s, t) element of $\mathbf{B}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})$ is given by $\mathbf{B}_{st}^{(\nu)} = b_{st}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})/b^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})$. The normalized variance of the ν^{th} moment estimator $\widehat{\mu}_\nu(x_0)$ can be expressed as:

$$N \times \widehat{\text{Var}}[\widehat{\mu}_\nu(x_0)] = s_\nu^2 \mathbf{u}'(x, \widehat{\boldsymbol{\theta}}_\nu) \mathbf{B}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi}) \mathbf{u}(x, \widehat{\boldsymbol{\theta}}_\nu). \tag{46}$$

The covariances of the moment estimators are given in (38) as:

$$\begin{aligned} & \widehat{\text{Cov}}[\widehat{\mu}_k(x), \widehat{\mu}_\ell(x)] \\ &= \mathbf{u}'(x, \widehat{\boldsymbol{\theta}}_k) (\widehat{\mathbf{V}}_k' \widehat{\mathbf{V}}_k)^{-1} \widehat{\mathbf{V}}_k' \widehat{\text{Cov}}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}] (\widehat{\mathbf{V}}_\ell (\widehat{\mathbf{V}}_\ell' \widehat{\mathbf{V}}_\ell)^{-1})' \mathbf{u}(x, \widehat{\boldsymbol{\theta}}_\ell) \end{aligned} \quad (47)$$

$k, \ell = 1, 2, 3; \quad k \neq \ell$

The covariance matrix $\text{Cov}[\mathbf{Z}^{(k)}, \mathbf{Z}^{(\ell)}]$ is given by (29), the elements of which $\sigma_{k\ell}(x)$ ($x = x_1, x_2, \dots, x_m$) can be estimated by (41) as explained in Appendix A.3.

Using (45), we can easily derive

$$\begin{aligned} & N \widehat{\text{Cov}}[\widehat{\mu}_k(x), \widehat{\mu}_\ell(x)] \\ &= \left(\sum_{i=1}^m \pi_i \widehat{\sigma}_{k\ell}(x_i) \left(\sum_{j=1}^{T_k} u_j(x, \widehat{\boldsymbol{\theta}}_k) d_{ji}^{(k)} \right) \left(\sum_{h=1}^{T_\ell} u_h(x, \widehat{\boldsymbol{\theta}}_\ell) d_{hi}^{(\ell)} \right) \right) \end{aligned} \quad (48)$$

where

$$d_{ji}^{(\nu)} = \sum_{t=1}^{T_\nu} \frac{b_{jt}^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})}{b^{(\nu)}(\mathbf{x}, \boldsymbol{\pi})} \times v_t^{(\nu)}(x_i), \quad (49)$$

and $u_j(x, \widehat{\boldsymbol{\theta}}_\ell)$ is defined in (30).

A.5 Numerical Evaluation of the Design Criterion

For convenience, we restate the design criterion in (22)

$$PM = N \sum_{\kappa \in \boldsymbol{\Omega}_x} \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)]. \quad (50)$$

In this section, we evaluate the design criterion numerically for given design $[\mathbf{x}, \boldsymbol{\pi}]$ in an effort to search for the optimal solution of (23). Recall that this evaluation is performed in the context where some simulation experiments have already been carried out. The key is to obtain the variances and covariances of the moment estimators $\{\widehat{\mu}_\nu(\kappa), \nu = 1, 2, 3\}$ for $\kappa \in \boldsymbol{\Omega}_x$, from which $\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)]$ can be approximately computed.

As we did in Yang et al. (2007), linear assumptions about the regression models for moment curves are made at the stage of designing experiments. More specifically, the unknown parameters, p_ν in the moment model (4) and q_ν in the variance model (5), are both assumed to be known and given as their best estimates obtained so far, \widehat{p}_ν and \widehat{q}_ν . This assumption reduces the transformed model (7) to a linear model with the exponent r_ν being a known parameter given by $\widehat{q}_\nu - \widehat{p}_\nu$. In addition, we suppose that the expanded experiments with m design points are designed to estimate moment models (4) whose polynomial order is equal to $m - 1$.

Under these assumptions, the moment models and transformed moment models are (with superscript ν dropped for simplicity):

$$\mu(x, \boldsymbol{\theta}) = \frac{\sum_{\ell=0}^{m-1} c_{\ell} x^{\ell}}{(1-x)^p} \quad (51)$$

$$\eta(x, \boldsymbol{\theta}) = \left(\sum_{\ell=0}^{m-1} c_{\ell} x^{\ell} \right) (1-x)^r \quad (52)$$

where both p and r are assumed to be known, and the unknown linear parameters are $\boldsymbol{\theta} = (c_0, c_1, \dots, c_{m-1})$.

As already noted, the design criterion (50) is a normalized variance measure, which depends on the design $[\mathbf{x}, \boldsymbol{\pi}]$. To search for the optimal $[\mathbf{x}, \boldsymbol{\pi}]$, we need to be able to evaluate (50) numerically for given $[\mathbf{x}, \boldsymbol{\pi}]$. With $\mu(x, \boldsymbol{\theta})$ and $\eta(x, \boldsymbol{\theta})$ defined as the linear models given by (51) and (52), we can numerically evaluate the normalized variances and covariances of the moment estimators for given $[\mathbf{x}, \boldsymbol{\pi}]$

$$N\text{Var}[\mu(\kappa, \widehat{\boldsymbol{\theta}}_{\nu})] \quad \nu = 1, 2, 3 \quad (53)$$

$$\text{and } N\text{Cov}[\mu(\kappa, \widehat{\boldsymbol{\theta}}_k), \mu(\kappa, \widehat{\boldsymbol{\theta}}_{\ell})] \quad k, \ell = 1, 2, 3; k \neq \ell, \quad (54)$$

by following the derivation in Appendix A.4.

