

AD-A185 763

PORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION unclassified			1b. RESTRICTIVE MARKINGS		
2a. SECURITY CLASSIFICATION AUTHORITY			3. DISTRIBUTION / AVAILABILITY OF REPORT unlimited		
2b. DECLASSIFICATION / DOWNGRADING SCHEDULE					
4. PERFORMING ORGANIZATION REPORT NUMBER(S)			5. MONITORING ORGANIZATION REPORT NUMBER(S)		
6a. NAME OF PERFORMING ORGANIZATION The Regents of the University of California		6b. OFFICE SYMBOL (if applicable)	7a. NAME OF MONITORING ORGANIZATION SPAWAR		
6c. ADDRESS (City, State, and ZIP Code) Berkeley, California 94720			7b. ADDRESS (City, State, and ZIP Code) Space and Naval Warfare Systems Command Washington, DC 20363-5100		
8a. NAME OF FUNDING / SPONSORING ORGANIZATION DARPA		8b. OFFICE SYMBOL (if applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER		
8c. ADDRESS (City, State, and ZIP Code) 1400 Wilson Blvd. Arlington, VA 22209			10. SOURCE OF FUNDING NUMBERS		
			PROGRAM ELEMENT NO.	PROJECT NO.	TASK NO.
			WORK UNIT ACCESSION NO.		
11. TITLE (Include Security Classification) An Empirical Investigation of Load Indices for Load Balancing Applications *					
12. PERSONAL AUTHOR(S) Domenico Ferrari and Songnian Zhou *					
13a. TYPE OF REPORT technical		13b. TIME COVERED FROM TO	14. DATE OF REPORT (Year, Month, Day) * July, 1987		15. PAGE COUNT * 14
16. SUPPLEMENTARY NOTATION					
17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)		
FIELD	GROUP	SUB-GROUP			
19. ABSTRACT (Continue on reverse if necessary and identify by block number) Enclosed in paper.					
<div data-bbox="264 1461 695 1640" data-label="Text"> <p>DISTRIBUTION STATEMENT A Approved for public release; Distribution Unlimited</p> </div> <div data-bbox="824 1351 1194 1630" data-label="Text"> <p>DTIC ELECTE D NOV 09 1987 S D</p> </div>					
20. DISTRIBUTION / AVAILABILITY OF ABSTRACT <input checked="" type="checkbox"/> UNCLASSIFIED/UNLIMITED <input type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS			21. ABSTRACT SECURITY CLASSIFICATION unclassified		
22a. NAME OF RESPONSIBLE INDIVIDUAL			22b. TELEPHONE (Include Area Code)		22c. OFFICE SYMBOL

Productivity Engineering in the UNIX† Environment

An Empirical Investigation of Load Indices for Load Balancing Applications

Technical Report

S. L. Graham
Principal Investigator

(415) 642-2059

"The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government."

Contract No. N00039-84-C-0089

August 7, 1984 - August 6, 1987

Arpa Order No. 4871



†UNIX is a trademark of AT&T Bell Laboratories

Accession For	
NTIS CRA&I	<input checked="checked" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

87 10 20 117

- 1 -

An Empirical Investigation of Load Indices for Load Balancing Applications†

Domenico Ferrari and Songnian Zhou

(Computer Systems Research Group)
Computer Science Division
Department of Electrical Engineering and Computer Sciences
University of California, Berkeley

The authors

Studied are

In this paper, we empirically evaluate the quality of several load indices in the context of dynamic load balancing. We have implemented a load balancer for Sun/UNIX† environments. In our experimental setup, six Sun-2 workstations were driven by job scripts, and job response times were measured while loads were being balanced and various load indices used to make job placement decisions. We study the effects on performance of the choice of load index, the averaging interval, the load information exchange period, and the characteristics of the workload. Measurements show that the performance benefits of load balancing are indeed strongly dependent upon the load index. Load indices based on resource queue lengths are found to perform better than those based on resource utilization, and the use of an exponential smoothing method yields further improvement over that of instantaneous queue lengths.

† This work was partially sponsored by the Defense Advanced Research Projects Agency (DoD), Arpa Order No. 4871, monitored by Space and Naval Warfare Systems Command under Contract No. N00039-84-C-0089, and by the National Science Foundation under grant DMC-8503575. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the Defense Research Projects Agency or of the US Government.

