

A Queueing Analysis for Job Assignment on Two-Type Heterogeneous Supercomputer System

C.H. Wang

Department of Industrial and Information Management
National Cheng Kung University
Taiwan

Y.T. Chen

Supercomputing Research Center, Computer and Network
Center
National Cheng Kung University
Taiwan

C.C. Hwang

Department of Engineering Science
National Cheng Kung University
Taiwan

Abstract—This paper introduces a queueing analysis of a two-tier service supercomputer system, where one type of processors offers service with a finite waiting buffer and the other offers unlimited waiting buffer. From the managerial viewpoint of system, a two-tier queueing model is developed to investigate the dynamic performance of the supercomputer system under finite buffer control. The queueing system is formulated as a state dependent quasi-birth-and-death (QBD) process with two-dimensional state space. System performance measures can be obtained through the matrix geometric solution for such a QBD process.

Keywords-finite buffer queue; matrix-geometric method; heterogeneous processors

I. INTRODUCTION

Supercomputers help researchers to discover real-world phenomenon through conducting large-scale scientific simulations except physical experiments. The implementations of heterogeneous system in a computer center has been lasting progress for past decades. The collection of multiprocessors is fed by a single common stream of batch jobs, where each job is dispatched to exactly one of the multiprocessor machines for processing. Examples of such distributed server systems include the Xolas distributed server at the MIT Lab for Computer Science, the Cray J90 distributed server at NASA Ames Research Lab, the Cray J90 distributed server at the Pittsburgh Supercomputing Center, and the Cray C90 distributed server at NASA Ames Research Lab (see Schroeder & Harchol-Balter [1] and reference therein).

Job scheduling on heterogeneous computer clusters is complicated, and how to utilize all heterogeneous server systems becomes an important research and managerial issue. In supercomputer centers today, designing a distributed server system (a collection of multiprocessors) often boils down to choosing the best task (computing jobs) assignment policy for the given supercomputer model and managerial requirements. As mentioned in Schroeder & Harchol-Balter [1], finding a good rule for assigning jobs to host machines remains an open question at many supercomputing sites. Over the last decades, many job assignment rules across a variety of distributed

computational resources have been studied, e.g., Schaar & Efe [2], Chlamtac et al. [3], Feitelson et al. [4], von Laszewski [5], Piro et al. [6], Tang et al. [7], etc.

Various managerial goals for job assignment are usually conflicting, such as queueing efficiency versus system utilization. Tang et al. [7] studied scheduling policies for balancing workload. Downey [8] investigated the decision making on job assignment to a space-sharing parallel computer based on observed workloads. Under First-Come-First-Served scheduling policy, von Laszewski [5] demonstrated the advantage of well utilizing compute systems while submitting jobs across a variety of cluster. Bucur & Epema [9] conducted experimental analysis for the computational power of the existing systems when scheduling rigid jobs on a multicluster systems. Piro et al. [6] improved job scheduling strategies based on historical data of computing workloads to balance job assignment efficiently. But seldom works were evaluated analytically or systematically from the managerial viewpoint of queueing system under supercomputing workloads.

The main contribution of this paper is to provide an analytic mechanism for managing heterogeneous systems in a supercomputer center. We study a real supercomputing system Cray XE6m with two type of processors, which is currently running at the Supercomputing Research Center (SRC), National Cheng Kung University (NCKU), Taiwan. To manage computing resources properly, a queueing analysis of two-tier service system with finite buffer control is conducted in this paper. One type of processors offers guaranteed delay time with a finite buffer space, and the other type offers service with unlimited waiting space. We develop a two-dimensional state dependent quasi-birth-and-death (QBD) queueing model for evaluating the impact of the finite buffer control on the system performance measures, e.g., average queue length and average waiting time.

