



US006223205B1

(12) **United States Patent**
Harchol-Balter et al.

(10) **Patent No.:** **US 6,223,205 B1**
(45) **Date of Patent:** **Apr. 24, 2001**

(54) **METHOD AND APPARATUS FOR ASSIGNING TASKS IN A DISTRIBUTED SERVER SYSTEM**

(76) Inventors: **Mor Harchol-Balter**, 47 Market St. #2, Cambridge, MA (US) 02139; **Mark E. Crovella**, 14 Collier Rd., Scituate, MA (US) 02066

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.

(21) Appl. No.: **09/023,893**

(22) Filed: **Feb. 13, 1998**

Related U.S. Application Data

(60) Provisional application No. 60/063,484, filed on Oct. 20, 1997.

(51) **Int. Cl.**⁷ **G06F 9/00**

(52) **U.S. Cl.** **709/105; 709/102**

(58) **Field of Search** 709/312, 203, 709/102, 105, 100, 101, 103, 104

(56) **References Cited**

U.S. PATENT DOCUMENTS

5,088,032	2/1992	Bosack	395/200
5,539,883	7/1996	Allon et al.	395/200.11
5,548,724	* 8/1996	Akizawa et al.	709/203
5,617,570	* 4/1997	Russell et al.	709/312
5,644,720	* 7/1997	Boll et al.	709/227

OTHER PUBLICATIONS

Almeida, J. et al., "Measuring the behavior of a World-Wide Web server," Seventh Conference on High Performance Networking (HPN), pp. 57-72, White Plains, NY, Apr., 1997. IFIP.

Almeida, V. et al., "Performance analysis of a WWW server," Proceedings of CMG '95, 1995.

Anderson et al., The magicrouter: An application of fast packet interposing. Technical report, UC Berkeley, 1996.

Arlitt et al., Web server workload characterization: The search for invariants. Proceedings of the 1996 SIGMETRICS Conference on Measurement and Modeling of Computer Systems, pp. 126-137, 1996.

Boyle, P., Web site traffic cops: Load balancers can provide the busiest web sites with non-stop performance. *PC Magazine*. Available online at <http://8.zdnet.com/pcmag/>, Feb. 18, 1997.

(List continued on next page.)

Primary Examiner—Majid A. Banankhah

(57) **ABSTRACT**

A distributed server system is disclosed which includes a load balancer and a plurality of host processors. The load balancer receives requests for service and distributes task assignments among the plurality of processors based upon the amount of work associated with the respective requests for service. More specifically, each host processor services requests for service within a predefined task size interval and the load balancer assigns to each host processor only those requests for service which involve task sizes within the particular task size interval associated with the respective processor. In the foregoing manner, the variability of the task sizes assigned to any given host processor is minimized and performance of the distributed server system is improved. In one embodiment of the invention, the thresholds defining the task size intervals served by respective host processors are selected using the task size distribution so as to attempt to allocate substantially the same amount of work to each of the host processors. In another embodiment of the invention, the thresholds defining the task size intervals served by the respective host processors are selected so as to intentionally vary the amount of work performed by the respective processors. This embodiment of the invention is intended to service requests for service in the nature of heavy tailed distributions. Task size interval assignments are skewed such that requests for service corresponding to smaller tasks are serviced by a host processor which is more lightly loaded than other processors. In this manner, mean slowdown metrics are improved.

39 Claims, 6 Drawing Sheets

