

1981

Progress Report Operational Analysis of Queueing Phenomena

Peter J. Denning

Jeffrey P. Buzen

Report Number:
81-370

Denning, Peter J. and Buzen, Jeffrey P., "Progress Report Operational Analysis of Queueing Phenomena" (1981). *Department of Computer Science Technical Reports*. Paper 299.
<https://docs.lib.purdue.edu/cstech/299>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

PROGRESS REPORT
OPERATIONAL ANALYSIS OF QUEUEING PHENOMENA

Peter J. Denning

Jeffrey P. Buzen

For NSF Grant MCS78-01729

University/Industry Cooperative Program

at

Purdue University

W. Lafayette, IN

and

BGS Systems, Inc.

Waltham, MA

CSD-TR 370

June 1, 1981

Introduction

Operational analysis is a new method of analyzing queueing phenomena, especially the phenomena encountered by analysts of computer system performance. Operational analysis deals with *data* collected from a system rather than the system itself. The data is a "behavior sequence" -- an event trace of system states during an observation period. Operational variables denote quantities measurable in a behavior sequence. Operational assumptions are testable in a behavior sequence.

With operational analysis, we have shown that many of the results known from stochastic analysis of queueing systems are more general (and more robust) than has been believed. This is because these results can be applied validly to behavior sequences that would occur with probability zero in the stochastic ensembles for which we can carry out analyses. For example, operational analysis has shown that the product form solution for closed queueing networks may hold for deterministic systems, which helps explain why these models work so well when applied by industrial engineers to job flow in factories.

Operational analysis has, more importantly, opened up new lines of research into the robustness of models. We believe that the most important instance of this is the *error analysis of assumptions*. It is possible to measure the error between the value of a state-dependent parameter observed to be state-dependent and the value that parameter is asserted to have by a modeling assumption, and then to express the consequent error between the model's estimate of a performance metric and the true value found in the data. Suri has recently completed an error analysis to evaluate the sensitivity of the product form solution to the assumption that the service times of devices are load-independent;¹ he has found that model estimates of utilizations and

¹ R. Suri, "Robustness of analytical formulae for performance prediction in certain nonclassical queueing networks," Technical Report No. 874, Division of Applied Sciences, Harvard University, Cambridge, MA 02138 (August 1980), 88pp.

mean queue lengths are not sensitive to this assumption. We believe this line of work is extremely important. It should eventually be coupled with an extensive experimental study to determine the class of physical systems for which variations in parameters (relative to assumed values) are within the ranges required for queueing network models to be accurate.

This progress report summarizes our work on operational analysis in this project, which is funded for cooperative work between Purdue and BGS Systems, Inc. The material is divided into four parts. First is a summary of work completed since the project was funded in 1978. Second is a summary of work now in progress and expected to be completed during 1981-82. Third is an annotated bibliography of papers that have been published. Fourth is an annotated bibliography of technical reports available from the principal investigators. Indicators in the text of the form [*i*] refer to the bibliography; references to the literature are in footnotes.

At Purdue, three students are assigned to the project: Jeff Brumfield, who is expected to complete his PhD thesis by December 1981; Subhash Agrawal, who passed PhD qualifying examinations in Spring 1981; and Andre Bondi, who passed his PhD qualifying examinations in fall 1980. Under the terms of the grant, Agrawal and Bondi are spending summer of 1981 on site at BGS systems to work on projects of the grant.

Work Completed

Prior to 1978, most of our research focussed on operational laws, bottleneck analysis, and the product form solution. A comprehensive review of this work was published in the paper by Denning & Buzen in the special issue of the *Computing Surveys* in September 1978 [1].

Operational techniques have been applied to the analysis of working set size and swapping demand for segment reference strings in virtual memory systems. This has been discussed in the paper by Denning published in the *IEEE Transactions on Software Engineering* in January 1980. [2]

The *Computing Surveys* paper focussed on queueing networks. We have also studied single-resource queues, under the assumption that only the arrivals, completions, and queue size are externally observable. (The queueing discipline, number of servers, and service periods are not externally visible.) We unified the heretofore somewhat obscure results on the relations among the arriver's distribution, the completer's distribution, and outside observer's distribution. We extended the techniques to obtain an operational analog of the Sevcik-Mitrani arrival theorem in queueing networks,² which leads to an operational formulation of mean value analysis. This work was reported in the papers by Buzen & Denning [3, 4].