II. ARCHITECTURE OF CRAY XE6M

Since Year 2012, a supercomputer Cray XE6m in the NCKU-SRC has been offering computing services for conducting academic research. In Year 2011, the XE6m model

Under the stability condition, the stationary probability vector is defined as

$$\boldsymbol{\pi}_n = [\pi_{n,0}, \pi_{n,1}, \dots, \pi_{n,K}],$$

where indices $n = 0, 1, \dots$, denote steady states of Type-1 queue. When $n \geq K$, the matrix geometric solution for such a QBD process can be obtained by

$$\boldsymbol{\pi}_{n+1} = \boldsymbol{\pi}_n \mathbf{R}, \quad (1)$$

where \mathbf{R} is the rate matrix. Like any regular QBD process, the rate matrix \mathbf{R} should satisfy $\mathbf{R}2\mathbf{A} + \mathbf{R}\mathbf{B} + \mathbf{C} = 0$ and can be solved by using one of many known algorithms. Interested readers may refer to Neuts [12] and references therein. For $0 \leq n \leq K$, the probability vector $\boldsymbol{\pi}_n$ can be obtained by solving a set of equations. From $\boldsymbol{\pi}\mathbf{Q} = 0$ and (1), the steady-state vectors $\boldsymbol{\pi}_0, \boldsymbol{\pi}_1, \dots$ and $\boldsymbol{\pi}_K$ can be solved from the boundary conditions (2)-(6) and the normalization condition (7) as follows:

$$\boldsymbol{\pi}_0 \mathbf{B}_{0,0} + \boldsymbol{\pi}_1 \mathbf{A} = 0, \quad (2)$$

$$\boldsymbol{\pi}_0 \mathbf{C}_{0,1} + \boldsymbol{\pi}_1 \mathbf{B}_{1,1} + \boldsymbol{\pi}_2 \mathbf{A} = 0, \quad (3)$$

$$\boldsymbol{\pi}_1 \mathbf{C}_{1,2} + \boldsymbol{\pi}_2 \mathbf{B}_{2,2} + \boldsymbol{\pi}_3 \mathbf{A} = 0, \quad (4)$$

⋮

$$\boldsymbol{\pi}_{K-2} \mathbf{C}_{K-2,K-1} + \boldsymbol{\pi}_{K-1} \mathbf{B}_{K-1,K-1} + \boldsymbol{\pi}_K \mathbf{A} \quad (5)$$

$$\boldsymbol{\pi}_{K-1} \mathbf{C}_{K-1,K} + \boldsymbol{\pi}_K (\mathbf{B}_{K,K} + \mathbf{R}\mathbf{A}) = 0, \quad (6)$$

$$\boldsymbol{\pi}_0 \mathbf{1} + \boldsymbol{\pi}_1 \mathbf{1} + \dots + \boldsymbol{\pi}_K (\mathbf{I} - \mathbf{R})^{-1} \mathbf{1} = 1, \quad (7)$$

After determining the stationary distribution, we can obtain the major system performance measures, including average queue length and average waiting time. In the case of large-scale supercomputing system with huge buffer size K , it would result in a large number of boundary states and a large number of phases of the QBD process. It would greatly increase the computational complexity and may cause the ill-conditioned matrices of the traditional iterative algorithm for the rate matrix. To overcome this challenge, an efficient algorithm proposed by Luh et al. [11] could be applied directly to the case of medium to large scale buffer size K for determining the stationary distribution in large-scale supercomputing systems.

IV. NUMERICAL EXPERIMENTS

We consider two scenarios in this numerical experiments, that is, heavy traffic load and light traffic load. We demonstrate the effect of finite buffer control on the system performances, i.e., average queue length and average waiting time. For the heavy traffic load, we have $\lambda = 1$, $\mu_1 = 0.6$ and $\mu_2 = 0.8$; on the other hand, we take $\lambda = 0.8$ for the light traffic load. The numerical experiments are conducted through computing language MATLAB on the PC platform with Intel® Core™ i7-3770 CPU @ 3.40 GHz and 32 GB RAM.

As the finite buffer size K varies from 2 to 15, Figure 2 shows the average queue length of jobs for Type-1 processors and Type-2 processors, individually. The effect of increasing the buffer size on deducing average waiting time is illustrated in Figure 3. As the finite buffer size is large enough, it can be found that the benefit of increasing the finite buffer on improving system performances would decrease.

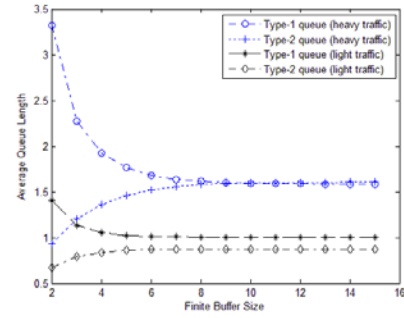


FIGURE II. AVERAGE QUEUE LENGTH VERSUS THE BUFFER SIZE.

