# Lehrstuhl für Rechnertechnik und Rechnerorganisation der Technischen Universität München

# Scalable Automated Online Performance Analysis of Applications using Performance Properties

## Karl Fürlinger

Vol. 33

# Scalable Automated Online
# Performance Analysis of Applications
# using Performance Properties

## Karl Fürlinger

# Abstract

Analyzing the performance of applications is an important step in the development process of scientific software for high-performance computing systems. Programmers need to ensure that their code performs well and does make good use of the available resources on such systems.

Unfortunately, the task of finding inefficiencies and identifying their reasons can be a time-consuming and difficult challenge – even for sequential programs. For large-scale parallel systems and programs the situation is complicated by several factors. Firstly, parallel programming is inherently more complicated than sequential programming and additional middleware layers may be present that affect the achieved performance. Secondly, the constantly increasing size and complexity of parallel systems and the problems solved on them leads to challenges concerning the scalable collection, visualization, and analysis of performance data.

This work presents an approach to performance analysis and a software tool called *Periscope* that tries to avoid some of the problems outlined above in order to effectively detect inefficiencies on large-scale parallel computing systems. The approach is based on the notion of performance properties that formally capture situations of inefficient execution. By formalizing what constitutes inefficient behavior, the detection process can be automated. The component which performs the automated performance analysis in Periscope is called an *agent* and several agents distributed over the parallel machine cooperate in the overall analysis process for an application.

The decomposition of the tool into a number of agents keeps the analysis process scalable. As an application's size increases in number of processors or SMP nodes used, the tool's size can also increase in number of analysis agents employed. The distribution also solves problems related to the management of large amounts of performance data. The agents analyze data close to the origin in the application and thus problems with a centralized collection, storage, and processing are avoided. Finally, the Periscope approach allows for online analysis. That is, the performance analysis process can be initiated at any time while the program is still running, as opposed to a post-mortem approach where performance data is analyzed after program termination.

# Acknowledgments

Many people have contributed to the success of this work. First and foremost, many thanks are directed to Prof. Michael Gerndt and Prof. Arndt Bode for providing such a supportive environment for conducting research. It is my belief that there are few places that can match the freedom and support for research that I had the opportunity to receive.

As any research, this work builds on the previous ideas of many people and I would like express my gratitude for their excellent work. Most notably, the work of Bernd Mohr and Felix Wolf was very helpful and enlightening in many situations. Additionally, numerous helpful people – too many to be named individually – have offered their advice when it was needed and supported me in many different ways.

Lastly, I owe the deepest gratitude to my parents for their continuous, unwavering support over so many years.

*Karl Fürlinger*
*Munich, Germany*
*May 2006*

# Contents

CONTENTS

# List of Figures