In answer to a challenge raised by critics of operational analysis, Brumfield & Denning undertook to find an operational counterpart of the Pollaczek-Khintchin formula for the mean queue length in a single-server queue. [11] The formula is obtained after applying a three-part homogeneity assumption that asserts the independence of the arrivals from queue length, length of service period, and average position in service periods. The formula is

$$\bar{n} = U \frac{1-U-Np(N)}{1-U-p(N)} + \frac{U(U-p(N))(CV^2+1)}{2(1-U-p(N))}$$

where U is the utilization, CV^2 is the squared coefficient of variation $\left(\frac{\sigma^2}{\bar{s}^2}\right)$ of service

² K. C. Sevcik and I. Mitrani, "The distribution of queueing network states at input and output instants," *J. ACM* 28, 2 (April 1981), 358-371. The theorem says that the queue length distribution seen by the arrivers to device i is the same as the overall distribution seen with one customer (the arriver) removed from the network.

periods, and $p(N)$ is the proportion of time the queue attains its maximum size (N jobs). If $p(N) = 0$, this reduces to the familiar form

$$\bar{n} = U + \frac{U^2(CV^2+1)}{2(1-U)}$$

Unlike the earlier analysis of single queues, this analysis requires that the service periods of jobs at the device be externally observable.

In a side project we applied simple principles of performance evaluation to specify the design of an efficient hardware-based process manager for a multiprocessor computer system. The resulting work by Denning, Dennis, and Brumfield will be published in 1981 in the *Communications of ACM*. [5]

In a keynote speech before the SIGMETRICS conference on modeling and measurement of computer systems (September 1981), Denning cites the IBM M44/44X project and the development of queueing network models as examples of experimental computer science at its best. [6]

Work in Progress

PROJECT 1: Unified Operational Theory of Queueing Systems (*J. Brumfield & P. J. Denning at Purdue*). For his PhD dissertation, Brumfield is collecting into one place all the major operational queueing results. This will include a review of the single-class queueing network results, the M/G/1 queue counterpart result [11], and new results for multiclass networks. The new results include the homogeneity assumptions needed to obtain the analog of the Baskett-Chandy-Muntz-Palacios theorem.³ There

³ F. Baskett, K. M. Chandy, R. R. Muntz, and F. G. Palacios, "Open, closed, and mixed networks with different classes of customers," *J. ACM* 22, 2 (April 1975), 248-268.

appear to be three different homogeneity assumptions possible for each device: the service function can depend on both class (τ) and local queue length by class ($n_{i\tau}$); it can be independent of class and depend only on the total queue length ($\sum_{\tau} n_{i\tau}$); or it can be independent of class and queue length. These possibilities correspond to the principal cases of the response time equation in mean-value analysis. The new results also include bottleneck analysis in multiclass networks.

PROJECT 2: Operational Sensitivity Analyses (*A. Bondi & P. J. Denning at Purdue*).

As noted above, Suri has completed a sensitivity analysis of the load-independent product form solution to errors in the assumption that the mean service times actually are independent of load. The objective of this project is to extend this analysis to determine the sensitivity of the product form to the basic homogeneity of the queueing network: that the job flow from device i to device j depends only on the queue length n_i at the source of the flow. Such results are a prerequisite to a complete experimental study to circumscribe the class of systems for which the product form gives accurate results.

PROJECT 3: Operational Renewal Theory (*A. Bondi and P. J. Denning at Purdue*).