† UNIX is a trademark of AT&T Bell Laboratories.

1. INTRODUCTION

In a loosely-coupled distributed system, the potential for resource sharing and its possible rewards are substantial. Two frequently cited advantages of resource sharing are the larger number of accessible resources, in terms of both type and quantity, and the higher reliability that may result from the multiplicity of available resources. In order to share these resources effectively, some measure of the loads being imposed on the resources has to be made available to the clients. The information about resource loads is part of the system's state, and is among the most rapidly changing aspects of it. Since the loads are likely to be changing all the time, load information tends to become stale rapidly. To quantify the concept of load, we use a *load index*, which preferably is a non-negative variable taking on a zero value if the resource is idle, and increasing positive values as the load increases. This paper is concerned with the quality of the possible load indices for hosts in a particular but important application of load indices, that of load balancing in distributed systems.

A job arriving at a host will very likely demand services from a number of resources (e.g., CPU and disks). Hence, it is important to define not only the load on a single resource in a host, but also that of the host viewed as a collection of resources. Since the resource consumption patterns of the jobs are likely to be different, it may not be meaningful to talk about "the load" of the host. For example, the CPU may be heavily congested, while the disks are not. In this case, to an incoming CPU-bound job the host's load is very high, whereas to an incoming I/O-bound job the host's load is low because it will not experience much queuing at the disks. This observation is formalized in [Ferrari86], where a job type-dependent load index based on the resource queue lengths is proposed and experimentally evaluated.

Load information is important since it can serve as the basis of the efforts to improve the system's performance by redistributing the loads. It is frequently observed that, in a distributed system, the loads of the hosts are not evenly distributed all the time. Livny and Melman pointed out that, for a queueing system consisting of multiple homogeneous service centers with Poisson arrivals of identical rates, the probability of some hosts being idle while some others have more than one job can be very significant; hence, redistributing the workload among the resources has the potential of improving performance [Livny82].

In order to evaluate the quality of a load index for load balancing, we specify a number of criteria, or desirable properties. These criteria, in turn, are dependent on the objective of load balancing, i.e., the *performance index* that is to be optimized by balancing the loads. In this research, we are mostly concerned with interactive computing environments, where the job response time and its predictability are very important measures of system performance. Therefore, we use the mean job response time as our performance index, supplemented by the standard deviation of the response times. A good load index should:

- 1) be able to reflect our qualitative estimates of the current load on a host;
- 2) be usable to predict the load in the near future, since the response time of a job will be affected more by the future load than by the present load;

$$li = \sum_{j=1}^N s_j \times q_j$$

where N is the total number of resources for which there is queueing in the host. This index was evaluated with measurement experiments under a production time-sharing workload [Zhou87b].

The index introduced in [Ferrari86] is response time oriented, and job dependent. Instead of a unique value at a particular moment in time, the load of a host differs for different jobs because of their varying resource demands, which are assumed to be known upon job arrival. This assumption enables us to predict the response time of a job more accurately, hence to make better load balancing decision. However, while we have found some simple relationships between the arguments of a job and the job's resource demands [Zhou87c], the assumption that the demands of a job are known in advance may be too strong in many cases. In this study, we investigate versions of the same load index in which the coefficients of the resource queue lengths are *job independent*, and only reflect the relative importance of the resources (with respect to a "basket" of jobs). For example, we can use unity as the coefficients to reduce the linear combination to the sum of the resource queue lengths, that is, in queueing modeling terms, "the number of jobs (or processes) in the system."

Our extensive measurements of production time-sharing workload show that the system load is changing quite rapidly [Zhou87b]. On top of a low-frequency main component, there are a number of high-frequency load components that may be regarded as "noise" rather than useful information. Using the instantaneous resource queue lengths may give excessive importance to such noise and lead to bad job transfer decisions. We used a smoothing algorithm to compute the time-averaged queue length and compared load balancing performance using smoothed queue lengths to that of the same scheme using instantaneous queue lengths.

3. SYSTEM AND WORKLOAD

In this section, we describe the experimental environment in which the measurements were taken, and the workloads used to drive the system.

