

Chapter 5

Communication Analysis and Performance Prediction of Parallel Applications on Large-Scale Machines

Yan Li

Intel Labs China, China

Jidong Zhai

Tsinghua University, China

Keqin Li

State University of New York, USA

ABSTRACT

With the development of high performance computers, communication performance is a key factor affecting the performance of HPC applications. Communication patterns can be obtained by analyzing communication traces. However, existing approaches to generating communication traces need to execute the entire parallel applications on full-scale systems that are time-consuming and expensive. Furthermore, for designers of large-scale parallel computers, it is greatly desired that performance of a parallel application can be predicted at the design phase. Despite previous efforts, it remains an open problem to estimate sequential computation time in each process accurately and efficiently for large-scale parallel applications on non-existing target machines. In this chapter, we will introduce a novel technique for performing fast communication trace collection for large-scale parallel applications and an automatic performance prediction framework with a trace-driven network simulator.

INTRODUCTION

Different applications in the high performance computing (HPC) field exhibit different communication patterns, which can be characterized by three key attributes: volume, spatial and temporal (Chodnekar, et al., 1997; Kim & Lilja, 1998). Proper understanding of communication patterns of parallel applica-

DOI: 10.4018/978-1-5225-0287-6.ch005

tions is important to optimize the communication performance of these applications (Chen et al., 2006; Preissl, et al., 2008a). For example, with the knowledge of spatial and volume communication attributes, MPIPP (Chen, et al., 2006) optimizes the performance of Message Passing Interface (MPI) programs on non-uniform communication platforms by tuning the scheme of process placement. Besides, such knowledge can also help design better communication subsystems. For instance, for circuit-switched networks used in parallel computing, communication patterns are used to pre-establish connections and eliminate the runtime overhead of path establishment. Furthermore, a recent work shows spatial and volume communication attributes can be employed by replay-based MPI debuggers to reduce replay overhead significantly (Xue, et al., 2009).

Previous work on communication patterns of parallel applications mainly relies on traditional trace collection methods (Kim & Lilja, 1998; Preissl et al., 2008b; Vetter & Mueller, 2002). A series of trace collection and analysis tools have been developed, such as ITC/ITA (Intel, 2008; Kerbyson et al., 2001), KOJAK (Mohr & Wolf, 2003), TAU (Shende & Malony, 2006), DiP (Labarta et al., 1996) and VAMPIR (Nagel et al., 1996). These tools need to instrument original programs at the invocation points of communication routines. The instrumented programs are executed on full-scale parallel systems and communication traces are collected during the execution. The collected communication trace files record type, size, source and destination etc. for each message. The communication patterns of parallel applications can be easily generated from the communication traces. However, traditional communication trace collection methods have two main limitations: huge resource requirement and long trace collection time. For example, ASCI SAGE routinely runs on 2000-4000 processors (Kerbyson, et al., 2001) and FT program in the NPB consumes more than 600 GB memory for Class E input (Bailey, et al., 1995). Therefore, it is impossible to use traditional trace collection methods to collect communication patterns of large-scale parallel applications without full-scale systems. Moreover, it takes several months to complete even on a system with thousands of CPUs. It is prohibitive long for trace collection and prevents many interesting explorations of using communication traces, such as input sensitivity analysis of communication patterns. Additionally, MPIP (Vetter, et al., 2001) is a lightweight profiling library for MPI applications and only collects statistical information of MPI functions. However, all these traditional trace collection methods require the execution of the entire instrumented programs, which restricts their wide usage for analyzing large-scale applications. Our method adopts the similar technique to capture the communication patterns at runtime as the traditional trace collection methods.

We have two observations on existing communication trace collection and analysis approaches:

1. Many important applications of communication pattern analysis, such as the process placement optimization (Chen, et al., 2006) and subgroup replay (Zhang, et al., 2009), do not require temporal attributes.
2. Most computation and message contents in message-passing parallel applications are not relevant to their spatial and volume communication attributes.

Motivated by the above observations, we describe a novel technique in this chapter, called FACT, which can perform fast communication trace collection for large-scale parallel applications on small-scale systems. Our idea is to reduce the original program to obtain a program slice through static analy-

24 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/communication-analysis-and-performance-prediction-of-parallel-applications-on-large-scale-machines/159041

Related Content

Sustainable Consumerism via Context-Aware Shopping

Johannes Klinglmayr, Bernhard Bergmair, Maria Anneliese Klaffenböck, LeanderB. Hörmann and Evangelos Pournaras (2017). *International Journal of Distributed Systems and Technologies* (pp. 54-72). www.irma-international.org/article/sustainable-consumerism-via-context-aware-shopping/188859

A Fault Tolerant Decentralized Scheduling in Large Scale Distributed Systems

Florin Pop (2010). *Handbook of Research on P2P and Grid Systems for Service-Oriented Computing: Models, Methodologies and Applications* (pp. 566-588). www.irma-international.org/chapter/fault-tolerant-decentralized-scheduling-large/40818

Heart Rate Estimation in Sports Based on Multi-Sensor Data for Sports Intensity Prediction

Feng Zhang (2022). *International Journal of Distributed Systems and Technologies* (pp. 1-12). www.irma-international.org/article/heart-rate-estimation-in-sports-based-on-multi-sensor-data-for-sports-intensity-prediction/307990

Securing Real-Time Interactive Applications in Federated Clouds

Michael Boniface, Bassem Nasser, Mike Surridge and Eduardo Oliveros (2012). *Grid and Cloud Computing: Concepts, Methodologies, Tools and Applications* (pp. 1822-1835). www.irma-international.org/chapter/securing-real-time-interactive-applications/64569

Programmability and Scalability on Multi-Core Architectures

Jaeyoung Yi, Yang J. Jang, Doohwan Oh and Won W. Ro (2010). *Handbook of Research on Scalable Computing Technologies* (pp. 276-294). www.irma-international.org/chapter/programmability-scalability-multi-core-architectures/36412