In his PhD thesis (1968), Denning derived equations for the mean working set size on the hypothesis that successive references to the same segment formed a renewal process; in this case, the interreference interval distribution was the recurrence distribution and the mean working set size was approximated by:

$$W(\tau) = \int_0^{\tau} (1-F(u))du,$$

where τ is the window size and $1-F(x)$ is the tail of the cumulative distribution of the interreference distribution. This result is closely related to the backward recurrence

problem in renewal theory. In 1974, Slutz and Traiger showed that a discrete form of this result could be derived without any renewal theory:⁴

$$\bar{w}(\tau) = \sum_{i=0}^{\tau-1} \sum_{k>i} f(k) ,$$

where $f(k)$ is the interreference distribution. We conjecture that this is an instance of a more general result. We believe that the limit equations of renewal theory have operational counterparts that do not depend on a concept of "forgetfulness" at recurrence times. The objective of this project is to work out the details.

PROJECT 4: Aggregate Server Models for Critical Section Delays (*J. Buzen and S. Agrawal at BGS Systems and P. J. Denning at Purdue*). Buzen, Liu, and Shum have recently developed the technique of "aggregate servers" for representing the delays in software systems that arise because of contention for critical sections. This project will investigate the accuracy of this technique by comparing it with exact solutions in special cases where these solutions are available, and by comparing with simulation results in other cases. Simplifications and alternative formulations of the basic technique will also be compared for relative accuracy.

PROJECT 5: Operational tests for Markovian Properties (*J. Bouhana and R. Bryant at U. Wisconsin, J. Buzen at BGS Systems, and P. J. Denning at Purdue*). A set of measurement data is assumed to have been generated by a stochastic process. By testing the data, we would like to draw inferences about the underlying stochastic process. We are particularly interested in inferences Poisson arrival processes and exponential service time distributions.

⁴ D. R. Slutz and I. W. Traiger, "A note on the calculation of average working set size," *Communications of ACM* 17, 10 (October 1974), 583-585.

A number of tests for these stochastic assumptions already exist. To the best of our knowledge, however, none of these tests makes use of the operational concepts of homogeneous arrivals and homogeneous services [3,4]. We propose to investigate how these operational concepts can provide new, potentially superior, tests for Poisson arrivals or exponential service times.

Bryant has already demonstrated that service times must, with probability one, be exponential if any underlying M/G/1 stochastic process generates an infinitely long behavior sequence that satisfies service time homogeneity exactly.⁵ We shall consider generalizations of this result to G/M/1 queues and, possibly, to G/G/1 queues. We shall also consider the conclusions about confidence levels that can be drawn if the behavior sequence is of finite length, and if operational homogeneity is only satisfied to within some ϵ greater than zero.

PROJECT 6: Operational M/G/1 queues (A. Bondi and J. Buzen at BGS Systems).

We have worked out a derivation of the queue length distribution for a single server queueing system. This distribution is expressed in terms of the average arrival rate and the distribution of service times. The analysis is entirely operational: distributions are interpreted as histograms of observed values, independent and dependent variables are defined in terms of directly measurable quantities, and all assumptions are expressed in terms of relationships among observables. Results analogous to those for stochastic M/G/1 queues are obtained.

⁵ R. Bryant, "On homogeneity and online = offline behavior in M/G/1 queueing systems," *Proc. Performance '80 Conference*, ACM SIGMETRICS (May 1980), 199-208.

PROJECT 7: Operational Priority Queues (A. Bondi and J. Buzen at BGS Systems).

Consider an $M/G/n$ queueing system with a preemptive-resume priority scheduler.

Assume customers at each priority level satisfy the assumptions of Poisson arrivals and exponential service times. Expressions for the mean queue lengths and mean response times are derived for cases where the mean service times are identical for all priority levels. Approximations for these quantities are derived for cases where the mean service times are not identical. In both cases, the operational counterparts of the stochastic queueing models are derived.

PAPERS

1. Peter J. Denning and Jeffrey P. Buzen, "The operational analysis of queueing network models," *Computing Surveys* 10, 3 (September 1978), 227-261.

Queueing network models have proved to be cost effective tools for analyzing modern computer systems. This tutorial paper presents the basic results using the operational approach, a framework which allows the analyst to test whether each assumption is met in a given system. The early sections describe the nature of queueing network models and their applications for calculating and predicting performance quantities. The basic performance quantities — such as utilizations, mean queue lengths, and mean response times — are defined, and operational relationships among them are derived. Following this, the concept of job flow balance is introduced and used to study asymptotic throughputs and response times. The concepts of state transition balance, one-step behavior, and homogeneity are then used to relate the proportions of time that each system state is occupied to the parameters of job demand and to device characteristics. Efficient methods for computing basic performance quantities are also described. Finally the concept of decomposition is used to simplify analyses by replacing subsystems with equivalent devices. All concepts are illustrated liberally with examples.