System

We implemented a dynamic load balancer for Sun/UNIX environments. The structure of the system is shown in Figure 1†. The UNIX user interface program, *cs*, is modified so that the commands typed in by the user are intercepted, and some of them are transferred to some remote host for execution when the local host is heavily loaded‡. At startup time, the C-shell reads in a configuration file that specifies a list of job types

† To distinguish our modified C shell from the standard one [Joy80], we call it *C-shell*. The *R-shell*, to be described below, shares the same software with the C-shell, but its only function is to receive remote jobs and execute them.

‡ Our system is based on a modified C shell implemented at Berkeley by Harry Rubin and Venkat Rangan for the Berkeley UNIX 4.3 BSD system running on VAX machines [Joy83, McKusick85].

home C-shell, and are terminated when the home C-shell exits. This scheme has the potential problem of R-shell proliferation. However, the code segments of all C-shells and R-shells on each host are shared, so that, when an R-shell is not active, almost no resources are consumed. Since files are retrieved from file servers, as the workstations are diskless, only the command line needs to be shipped, and the cost of file access is essentially the same from all hosts.

Load balancing algorithms have a strong influence on performance. We implemented and studied a number of algorithms using different methods for load information exchange and job placement [Zhou87a]. For this study of load indices, however, we just selected one of the best realizable algorithms, that is, the one we called GLOBAL. For every time period P , the LIM on each host extracts load information from the local kernel to compute the local host's load index. If the new value of the load index is significantly different from the previous one, the new value is sent to the *master* LIM, which collects load information from every host and broadcasts the entire load vector in each period P . When a job whose name is on the eligibility list is submitted to a host, the local LIM is contacted for job placement. If the local load is high, the host perceived by the local LIM to have the least load is selected, and the job is sent there.

The implementation described above provides a transparent, low-cost, and general-purpose load balancer whose installation requires no changes to the kernel† or to the application programs. Since the emphasis of this paper is on the measurement experiments we performed on the system, we will not describe the design and implementation issues in more detail. The interested reader is referred to [Zhou87a].

Workload

Workload characterization and selection are crucial to a measurement study. Although artificial workloads considerably increase the repeatability of experiments, they ought to represent natural workloads reasonably well, so as to strengthen our confidence in the results. We traced a production VAX-11/780 machine running under the Berkeley UNIX 4.3BSD system [Joy83, McKusick85] for an extended period of several months, and analyzed the types and frequencies of the commands executed by the system. On the basis of such an analysis, we selected 30 frequently executed commands, listed in Table 1, and used them to construct job scripts, i.e., sequences of commands.

To obtain various levels, or intensities, of a host's load, we ran a variable number of jobs in the background. The artificial workloads were not intended to represent the typical workloads of personal workstations, but rather those of small (i.e., not very powerful) time-sharing systems. Workstations were used because of their being available in our distributed systems laboratory. We simulated user think times by the "*sleep*" command. The scripts are classified into three levels: light (L), moderate (M), and heavy (H), with a number of distinct scripts constructed for each level, so that hosts subjected to the same level of workload always use different scripts. The ranges of CPU utilization and mean load index values of the three levels of scripts are shown in Table 2. Each script runs for about 30 minutes on a Sun-2 workstation. Job and system performance statistics, such as

† For our experiments, to obtain accurate values of resource queue lengths and to perform the smoothing operations efficiently, some code had to be added to the kernel. No functional changes were made, however.

interval (CI) of the values of the performance indices over these replications.

4. DESIGN AND RESULTS OF THE EXPERIMENTS

Experimental Factors

Four factors were identified to be of interest in the study of load indices:

- 1) **Load index.** We used as load indices the following quantities: the instantaneous CPU queue length; exponentially averaged CPU queue length; the sum of averaged CPU, file and paging/swapping I/O, and memory queue lengths†; and the average CPU utilization over a recent period. Inside the kernel, we kept variables for the queue length of each resource type. The length of each queue was sampled every 10 ms by the clock interrupt routine, and used to compute the one-second average queue length, q_i . Exponential smoothing was used to compute the average queue length over the last T seconds:

$$Q_i = Q_{i-1}(1-e^{-T}) + q_i e^{-T}, \quad i \geq 1$$

$$Q_0 = 0$$

- 2) **Averaging interval T .** For exponentially smoothed values of a resource queue length, and for the average CPU utilization, the interval T over which the average is computed conceivably affects the quality of the index, and hence the system's performance.
- 3) **Workload.** There may be interactions between the load index chosen and the workload the system is subjected to. Using the three suites of host workloads described in the previous section, we were able to construct several combinations of system workload for the six workstations in our system. The canonical workload consisted of two heavy, two moderate, and two light scripts (2H, 2M, 2L). We also studied the indices under a more balanced workload, with all six workstations driven by moderate scripts (6M).
- 4) **Exchange interval P .** The GLOBAL algorithm employs periodic updates of load information. If P is too short, the overhead may be too high, but, if P is too long, then job placements are based on stale information, and performance may deteriorate, and system instability may result.