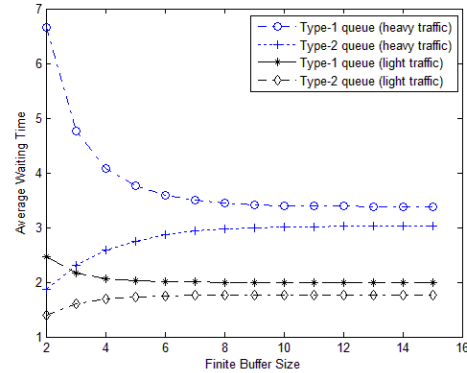


FIGURE III. AVERAGE WAITING TIME VERSUS THE BUFFER SIZE.

V. CONCLUSIONS

We present a queueing model based on a real supercomputer system and observed workloads at NCKU-SRC. The state dependent QBD process derived in this paper allows us to investigate the dynamic performance of such a two-tier service supercomputer system. We draw a conclusion that proper buffer control can balance queuing efficiency and the cost of enlarging the finite buffer size. According to the presented queueing analysis, managers could assign arriving jobs to appropriate type of processors in order to achieve the better system performance. The queueing analysis conducted in this study could help in developing a management scheme for job assignment on heterogeneous supercomputer systems with finite buffer control.

ACKNOWLEDGEMENTS

The authors would like to thank Prof. Yuefan Deng and Dr. Sing-Wu Liou for useful discussions on this work. This research was partially supported by the National Science Council, Taiwan, R.O.C., under grant number NSC 102-2911-I-006-301.

REFERENCES

- [1] Schroeder, B. & Harchol-Balter, M., Evaluation of task assignment policies for supercomputing servers: The case for load unbalancing and fairness, *Cluster Computing*, vol. 7, pp. 151-161, 2004.
- [2] Schaar, M.A. & Efe, K., Effective queueing strategies for co-scheduling in a pool of processors, *Computer Communications*, vol. 19, pp. 743-753, 1996.

- [3] Chlamtac, I. Kienzle, M.G. & Szabo, C., Characterizing the behaviour of high-speed interconnection systems with distributed control, *Telecommunication Systems*, vol. 6, pp. 91-115, 1996.
- [4] Feitelson, D.G., Rudolph, L., Schwiegelshohn, U., Sevcik, K.C. & Wong, P., Theory and practice in parallel job scheduling, in *Proceedings of the 3rd Workshop on Job Scheduling Strategies for Parallel Processing*, pp. 1-34, 1997.
- [5] von Laszewski, G., A loosely coupled metacomputer: co-operating job submissions across multiple supercomputing sites, *Concurrency Practice and Experience*, vol. 11, pp. 933-948, 1999.
- [6] Piro, R.M., Guarise, A., Patania, G. & Werbrouck, A., Using historical accounting information to predict the resource usage of grid jobs, *Future Generation Computer Systems*, vol. 25, pp. 499-510, 2009.
- [7] Tang, W., Ren, D., Lan, Z. & Desai, N., Toward balanced and sustainable job scheduling for production supercomputers, *Parallel Computing*, vol. 39, pp. 753-768, 2013.
- [8] Downey, A.B., Using queue time predictions for processor allocation, in *Proceedings of the 3rd Workshop on Job Scheduling Strategies for Parallel Processing*, pp. 35-57, 1997.
- [9] Bucur, A.I. & Epema, D.H., The influence of the structure and sizes of jobs on the performance of co-allocation, *Lecture Notes in Computer Science*, vol. 1911, pp. 154-173, 2000.
- [10] Kerbyson, D.J., Barker, K.J., Vishnu, A. & Hoisie, A., A performance comparison of current HPC systems: Blue Gene/Q, Cray XE6 and InfiniBand systems, *Future Generation Computer Systems*, vol. 30, pp. 291-304, 2014.
- [11] Luh, H., Zhang, Z. & Wang, C.H., A computing approach to two competing services with a finite buffer effect, in *Proceedings of the 8th International Conference on Queueing Theory and Network Applications*, pp. 15-21, 2013.
- [12] Neuts, M.F., *Matrix-Geometric Solutions in Stochastic Models*. The John Hopkins University Press, 1981.