2. Peter J. Denning, "Working sets past and present," *IEEE Transactions on Software Engineering* SE-6, 1 (January 1980), 44-84.

A program's working set is the collection of segments (or pages) recently referenced. This concept has led to efficient methods for measuring a program's intrinsic memory demand; it has assisted in understanding and in modeling program behavior; and it has been used as the basis of optimal multiprogrammed memory management. The total cost of a working set dispatcher is no larger than the total cost of other common dispatchers. This paper outlines the argument why it is unlikely that anyone will find a cheaper nonlookahead memory policy that delivers significantly better performance.

3. Jeffrey P. Buzen and Peter J. Denning, "Measuring and calculating queue length distributions," *IEEE Computer* (April 1980), 33-44.

This paper continues the tutorial of the *Computing Surveys* paper [1] by focussing on the analysis of a single queue. By measuring arrivals, completions, and holding times for each possible queue length, one can calculate the arriver's, completer's, and overall (outside observer's) distribution. Arrivals are homogeneous in a behavior sequence if the arrival counts are equal for all queue lengths. Services are homogeneous in a behavior sequence if the completion rates are equal for all queue lengths. The basic results are three: 1) the arriver's distribution and completer's distribution are identical for any flow balanced behavior sequence; 2) the arriver's distribution is of the same form as the overall distribution for any behavior sequence with homogeneous arrivals; and 3) the response time is the mean service time multiplied by the mean length of the queue just after an arrival for any behavior sequence with homogeneous services. These results are extended to demonstrate the arrival theorem for closed queueing networks: the queue length distribution seen by an arriver is identical to the overall distribution with one less job (the arriver) in the network. A simple consequence of the arrival theorem is the set of equations for mean value analysis. The concepts and methods are well illustrated with examples. One example demonstrates an operational sensitivity analysis, showing the error introduced by deviations from the homogeneous arrivals assumption.

4. Jeffrey P. Buzen and Peter J. Denning, "Operational treatment of queue distributions and mean value analysis," *Computer Performance* 1, 1 (June 1980), 8-15.

This paper contains the formal derivations of the results in the *Computer* paper [3]. The relations among the queue length distributions seen by an arriving job, a completing job, and an outside observer are derived using operational analysis. A simplified derivation of the Sevcik-Mitrani arrival theorem is presented and used as the basis for discussing the Reiser and Lavenberg mean value analysis. Two results are presented: an algorithm for computing queue-length distributions from conditional throughputs in closed, product-form queueing networks; and an operational bound on the errors that can rise in certain theorems when homogeneity is violated.

5. Peter J. Denning, T. Don Dennis, and Jeffrey A. Brumfield, "Low contention semaphores and ready lists," Purdue University, Computer Sciences Department, Technical Report CSD-TR-332 (April 1981), accepted for publication in *Communications of ACM*.

A method for reducing semaphore and ready list contention in multiprocessor operating systems is described. Its correctness is established. Its performance is compared with conventional implementations. A ready list implemented as a ring network is proposed and evaluated.

6. Peter J. Denning, "Performance Evaluation: Experimental Computer Science at its Best," *Proc. ACM SIGMETRICS Conference* (September 1981), to appear.

Experimental science classifies knowledge derived from observations. The experimenter sets up an apparatus, uses it to collect data, and analyzes the data to sustain or refute hypotheses. The result of one line of investigation can be a model that becomes the apparatus for another line of investigation. Two examples from the area of performance evaluation are used to illustrate. The M44/44X project at IBM Watson Research Lab in the mid 1980s evaluated concepts of time sharing, especially about memory policies and program behavior, by implementing and measuring them on an IBM 7044. The study of queueing network models since 1971 illustrates how strong interaction between theory and experiment can produce a model (the Bard-Schweitzer mean value equations) that is sufficiently simple to serve as the starting point for new lines of investigation of system models.

