Lehrstuhl für Rechnertechnik und
Rechnerorganisation
der Technischen Universität München

# Data Locality Optimization of Shared Memory Programs on NUMA Architectures Using an Integrated Tool Environment

## *Jie Tao*

Vollständiger Abdruck der von der Fakultät für Informatik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktors der Naturwissenschaften (Dr. rer. nat.)

genehmigten Dissertation.

Vorsitzender:         Univ.-Prof. Dr. H. M. Gerndt
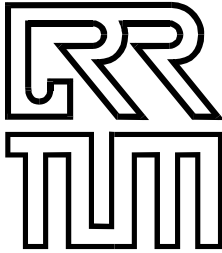
Prüfer der Dissertation:
                      1. Univ.-Prof. Dr. A. Bode

                      2. Univ.-Prof. Dr. E. Jessen

Die Dissertation wurde am 20.06.2002 bei der Technischen Universität München eingereicht und durch die Fakultät für Informatik am 29.08.2002 angenommen.

Vol. 26

# Data Locality Optimization of Shared Memory Programs on NUMA Architectures Using an Integrated Tool Environment

## Jie Tao

Printed in Germany.

# Abstract

Due to their excellent price-performance ratio, clusters built from commodity nodes have become broadly adopted and increasingly popular as platforms for parallel processing. Among them, the clusters of standard PCs interconnected with high-speed system area networks (SANs) are especially attractive and have been widely established. At the same time, the developments in interconnection technologies also formed the basis for the rise of Non-Uniform Memory Access (NUMA) architectures, i.e. systems with physically distributed memories, but with a global address space allowing an efficient but non-uniform access to any memory location in the system. These kinds of systems, especially when offered as non cache coherent NUMA for loosely coupled commodity architectures, can easily be implemented in a straightforward manner without major hardware efforts. They form a favorable architectural tradeoff by combining the scalability and cost–effectiveness of standard clusters with a shared memory support close to symmetric multiprocessors.

The non-uniform memory access characteristic, however, introduces a distinction between local and remote memory causing different memory access latencies. In systems with such characteristics, a remote memory access can take up to an order of magnitude longer than a local one. This leads to the fact that many shared memory applications initially do not achieve a good parallel speedup when running on NUMA-like architectures due to excessive remote memory accesses.

This thesis targets such inefficiency problems of NUMA-based shared memory programs. For this purpose, a comprehensive and integrated tool environment has been built which aims at improving the data locality of running applications, by combining single frameworks enabling both program tuning and runtime manipulation. This environment comprises a low-level data acquisition system, a distributed tool middleware, and a set of performance tools. Based on the hardware monitoring facility, which is capable of observing all memory transactions performed on the interconnection fabric, the data acquisition system provides information about an application's memory access behavior as well as information about e.g. synchronization primitives and address mapping necessary for data placement. This information is then aggregated across the distributed system and made accessible to the tools through an established on-line monitoring interface specification serving as middleware. Tools further process the acquired performance data and use it to steer the execution of programs with a result of an optimization for the runtime data layout.

Currently, two such tools have been implemented: a Data Layout Visualizer (DLV) and an Adaptive Runtime System (ARS). DLV is used to present the monitoring data in an easy-

to-use fashion allowing a programmer to understand the dynamic behavior of an application and to detect and correct communication hot spots. ARS is used to transparently modify the data layout during the execution of a program without involving the user. Both tools are capable of introducing a better data locality and hence a better performance.

In addition, the tool environment includes an event-driven multiprocessor simulator called SIMT. SIMT was originally designed to supply the monitoring information in case no hardware monitor is available. It has been developed, however, to be fully flexible with respect to the target architectures and can thereby be used as a general tool for system design and performance evaluation of PC clusters with NUMA organization.

The monitoring and tool environment has been evaluated using a set of shared memory applications selected from standard benchmarks and kernels. Most of them show a high improvement in terms of absolute execution time and speedup, proving the feasibility and effectiveness of this monitoring approach. Furthermore, the monitoring approach provides the basic means to ensure interoperability between shared memory as well as existing paradigm independent tools. This strongly contributes to the flexibility of the proposed approach and also enables the seamless integration into existing systems.

# Acknowledgments

At this minute I am very excited because I have finished this work. I know at the same time that it would not have been possible without the support of many others. I would like to take this opportunity to express my sincere appreciations of their help.

First of all I would like to thank my advisor, professor Arndt Bode, for providing me with an excellent research environment and giving me the greatest freedom to conduct my research. I am also appreciate of his help in overcoming all the possible financial and administrative issues. In addition, he had always time for me despite his busy schedules.

I would like to thank professor Eike Jessen for taking the time to review this thesis and giving valuable comments and detailed improvements regarding the language.

I would like to express my special gratitude to Dr. Wolfgang Karl, the leader of my research group. He not only directed me to the right way for pursuing my work but also gave me the most considerate help throughout the whole course of this research work. In addition, I especially appreciate his friendship and the time and enthusiasm for helping me to feel at home in Germany.

Besides all my co-workers and the technical and administrative staff at LRR-TUM, I would like to specially mention Dr. Martin Schultz who worked together with me on the same project. His suggestions and good ideas directly contributed to the success of this work. I am especially grateful for his patience and tolerance to my poor German. He also deserves this special thanks for directing me to write English not in a Chinese way.

My stay in Germany has been supported by the German foundation Friedrich-Ebert-Stiftung. It deserves my special thanks.

Finally, I would like to express my deepest gratitude to my parents who have shown great understanding for my impossibility to live around them. I thank my brother and sisters who share my responsibility for taking care of our parents.

*Jie Tao*
*Munich, Germany*
*June 2002*

# Contents

# List of Figures

# List of Tables