Measurement Results

We shall first study the indices and the averaging interval T by fixing the workload at its canonical level, and the exchange interval at 10 seconds. We will then use the more balanced workload 6M to examine the interactions between load indices and workload. Finally, we will study the effect of load exchange interval P on performance.

† For simplicity, we treated the disk queues as a single aggregate queue for I/O operations. For the memory queue, we identified a number of places inside the kernel where processes queue up for various types of memory resources (e.g., buffer space, page table), and treated all these as a single memory queue.

Comparing the queue-length-based indices with each other, we notice that the exponentially smoothed indices can perform best, but, if the averaging period T is too long (e.g., ≥ 20 s), performance may even become worse. Earlier in this paper, we have pointed out that, by averaging the queue lengths, the adverse effect of the high-frequency "noise" in the load can be reduced. This is reflected by improved performance. However, since the system load is changing all the time, averaging over too long a period will emphasize too much the past loads, which have little correlation with the future ones. The optimum averaging interval is clearly dependent upon the dynamics of the workload: the faster the load changes, the shorter the interval should be. In a measurement study of production workloads on a VAX-11/780 running Berkeley UNIX 4.2BSD [Zhou87b], we found that the average net change in CPU queue length in 30 seconds was 2.31, when the average CPU queue length itself was 4.12. This suggests that T should be substantially shorter than 30 seconds.

The performance difference between the cases in which indices based on CPU queue alone are used, and those in which indices consider I/O and memory contention also, is not significant, suggesting that the CPU is the predominant resource in our hosts. We found that the I/O and memory queue lengths were generally much shorter than that of CPU; that is, the former are much less contended for. It should be pointed out that our systems support general computing in a research environment; with other types of workload, e.g., database-oriented one, the contention profile of the various resource types may be substantially different. However, to achieve near-optimal performance, we do not have to consider all the resources in the system, but rather only those with significant contention. We also studied more general forms of linear combinations of queue lengths by using coefficients other than unity, but no significant changes in performance were observed. This, again, is probably due to the dominating influence of the CPU queue.

The load average shown in Table 3 is an index provided by a UNIX command; it is the exponentially smoothed number of processes ready to run, or running, or waiting for some high-priority event (e.g., disk I/O completion). A number of load balancers constructed in the past in the UNIX environment have used the load average as their load index (e.g., [Bershad85]). This research shows that significant further improvement can be obtained by using indices that more accurately reflect the current queueing at the resources.

The performances produced by the indices under the more balanced workload 6M is shown in Table 4. Since the workload is now more balanced and moderate, the amount of improvement in response time is not as much as that under the canonical workload; however, the relative rankings of the indices are quite similar. This suggests that the above analyses of the qualities of the indices and the appropriate values for T remain valid under a more balanced, moderate workload. It is worth noting that, in this case, due to the smaller improvement, using a poor load index (e.g., load average or 60 s CPU utilization) may yield little or no performance improvement.

Finally, we study the influence of the load exchange period P . Figure 2 shows the mean job response time as a function of P , and with the other three factors fixed. The brackets around the data points show their 90% confidence intervals. When the exchange period P is very short, the load information used in job placements is generally up to date, but this positive influence is outweighed by high message overhead. Conversely, if P is too long, the information may get stale, the quality of job placements deteriorates, and