TECHNICAL REPORTS

7. Jeffrey P. Buzen, Peter J. Denning, Donald B. Rubin, and Linda S. Wright, "Operational Markov Chains," Technical Report, BGS Systems, Inc. (January 1979).

If a sequence of observations is assumed to have been generated by a Markov chain, various relationships among the elements of the sequence can be derived. In the paper, operational analysis is used to prove that some of these relationships hold under very general conditions. For example, the eigenvalue relation $P = \pi P$ holds for any state sequence in which the initial and final states are the same, provided that $P = (p_1, \dots, p_n)$ is the vector of proportions of time each of the respective states $1, \dots, n$ is occupied and $\pi = [\pi_{ij}]$ is the matrix of observed one-step interstate transition frequencies. The eigenvalue equation is an example of a relation that holds operationally in cases where the conventional Markov chain model could not reasonably be justified. An analysis of the sensitivity of these relationships to these more general operational assumptions is also presented. For example, the error introduced by measuring matrix π and using the normalized solution of $\mathbf{x} = \pi \mathbf{x}$ to estimate the true values of P is of order $O(n/T)$, where T is the length of the observation sequence.

8. Peter J. Denning and Jeffrey P. Buzen, "Questions and answers about operational analysis," Purdue University, Computer Sciences Department, Technical Report CSD-TR-318 (November 1979), 57pp.

This paper addresses a number of criticisms of operational analysis that have been stated over the past few years. It is organized as a series of statements or questions, each followed by a response. Related sets of questions are grouped together and preceded by a short background essay. Here is the table of contents of the report:

- 1 OVERVIEW OF OPERATIONAL ANALYSIS
- 2 SCOPE OF OPERATIONAL ANALYSIS
 - 2.1 *Operational results are tautologies.*
 - 2.2 *Operational analysis cannot handle queueing disciplines.*
 - 2.3 *Operational analysis cannot handle service time distributions.*
 - 2.4 *Operational analysis is not well suited for prediction because it is difficult to determine what the values of $S(n)$ will be after a change is made to the speed of a server or to the arrival rate.*
 - 2.5 *Operational analysis cannot deal with cases where the one step behavior assumption is violated.*
 - 2.6 *Systems that violate certain operational invariance assumptions represent cases where operational analysis fails.*

3 MEASUREMENT

- 3.1 *Because tests of operational assumptions such as flow balance and homogeneity are almost certain to reveal that the assumptions do not hold, operational analysis does not increase one's confidence in the validity of the results.*
- 3.2 *Operational laws are not useful for prediction because measured values can change randomly from one observation to the next. Distribution-free stochastic theorems are invariant and therefore useful for prediction.*
- 3.3 *Do operational results depend on the way workloads (job classes) are defined and their parameters measured?*
- 3.4 *Operational measurements over short intervals are misleading because they may differ significantly from the truth, owing to end effects.*
- 3.5 *Because it is unlikely that the initial and final states of an observation will be the same, flow balance is an unrealistic assumption.*
- 3.6 *Because it is often infeasible to test whether homogeneity assumptions are satisfied, and unlikely that they are, homogeneity assumptions are no more practical than Markovian assumptions.*

4 NATURE OF OPERATIONAL ASSUMPTIONS

- 4.1 *Only Poisson arrival processes and exponential service distributions can generate one-step behavior sequences.*
- 4.2 *Only Poisson arrivals, exponential services, or Markovian networks can satisfy the respective operational assumptions of homogeneity.*
- 4.3 *Flow balance is equivalent to stochastic steady state.*
- 4.4 *The "online = offline" condition is equivalent to the stochastic Markovian assumption.*
- 4.5 *Operational laws are equivalent to distribution free stochastic theorems.*

5 PREDICTION AND UNCERTAINTY

- 5.1 *Since stochastic analysis deals directly with random variables, does it not have an advantage in applications involving prediction?*
- 5.2 *It is said that operational analysis is based only on testable assumptions. In using an operational result to make a prediction, one is assuming that assumptions like flow balance and homogeneity will hold. Because there is no way to prove this, a basic tenet of operational analysis has been violated.*
- 5.3 *Since the length of the observation period is part of an operational analysis, the values of parameters measured in a short period may bear no relationship to the values of the same parameters in a long period. In particular, these values may be confounded by end effects and natural variations.*