In addition, based on the data set already collected, $\{\mu(\kappa, \widehat{\boldsymbol{\theta}}_{\nu}); \kappa \in \boldsymbol{\Omega}_x, \nu = 1, 2, 3\}$, the distribution parameters of $G_{\kappa}(t; \widehat{a}(\kappa), \widehat{b}(\kappa), \widehat{k}(\kappa))$ can be estimated. Applying the delta method as illustrated in Section 5.2, we can approximately calculate $N\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)]$, which is a component of the sum formulation (50).

A.6 Incrementing the Number of Runs

As mentioned in Section 6.1.2, we need to determine ΔN , the number of replications to be added to the current design before solving the optimization problem (23) to obtain the optimal design for further simulation experiments. As explained in Yang et al. (2007), adding ΔN replications should help drive the precision level of the chosen estimators, in this case $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$ and $\widehat{\mathcal{C}}_{\alpha_L}(x_L)$, down to the target $100\gamma\%$ (the definition of $100\gamma\%$ is given in Section 6.2). Conditional on the size of the current estimated relative error $100\widehat{\gamma}_c$, which indicates how precise the current estimate is, we set an intermediate target, say $100\gamma_1\%$ ($\gamma \leq \gamma_1 < \gamma_c$), for the expanded data set to achieve. Details regarding the choice of $100\gamma_1\%$ are discussed in Yang et al. (2004, Appendix 3).

Regarding the determination of ΔN , we follow the same logic as in Appendix 3 of Yang et al. (2007). In the remainder of this section, we illustrate how a rough estimate of ΔN is computed for achieving the desired precision level $100\gamma_1\%$ under the following assumptions:

- (i) To obtain a rough estimate of ΔN needed for the follow-up design given the current data, the design points \mathbf{x} to be included in the follow-up design have to be known. As noted in Section 6.2, there are two types of design augmentation involved in our procedure: adding more replications to the current design points; and expanding the current design by adding new design points to the initial design composed of the two points x_L and x_U . In the first case, the location of the design points is known. In the second case, we simply assume that the design points in the augmented design will be evenly spaced throughout the throughput range.
- (ii) The number of unknown parameters included in the estimated moment models $\{\widehat{\mu}_\nu(x), \nu = 1, 2, 3\}$ remains the same after the design augmentation.
- (iii) The allocation of the total $N = N_c + \Delta N$ replications is not constrained by the fact that the N_c replications already performed cannot be reassigned.

The process for obtaining an estimate of N consists of two steps.

1. With \mathbf{x} being given in Assumption (i), solve the following optimization problem for the optimal solution $\boldsymbol{\pi}^*$:

$$\begin{aligned} \min_{\boldsymbol{\pi}} PM &= N \sum_{\kappa \in \Omega_x} \text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(\kappa)] \\ \text{s.t. } \sum_{i=1}^m \pi_i &= 1. \end{aligned} \quad (55)$$

The objective PM is a function of $\boldsymbol{\pi}$ (independent of N) given the location of the design points \mathbf{x} . Under the assumptions above, the performance measure PM can be evaluated for given $\boldsymbol{\pi}$ as illustrated in Appendix A.5. Thus a numerical search will lead to the optimal solution $\boldsymbol{\pi}^*$.

2. Assuming the follow-up design is $[\mathbf{x}, \boldsymbol{\pi}^*]$, find the number of replications N that satisfy

$$\frac{2 \times \widehat{\text{SE}}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)]}{\widehat{\mathcal{C}}_{\alpha_U}(x_U)} \leq 100\gamma_1\%. \quad (56)$$

Based on the current data set, an estimate of $\mathcal{C}_{\alpha_U}(x_U)$ is obviously available. Given design $[\mathbf{x}, \boldsymbol{\pi}^*]$, the normalized variance $N\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)]$ can also be estimated: first, we obtain $\{N\text{Var}[\widehat{\mu}_\nu(x_U)], \nu = 1, 2, 3\}$ and $\{NCov[\widehat{\mu}_k(x_U), \widehat{\mu}_\ell(x_U)], k, \ell = 1, 2, 3; k \neq \ell\}$ by following the derivation in Appendix A.4; then we estimate $N\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)]$ by resorting to the delta method. The standard error is $\widehat{\text{SE}} = \sqrt{N^{-1} \times (N\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)])}$ with N

being the unknown variable to be determined (the normalized variance $N\text{Var}[\widehat{\mathcal{C}}_{\alpha_U}(x_U)]$ is independent of N). Therefore, by solving the inequality (56), we obtain N_{min} , the smallest number of replications required to satisfy (56). Then $\Delta N = N_{min} - N_c$.

Once the desired relative precision has been achieved on $\widehat{\mathcal{C}}_{\alpha_U}(x_U)$, we will switch to use the precision of $\widehat{\mathcal{C}}_{\alpha_U}(x_L)$ to drive the procedure until $\widehat{\mathcal{C}}_{\alpha_U}(x_L)$ also achieves the prespecified precision level. The derivation for determining ΔN presented above can be applied directly to $\widehat{\mathcal{C}}_{\alpha_U}(x_L)$.