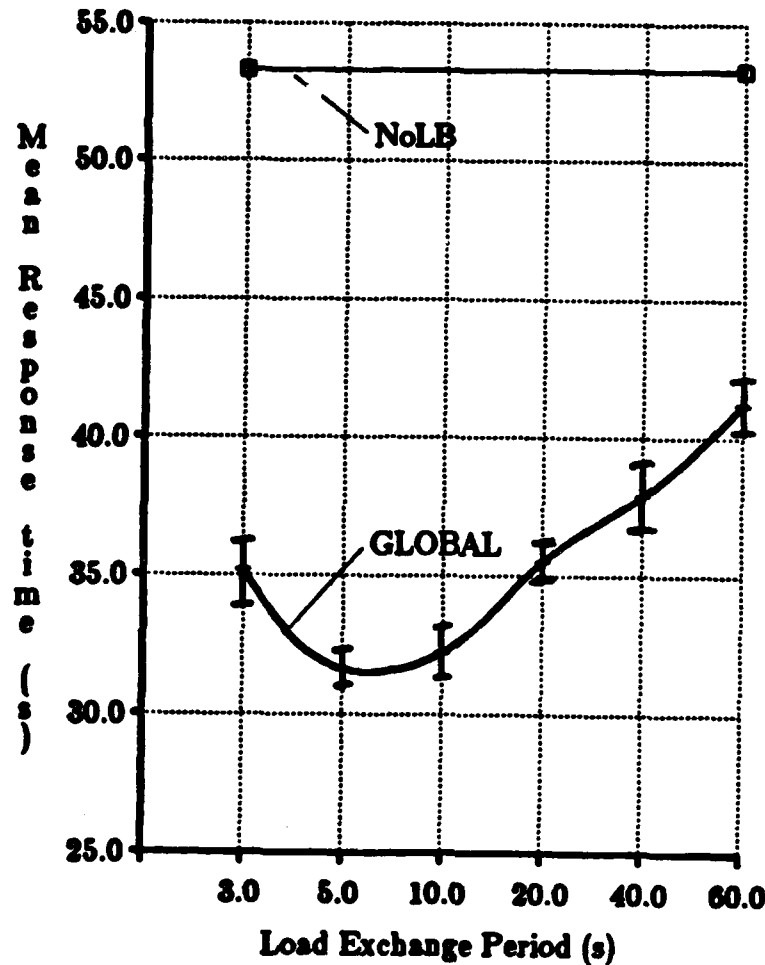


Figure 2. Mean process response time under various load exchange periods P
(Canonical workload, load index 4 s CPU+I/O+Mem ql).

criteria reasonably well: the queue length is an accurate measure of a resource's load, and smoothing over a short interval into the past gives predictive capabilities to the value of the index, as well as stability against the noise in the load waveform. Queue-length-based load indices also appear to be more adaptable to a heterogeneous environment, but more studies are needed to substantiate this conjecture.

Our results support indices compatible with the one proposed in [Ferrari86], as they can be seen as degenerate forms of that index. However, the comparisons performed in this study are far from being complete. We decided to use the same load balancing algorithm for all the indices, so that the qualities of the load indices may be directly comparable. On the other hand, the algorithm limited the varieties of load indices that could be studied. We demonstrated, using a particular set of workloads and in a particular computing environment, that linear combinations of resource queue lengths may be good load indices. No proof, however, is offered that they are the best.

Systems, pp. 54-69, May 1986.

[Livny82]

M. Livny and M. Melman, "Load Balancing in Homogeneous Broadcast Distributed Systems," Proc. ACM Computer Network Performance Symposium, pp. 47-55, April 1982.

[McKusick85]

K. McKusick, M. Karels, and S. Leffler, "Performance Improvements and Functional Enhancements in 4.3 BSD," Proc. Summer USENIX Conference, June 1985, Portland, OR, pp. 519-531.

[Wang85]

Y. Wang and R. Morris, "Load Balancing in Distributed Systems," IEEE Trans. Comp. Vol.C-34, No.3, pp. 204-217, March 1985.

[Zhou86]

S. Zhou, "A Trace-Driven Simulation Study of Dynamic Load Balancing," Tech. Rept No. UCB/CSD 87/305, September 1986, also submitted for publication.

[Zhou87a]

S. Zhou and D. Ferrari, "An Experimental Study of Load Balancing Performance," Tech. Rept No. UCB/CSD 87/336 January 1987, also submitted for publication.

[Zhou87b]

S. Zhou, "An Experimental Assessment of Resource Queue Length as Load Indices," Proc. Winter USENIX Conference, Washington, D.C., pp. 73-82, January 21-24, 1987.

[Zhou87c]

S. Zhou, "Predicting Job Resource Demands: a Case Study in Berkeley UNIX," in preparation.