5.4 *In making predictions, neither stochastic nor operational analysis has an advantage.*

5.5 *Operational results can lead to inconsistencies when used for prediction.*

6 DETERMINISM AND OPERATIONAL ANALYSIS

6.1 *Because operational performance metrics are exact for behavior sequences satisfying the assumptions, operational analysis is part of deterministic analysis.*

6.2 *Since operational results are "deterministically correct," they cannot be used for prediction.*

7 IS OPERATIONAL ANALYSIS REALLY A NEW DISCIPLINE?

7.1 *Operational analysis is incapable of producing new results.*

7.2 *What new operational results have been derived?*

7.3 *Why are not confidence intervals the stochastic counterparts of operational sensitivity analyses?*

7.4 *It is necessary to carry out a stochastic derivation in order to obtain a truly rigorous result.*

7.5 *Stochastic analysts have proved all the operational results informally. There is nothing new.*

7.6 *Operational analysis is a branch of renewal theory.*

7.7 *Operational analysis provides little information about the underlying stochastic process.*

7.8 *Once an equation has been derived in stochastic analysis, it is applied by substituting measured values for its parameters. Operationally, there is no difference between operational and stochastic analysis.*

9. Peter J. Denning, "On increasing confidence in confidence intervals," Purdue University, Department of Computer Sciences, Technical Report, CSD-TR-324 (January 1980), 20pp.

A set of N data elements x_1, \dots, x_N has r^{th} moment $E(x^r) = (x_1^r + \dots + x_N^r)/N$ and variance of the r^{th} moment $\text{Var}(x^r) = E(x^{2r}) - E^2(x^r)$. The r^{th} moment of the data elements in an arbitrary subset of k elements is used to estimate $E(x^r)$. Over all the choices of the sample, the mean error is zero and the mean square error is $(N-k)\text{Var}(x^r)/k(N-1)$. A little-known theorem by Madow shows that the frequency distribution of values of the estimator is approximately normal with mean $E(x^r)$ and variance $\text{Var}(x^r)$. All these results are proved without assuming statistical independence among sampled data elements. The conclusion is that confidence-interval calculations based on the normal distribution actually apply in more cases than is commonly believed.

10. James P. Bouhana, Jeffrey P. Buzen, and Allen I. Levy, "Data collection for analytic modeling of UNIVAC 1100 systems," Technical Report, BGS Systems, Inc. (May 1981).

A methodology is presented for extracting parameters for queueing network models from standard Univac performance reports. The parameters are described using the notation and conventions of the BEST/1 modeling package. Three types of queueing models corresponding to different system workloads (batch, demand, and TIP) are discussed. Data sources that yield the required model parameters are then described. Principal data sources are the SIP/PAR reports and LASSO; reports from IOTRACE are needed only if TIP is being used. In instances where the nature of reported data does not conform directly with required model input, algorithms are given so that suitable approximations to the required input may be derived. Many of the techniques presented here can be generalized to other cases.

11. Jeffrey A. Brumfield and Peter J. Denning, "Operational analysis of queues with general service times," Purdue University, Department of Computer Science, Technical Report CSD-TR-357 (January 1981), 33pp.

Formulas relating the mean queue length, utilization, and coefficient of variation of service time of a queue during a given observation period are derived using operational analysis. The main formula, a counterpart of the Pollaczek-Khintchin formula for $M/G/1$ queues, relies on four homogeneity assumptions: 1) the queue length at the start of a service period is independent of the length of the service period; 2) the arrival rate is not conditioned on queue length; 3) the arrival rate is not conditioned on length of service period; and 4) the total of forward residuals is the same as the total of backward residuals. A forward (backward) residual is the time between an arrival and the nearest end (start) of a service period. A long behavior sequence from an $M/G/1$ queue will satisfy these assumptions. Alternate formulas, depending on fewer assumptions, are derived. Simulation experiments compare the robustness of these estimators against the stochastic Pollaczek-Khintchin formula.